

A SMART AND CONNECTED HEALTHCARE DELIVERY PROCESS: FROM PREDICTION TO DECISION SUPPORT

By

Sujee Lee

A dissertation submitted in partial fulfillment of
the requirements for the degree of

Doctor of Philosophy

(Industrial and Systems Engineering)

at the

UNIVERSITY OF WISCONSIN–MADISON

2020

Date of final oral examination: 12/16/2020

The dissertation is approved by the following members of the Final Oral Committee:

Jingshan Li, Professor, Industrial and Systems Engineering

Laura Albert, Professor, Industrial and Systems Engineering

Douglas A. Wiegmann, Associate Professor, Industrial and Systems Engineering

Gabriel Zayas-Cabán, Assistant Professor, Industrial and Systems Engineering

Marlon Mundt, Associate Professor, Family Medicine and Community Health

To My Heavenly Father

Acknowledgments

First and foremost, I would like to express my sincere gratitude to my advisor, Professor Jingshan Li, for his tremendous support and guidance throughout my graduate study at the University of Wisconsin-Madison. Without his supervision, none of the research in this dissertation would be possible. I am deeply indebted to his patience and dedication over the past years and I am so fortunate to have him as my advisor. He will always be my greatest role model in my future academic life.

I wholeheartedly thank my committee members, Professor Laura Albert, Professor Douglas A. Wiegmann, Professor Gabriel Zayas-Cabán, and Professor Marlon Mundt. Their invaluable feedback on my dissertation was essential to substantially improve this manuscript. Thanks to their encouragement and consideration, I was able to successfully finish my degree and pursue the academic career even in the COVID-19 pandemic. Also, I would like to thank Dr. Philip Bain and Dr. Albert Musa from SSM Health Dean Medical Group, and Dr. Christine Baker from SSM Health St. Mary's Hospital, who provided the valuable collaboration opportunities throughout my study.

I also thank my dear friends and colleagues. My academic siblings, Wenjun Zhu, Hyo Kyung Lee, Xiang Zhong, Cong Zhao, Feng Ju, and Xiaolei Xie are always my great friends and colleagues. Congfang Huang, Hoyoung Yoo, Seongkyung Cho, Soovin Yoon, Yerang Park, and Yunji Park are my best friends who supported my Ph.D life and made my days in Madison full of warm and happy memories.

I am particularly grateful for the unconditional love and support from my family. My parent and my brother have always supported and encouraged me to pursue my goals.

More than anything, I thank Sanghoon Shin who has always been supporting me by my side. He always gives me strength whenever facing challenges in my life in Madison from a distance.

This dissertation would not have been possible without the support of the institutions. This manuscript is based upon research supported by the National Science Foundation, the Agency for Healthcare Research and Quality, and Wisconsin Distinguish Graduate Fellowship, and the ISyE department.

Lastly, thank God, my source of inspiration, wisdom, knowledge and understanding. His faithfulness and love have made me excel and be successful in all my academic pursuits. This dissertation is sincerely dedicated to Him.

Contents

Contents	iv
List of Tables	vii
List of Figures	viii
Abstract	x
1 Introduction	1
1.1 Healthcare Systems with Emerging Data Integration	1
1.2 COPD Patient Readmission Risk and Intervention	3
1.3 Opioid Consumption Prediction and Optimal Prescription for TJR Patients	6
1.4 Continuous Delivery of Care: Workflow in Primary Care	8
1.5 Organization of the Document	10
2 Literature Review	12
2.1 COPD Patient Readmission Reduction	12
2.2 Opioid Prescriptions to TJR Patients	16
2.3 Workflow in Primary Care	18
3 Prediction of COPD Patient Readmission Risks	21
3.1 Introduction	21
3.2 Methods	22
3.3 Data Collection and Processing	25
3.4 Learning a Causal Bayesian Network	28
3.5 Model Validation and Results	31
3.6 Discussion	34
3.7 Conclustions	38
4 Prediction of TJR Patients Opioid Consumption	40

4.1	Introduction	40
4.2	System Description	41
4.3	Data Collection	42
4.4	Data Preprocessing	46
4.5	Modeling	48
4.6	Results	55
4.7	Conclusions	61
5	Intervention Planning for COPD Patients Post Discharge	62
5.1	Introduction	62
5.2	COPD Intervention Process	63
5.3	Optimality Analysis	67
5.4	Case Study	76
5.5	Discussions	79
5.6	Conclusions	88
6	Dynamic Intervention Decision for Reducing COPD Readmissions	90
6.1	Introduction	90
6.2	Causal Network Markov Decision Process	91
6.3	Solving CNMDP Problems	101
6.4	Structural Properties	104
6.5	Case Study	107
6.6	Conclusions	111
7	Optimization of Opioid Prescription Post TJR Surgery	113
7.1	Introduction	113
7.2	An Integrated Opioid Prescription Optimization Framework	114
7.3	Stochastic Programming for Opioid Prescription	117
7.4	Results and Solutions	124
7.5	Sensitivity Analysis	131
7.6	Conclusions	134
8	Improving Care Quality in Primary Care Clinics	136
8.1	Introduction	136
8.2	System Description and Problem Formulation	137
8.3	Workflow Modeling	139
8.4	Performance Analysis	144

8.5	Extensions	154
8.6	Case Study	163
8.7	Conclusions	169
9	Conclusions and Extensions	171
9.1	Summary of Contributions	171
9.2	Extensions	174
	Bibliography	196

List of Tables

3.1	Variables and their discretized values in the COPD patient records	26
3.2	Prediction Accuracy of the resulting Bayesian Network and other machine learning methods	33
4.1	Cohort characteristics table	45
4.2	Classification Results with threshold 200 MME	58
4.3	Classification Results with threshold 700 MME	58
4.4	Top 7 important features in the resulting models with threshold 200 MME . . .	59
4.5	Top 7 important features in the resulting models with threshold 700 MME . . .	60
5.1	Case study parameters ($N = 50, \rho = 0.3, C_d = 5000$)	76
7.1	An example of practicing an optimal solution on the collected data	129
8.1	Accuracy of total working time formula	159
8.2	Accuracy of patient waiting time formula	161
8.3	Accuracy of documentation waiting time formula	162
8.4	Performance measures in a morning shift ($N = 8$)	164
8.5	Performance measures in an afternoon shift ($N = 6$)	165
8.6	Performance measures in a full day ($N = 14$)	166
8.7	Sensitivity Analysis with increased number of patient visits	167
8.8	Sensitivity Analysis with varying patient service and documentation times . . .	168
8.9	Sensitivity Analysis with increased patient visits and reduced service and documentation times	168
8.10	Sensitivity Analysis with varying CVs	169
.1	Conditional Probability Table $p_G(y x_p)$	182
.2	Conditional Probability Table $p_G(x_p x_a, x_e)$	182
.3	Conditional Probability Table $p_G(x_a x_w)$	182
.4	Conditional Probability Table $p_G(x_e x_s)$	182

List of Figures

1.1	A smart and connected healthcare delivery process	2
1.2	The structure of the thesis	11
3.1	Arcs that are ruled out from the search space by data-specific domain knowledge	30
3.2	Possible arcs that can be forced to be added in candidates of graph structure . .	31
3.3	The full graph structure of resulting Bayesian network using domain knowledge	32
3.4	The partial graph structure of resulting Bayesian network that is not a causal network	35
3.5	The decreased readmission probabilities after manipulating each variable through an intervention	37
3.6	The increased readmission probabilities after manipulating each variable through an intervention	38
4.1	Flow diagram of inclusion and exclusion of TJR patient data in the study	43
4.2	Diagram of semi-supervised learning process	49
4.3	Probability of pseudo-labeling failure	52
4.4	Posterior distributions of regression coefficients	56
4.5	Histogram of opioid consumption of originally labeled data (left) and predicted values (right)	57
5.1	Intervention flow model	64
5.2	Projection of the problem in a 2D space	68
5.3	Illustration of candidates for optimal combination of interventions	74
5.4	The candidates for optimal combinations of interventions, when rehab is the most cost-effective intervention	75
5.5	The candidates for optimal combinations of interventions, when PCP is the most cost-effective intervention and taking both is next	75
5.6	The candidates for optimal combinations of interventions, when PCP is the most cost-effective intervention and rehab is next	76

5.7	Case study illustration	78
5.8	ICER illustration	81
5.9	The resulting optimal readmission rate with the change of PCP cost	83
5.10	An optimal policy transition by the PCP cost shift	83
5.11	The resulting optimal readmission rate with the change of rehab cost	84
5.12	An optimal policy transition by the rehab cost shift	85
5.13	The resulting optimal readmission rate changes by the change of each interven- tion's readmission rate	86
5.14	The resulting optimal readmission rate changes by the change of each interven- tion's readmission rate	86
6.1	Structure of CNMDP	96
6.2	A full diagram of the learned causal network for COPD readmission (Figure 3.3 in Chapter 3)	108
6.3	CNMDP model for COPD readmission	108
6.4	Resulting optimal policy depending on different efficacy	111
7.1	Patterns of opioid prescribed amount at discharge and usage within 14 days . .	115
7.2	An analytical framework of classification and stochastic programming models .	115
7.3	Results of opioid demand distribution estimation	125
7.4	Optimal solutions of Problem 7.5 for Group L	126
7.5	Aggregating the optimal solutions into final solution of the original problem . .	127
7.6	Comparison of the proposed framework with one-size-fits-all prescriptions . .	128
7.7	Relationship between classification accuracy and resulting density estimation .	133
7.8	Relationship between classification accuracy and resulting optimal objectives .	135
8.1	Performance comparison in all models	152
8.2	Performance measures with non-exponential distributions	156

Abstract

Healthcare delivery is facing a paradigm change to embrace rapid development in information technology, data analytics, artificial intelligence, as well as numerous medical devices and treatments to achieve smart and interconnected care.

In a smart and interconnected healthcare system, integration of data analytics, system modeling, optimal decision-making, and care intervention is necessary and important. Based on the collected data, including the patient's demographic information, disease history, physical exam, and diagnostic test, the smart and patient-specific intervention decision will be formed, and proper care practice will be delivered. Through this, all activities of prevention, diagnosis, treatment, clinical visits, and home care are all connected together.

This dissertation is dedicated to providing analytical frameworks for such a smart and connected healthcare delivery process to address issues related to classification, prediction, intervention, and care service by integrating machine learning, optimization, and system modeling techniques. Specifically, (1) through data collection and preprocessing, predictive models are developed to stratify patients, determine the patient's status, or predict risks by means of machine learning algorithms. (2) Based on patient identification, through modeling the post-discharge care process, intervention plans and policies are evaluated and optimal decisions are proposed. (3) Finally, to implement the intervention and treatment plan, care delivery policies are studied to improve care quality. Through these steps, an integrated and comprehensive framework can be established to connect data analytics, intervention planning, and care services in a closed loop.

In order to show the significance and applicability of such frameworks for the smart and connected healthcare system, this dissertation introduces analytical frameworks applied on readmission risk management for COPD (Chronic Obstructive Pulmonary Disease) patients and opioid prescription optimization for TJR (Total Joint Replacement) patients. In order to provide models of continuous care delivery, a workflow policy development study for primary care physicians is also introduced.

Specifically, for reducing COPD readmissions, two different sub-frameworks integrating machine learning models and operations research methods are developed. The first sub-framework classifies COPD patients into the high or low risk of readmissions, then based on the risk group, an intervention resource allocation is determined through linear programming and graphical analysis. In the second sub-framework, a training procedure for a causal Bayesian network is proposed and the resulting causal network describing relationships between factors and readmission is integrated into a Markov decision process to provide a dynamic intervention planning. For opioid prescription optimization, a novel approach for semi-supervised learning is proposed and the resulting classification model predicts patients' expected opioid consumption levels. Then, a stochastic program is introduced to decide how many opioids should be prescribed to each class of patients to reduce opioid leftovers and thereby curtail the opioid crisis. Finally, as primary care is in charge of continuous care delivery regardless of patients' underlying diseases and conditions, workflow models for primary care physicians are developed by utilizing stochastic process modeling techniques.

In summary, the work developed in this dissertation provides novel frameworks enabling smart and interconnected care to treat patients in need and resolve issues in the U.S. healthcare system.

Chapter 1

Introduction

1.1 Healthcare Systems with Emerging Data Integration

Healthcare delivery is facing a paradigm change to embrace rapid development in information technology, data analytics, artificial intelligence, as well as numerous medical devices and treatments to achieve smart and interconnected care.

In a smart and interconnected healthcare system, integration of data analytics, system modeling, optimal decision-making, and care intervention is necessary and important. Based on the collected data, including the patient's demographic information, disease history, physical exam, and diagnostic test, the smart and patient-specific intervention decision will be formed, and proper care practice will be delivered. Through this, all activities of prevention, diagnosis, treatment, clinical visits, and home care are all connected together.

In this dissertation, we consider such a smart and connected healthcare delivery process to address issues related to classification, prediction, intervention, and care service by integrating machine learning, optimization, and system modeling techniques. Specifically, (1) through data collection and preprocessing, predictive models are developed to stratify patients, determine the patient's status, or predict the risk (risk of readmission or disease exacerbation) by means of machine learning algorithms. (2) Based on patient identification,

through modeling the post-discharge care process, intervention plans and policies are evaluated and optimal decisions are proposed. (3) Finally, to implement the intervention and treatment plans, care delivery policies (in terms of physician workflow models) are studied to improve care quality and the result will be fed back to improve the analysis. Through these steps, an integrated and comprehensive framework can be established to connect data analytics, intervention planning, and care services in a closed loop. An illustration of such a smart and connected healthcare framework is presented in Figure 1.1.

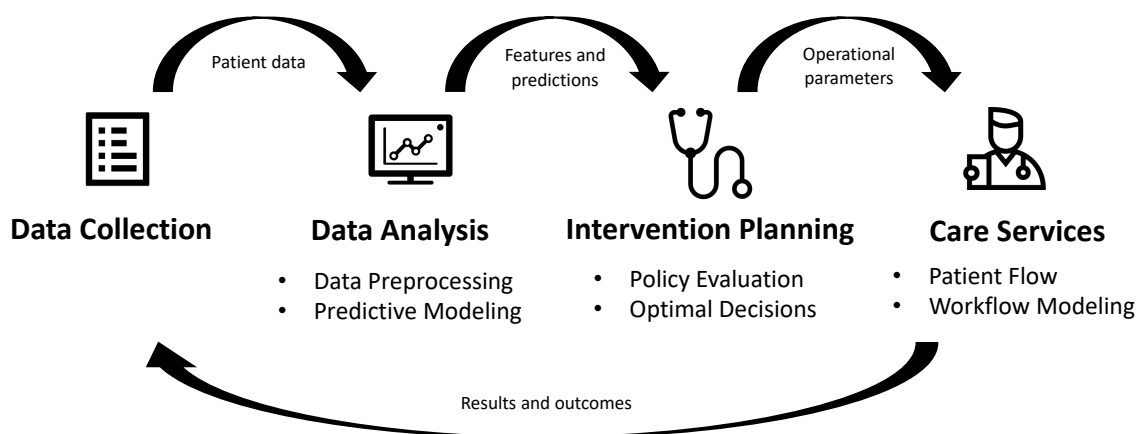


Figure 1.1: A smart and connected healthcare delivery process

Such a process of care is implemented in specific healthcare issues such as COPD (Chronic Obstructive Pulmonary Disease) readmission reduction, opioids prescription optimization for TJR (Total Joint Replacement) patients, and primary care physician's workflow. For reducing COPD readmissions, two different sub-frameworks are developed. The first sub-framework classifies COPD patients into high or low risk of readmissions, then based on the risk group, an intervention resource allocation is determined. In the second sub-framework, a causal model describing relationships between factors and readmission is integrated to a Markov decision process to propose a dynamic intervention planning. For opioid prescription optimization, a semi-supervised classification model predicts patients' expected opioid consumption levels, and a stochastic program decides how many opioids

should be prescribed to each class of patients to reduce opioid leftovers. Finally, as primary care is in charge of continuous care delivery regardless of patients' underlying diseases and conditions, workflow models for primary care physicians are developed by integrating patient data and analysis into operational parameters.

The following three subsections address the research background and motivations for the specific implementations, COPD readmission reduction, opioid prescription optimization, and physician's workflow studies.

1.2 COPD Patient Readmission Risk and Intervention

COPD is a chronic lung disease that makes a patient hard to breathe. It is used as an umbrella term to describe progressive lung diseases. Increased breathlessness, frequent coughing (with and without sputum), wheezing, and tightness in the chest are the major signs and symptoms. Smoking, fumes, chemicals, dust, and genetic factors are the main causes of COPD [1].

As an increasingly common disease that has affected 30 millions of Americans (over half of them have symptoms but do not know it), COPD is becoming one of the leading causes of death in the US. It is also a worldwide issue. In 2015, COPD affected about 174.5 million, i.e., 2.4%, of the global population, resulting in 3.2 million deaths, up from 2.4 million deaths in 1990 [2,3]. Such numbers are projected to increase further due to higher smoking rates in developing countries and aging populations worldwide [4]. More than \$2.1 trillion cost is estimated due to COPD [5], with \$1.9 trillion direct cost such as medical care. In the US, such costs are estimated at \$50 billion, most of which are due to exacerbation [6].

Although COPD is treatable, readmissions after discharge can occur, which not only cause pain and discomforts to the patients, but also incur substantial cost. It is estimated that the average cost of a readmission from a diagnosis including COPD is \$10,900. Moreover, starting from October 2014, hospitals in the US are being penalized by the Centers for

Medicare and Medicaid Services (CMS) on excessive 30-day unplanned readmissions for COPD [7]. The three-year hospital discharge data for COPD readmissions will be compared to the national average, and a percentage of a hospital's aggregate Medicare payments can be penalized. Such an initiative has intrigued hospitals and doctors to reduce costs, eliminate waste, and standardize evidence-based protocols to achieve patient-centered and value-driven care. Thus, reducing COPD readmissions has been a major effort in hospitals across the country.

Although COPD readmission has received substantial research attentions, most of them are qualitative-oriented or pilot study-based (see reviews [8–10]). The limited quantitative analyses mainly focus on identifying the risk factors correlated to COPD readmission and predicting the risks using statistical methods (see, for instance, [11–14]). To the best of our knowledge, such risks and the associated factors have seldom been used in quantitative models to supervise the intervention or follow-up plans.

Moreover, most prediction models identify correlations between predictors and readmission risk, rather than causal relationships. To design effective intervention plans, the impact of manipulations of risk factors should be known. Thus, causal models are needed. Also, as mathematical models can provide a fresh look of the intervention process, developing an optimization framework, particularly, based on causal prediction models, is of significant importance to support decision making by providing guidelines for patient-centered interventions.

In addition to the lack of causal models and subsequent models for intervention planning, numerous challenges and obstacles exist in introducing and implementing effective follow-up interventions. A low medical compliance level is one of them, which is due to prevalence of psycho-social issues, such as low socioeconomic status, predisposition to unhealthy behaviors and life style, heavy comorbidity burden, complexity in home medical equipment [15,16]. Such a low level of compliance results in increased risks of readmission. Thus, developing and implementing effective interventions with improved

medical compliance become a compelling need for the hospitals.

Furthermore, in the clinical environment, medical intervention decisions are often made sequentially during the course of disease progression and treatment process. Likewise, post-discharge interventions should be tailored to each individual patient and they should be updated by considering the patient's health condition and characteristics during the monitoring period in order to achieve better outcomes. However, only limited studies are devoted to designing personalized intervention policies and the dynamic nature of post-discharge care is not addressed. Therefore, there is a need to develop personalized intervention plans and update intervention decisions dynamically.

To overcome aforementioned limitations of existing studies, we develop a causal Bayesian network which can not only predict future readmission risks, but also identify factors that account for readmissions and should be managed. Moreover, to reduce readmission risks while improving the compliance level of COPD patients to interventions, we propose an incentive framework for hospitals to reimburse out-of-pocket and transportation costs for patients visiting primary care physicians (PCPs) and respiratory rehabilitation centers in their follow-up plans based on their expected readmission risks. It is expected that the increased compliance level can lead to reduced readmissions, whose savings in readmission cost and penalties will exceed the incentive cost. Also, as causal Bayesian network models are capable of finding the causal relationships and Markov decision process (MDP) models can be used to find optimal solutions to seek stochastic and dynamic decisions, a causal Bayesian network-based Markov decision process (CNMDP) model is proposed to provide personalized dynamic decision support for intervention planning.

1.3 Opioid Consumption Prediction and Optimal Prescription for TJR Patients

The United States is in the midst of the opioid epidemic, with increasing prevalence of opioid misuse, overdose, addiction, and death [17, 18]. Opioid overdose becomes the leading cause of unintentional deaths and declining lifespan expectancy [19]. From 1999 to 2017, almost 400,000 people in the United States have died from opioid overdose [20–22]. As the opioid epidemic continues, a recent focus has been on clinical prescribing practices or narcotic pain medications [23].

Opioids, including prescription opioids, heroin, and synthetic opioids are a class of drugs used to reduce pain. Among them, prescription opioids, such as oxycodone, hydrocodone, codeine, and morphine, are typically prescribed by physicians to treat chronic or acute pain, or manage pain-related serious health conditions, for instance, cancer [24,25]. However, serious risks and side effects have been associated with the prescription opioids. A substantial number of patients using prescription opioids for chronic pain may develop opioid use disorder [26]. It is reported that 55% of opioid deaths are related to prescription opioids [20–22]. Up to 2% of the population is estimated to be currently using prescription opioids non-medically, and over 55% of prescription opioid misusers obtain opioids from a friend or relative who receives the prescription from a physician [27]. Since one of the major sources of opioid overdose is hospital over-prescription [28], to fight against the opioid crisis, reducing over-prescription at hospitals is of significant importance.

Post-surgery opioid prescription is one of the leading contributors to over-prescription [29]. In the United States, orthopedic surgeons are responsible for 7.7% of all opioid prescriptions and are the third highest opioid prescribers [30]. In a surgery setting, an opioid-naive patient can be easily exposed to the drug [31]. It has been reported that among the surgery patients who have not taken opioids at all on their last day of hospitalization, 45% have been prescribed opioids at discharge [32]. Moreover, several studies have found

that most patients do not consume all of the opioid pills they are prescribed [33,34]. These facts imply that surgery patients can be prescribed more opioids than their needs, or be prescribed even when no opioid is necessary, during the post-surgery recovery period. An even more concerning fact is that up to 92% of patients do not properly dispose of their excess opioid pills, and the excess opioids pills could be passed on to surrounding communities of patients and place them at great risk of abuse. [35]. Thus, unnecessary opioid pills in the prescription should be reduced, which can help curtail the problems related to opioid crisis.

Along this line, the Centers for Disease Control and Prevention (CDC) provides a guideline for opioid prescription [36]. However, such a guideline is only useful for a general and broad purpose, since it introduces a one-size-fits-all recommendation regardless of the patient type or surgical procedure. In fact, both the surgical operations and patients are different from each other in terms of pain mechanism, intensity and tolerance. Therefore, an opioid prescribing strategy should be variable subject to different demands, i.e., adaptable with respect to the patient and surgery type and tailored to each patient's need [29].

Among different types of surgeries, total joint replacement (TJR), including both total knee replacement (TKR) and total hip replacement (THR), is a common procedure performed for more than a million cases each year [37]. Such a number is expected to grow substantially due to the increasing demands of mobility, advanced diagnosis and treatment of arthritis, and longer life expectation [38]. Management of pain following total joint arthroplasty is an important aspect of this procedure. Adequate pain control allows faster rehabilitation and reduces the risk of postoperative complications, and opioids are very effective analgesics and constitute the foundation of the management of moderate-to-severe acute postoperative pain. Thus, the TJR patients are often prescribed with oral opioid medication after surgery to help manage and control pain after hospital discharge.

There have been numerous studies attempted to identify the risk factors of persistent post-surgical pain, pain severity, and prolonged opioid use after TJR [39–48]. Such studies

mainly focus on the long term pattern such as persistent opioid use after 6 or 12 months. However, in order to avoid over-prescription and reduce the number of superfluous pills, it is necessary to understand the actual opioid use patterns in a short time period. To our best knowledge, most existing literature does not provide any research that actually tracks the exact amount of used opioids after discharge, and there are only few studies monitoring TJR patient's opioid usage. Further, there is no study using a mathematical model to determine the optimal prescription of opioids for TJR patients.

Therefore, we focus on understanding and modeling the short-term actual opioid usage for TJR patients after discharge. The goals of the study are, first, collecting the information of TJR patients' short-term opioid usage; second, developing a classification model using semi-supervised learning to identify patients who may need less or more opioids after being discharged from TJR surgeries; and third, helping decision making on the appropriate amount of opioid prescription tailored to each patient. To suggest the optimal prescriptions meeting the patients' needs while leaving minimal unused opioids, a framework that integrates a stochastic programming with the classification model is developed. The research efforts, in turn, will help reduce the overall opioid over-prescription and overdose.

1.4 Continuous Delivery of Care: Workflow in Primary Care

As the backbone of nation's healthcare system, primary care is facing significant challenges in recent years [49–51]. Physicians have become increasingly busy to deal with overwhelming amounts of tasks, regulation pressures, and electronic health record (EHR) usages [52].

A general internist in a typical primary care clinic may spend his/her time in the following major categories: First, time in exam room meeting with patients; second, time in his/her office between patient visits; and last, time after hours catching up remaining work

of the day. Thus, the majority of a physician's time is on face-to-face encounters with patients and on computer and desk work [53], with a small portion of time on interacting with other medical staff. Particularly, working on computers/tablets to document information related to patient visits, communicate with nurses, medical assistants or other staff electronically, carry out administrative work, and reply to in-basket electronic messages from patients, etc., has increased the physician's workload [53,54].

Substantial amount of research has been devoted to analyzing physician workload. Most of them are empirical, qualitative, or simulation based studies, which are useful, but may not be easy to find system properties that can provide quantitative guidance or predictions for continuous improvement, and they may require more efforts and computation times to evaluate different parameter settings. As mathematical formulation and derivation can quickly analyze physician workflow, evaluate workload, and provide predictable suggestions and insights for continuous improvement, developing such models becomes important. However, only limited analytical work exists, mainly focusing on long term planning without much consideration of EHR documentation activities. Since most primary care clinics operate with fixed working hours, studying only the long term behavior may not be suitable. In addition, since EHR and computer related tasks have occupied a substantial amount of physician's time in many primary care clinics [53,54], including such tasks in a physician's workflow model is critically needed.

To overcome the aforementioned issues, in this study, we use terminating Markov chain models to study a physician's workflow management policies, in which the absorbing or terminating state indicates that all tasks of the day have been finished by the physician. Specifically, as meeting with patients and doing computer work occupy most of the physician's working time, we simplify and abstract the workflow into direct clinical face-to-face encounters and EHR documentation tasks. The former includes physical examination, diagnosis, and procedures, etc., while the latter consists of EHR writing, review, assessing, and all computer works associated with patient care. By integrating patient data and

analysis into operational and service distribution parameters, three workflow management policies are considered according to either stopping or continuing ongoing EHR documentation during a new patient's arrival (referred to as preemptive priority (PEP) and non-preemptive priority (NPP) policies, respectively), or documenting EHR tasks in batches (BDC policy). Closed formulas to evaluate the efficacy of such workflows in terms of physician working time, patient and documentation waiting times, are derived and the corresponding performance measures are compared. Both Markovian and non-Markovian scenarios are studied. The main contribution of the study is in developing analytical models to study physician workflow management policies and evaluate the efficacy of primary care clinics. Finally, a case study at a primary care clinic is carried out to illustrate the applicability of the models.

1.5 Organization of the Document

In the course of the proposed health management process, each analytical component is a subject of research. Each chapter is devoted to each analytical component within the care delivery where a set of components are integrated together to form an analytical framework to solve a specific issue in the healthcare system. The analytical component and the corresponding stage of process and application for each chapter are depicted in Figure 1.2.

The rest of this document is organized as follows. Chapter 2 reviews the related literature for predictive models, intervention planning, and policy design and evaluation for the healthcare topics addressed in the dissertation, and the modeling techniques and methods utilized. Chapter 3 introduces a causal Bayesian network not only to predict COPD patients readmission risk, but also to find variables that need to be intervened to reduce the readmission risk. In Chapter 4, a survey study and development of a semi-supervised classification model for predicting the expected postoperative opioid consumption levels of TJR patients

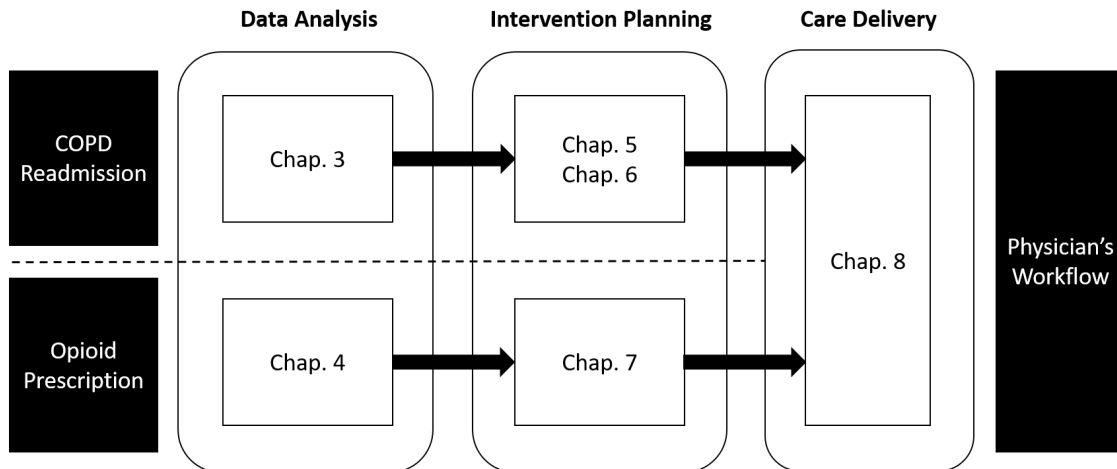


Figure 1.2: The structure of the thesis

are presented. By incorporating the classification model developed in Chapter 3, Chapters 5 and 6 provide frameworks for post-discharge intervention decision to reduce hospital readmissions of COPD patients. Chapter 5 proposes a framework where an optimization model determines intervention resource allocation based on the classification results. In Chapter 6, the causal Bayesian network is integrated into a Markov decision process to enable dynamic intervention decisions by updating the patient's information at every decision epoch. Chapter 7 presents an opioid prescription optimization framework that minimizes unused opioids as well as the number of refills using a stochastic programming. In the framework, the opioid consumption prediction model in Chapter 4 is integrated to diversify the prescription plans tailored to the consumption levels. Chapter 8 proposes terminating Markov chain models to design and manage primary care physician's workflow policies. Finally, the summary and possible extensions are presented in Chapter 9. All proofs and derivations can be found in the Appendices.

Chapter 2

Literature Review

The goal of the research in this dissertation is to develop rigorous engineering approaches to model and solve problems in health care delivery systems to improve patient care. Thus, the literature review focuses on the following aspects: Section [2.1](#) highlights the research efforts and the relevant analytical methods on COPD readmission reduction. Section [2.2](#) reviews the studies regarding opioids prescription plans and the related factors for TJR patients. The existing research on primary care physician's workflow is summarized in Section [2.3](#).

2.1 COPD Patient Readmission Reduction

This section reviews the existing studies related to COPD readmission predictions as well as the methods utilized to develop the COPD readmission reduction frameworks in the dissertation. Particularly, the healthcare studies utilizing a causal Bayesian network and a Markov decision process are reviewed as such methods are parts of the frameworks introduced in Chapter [3](#) and Chapter [6](#), respectively.

2.1.1 COPD Readmission and Interventions

Substantial efforts have been devoted to reducing hospital readmission (e.g., reviews [55–57]) for patients having heart failure [58], surgery [59–61], pneumonia [62], kidney injury [63], diabetes [64] and sepsis [65] as well as COPD [8].

Since COPD is one of the leading causes of death, a significant amount of research has been focused on finding the factors that impact COPD patients' readmission rates. A systematic review of available studies is provided in [8] to identify risk factors for readmission of COPD exacerbation patients. Prior readmission, dyspnea, and oral corticosteroids are identified as significant factors for risk prediction. Low health status, poor quality of life due to bad health conditions, and lack of routine physical activity, are also identified as related factors that increase readmission risks. In other studies, similar or additional factors have been identified, such as heavy smoking, comorbidity, depression, severity of disease, social support and socioeconomic characteristics, etc. (see, for instance, [66–69]).

By considering these risk factors, numerous intervention methods have been proposed. A review of such activities is presented in [9], which includes patient education about the use of respiratory inhalers, instruction on steps for alleviating symptoms, phone or page hotline, general health counseling, coordination with PCP, home visits, and follow-up phone calls. It also includes smoking cessation counseling, social services referral, and assessment of comorbidities, etc. However, the effects of these efforts may not be significant. As argued in [70], one of the reasons is lack of emphasis on patients with the greatest risk of readmissions. Another reason is due to compliance. In [15], it is indicated that the level of compliance for COPD patients is lower than the compliance rates recorded for other diseases, despite the severity of the medical condition. It is also stated in [16] that adherence in COPD patients is generally poor and becomes a big concern. It can be associated with a significant health and economic burden in those patients.

To segment patients and identify patients in higher risk of readmission, models to predict the risk of readmissions have been developed. For example, a review of five

statistical models for COPD readmission prediction is presented in [13], including logistic and Bayesian logistic regression, logistic regression tree, generalized estimating equations, and generalized linear mixed model. A simple model to estimate readmission risks and identify high readmission risk patients is introduced in [11] using logistic regression. However, based on 7843 citations and 30 studies, it is found in [12] that most updated readmission risk prediction models are only beneficial in limited scenarios. In [14], it also claims that there is no well-acknowledged clinically useful score to predict COPD exacerbation in a short term.

Although various prediction algorithms have been applied to this area, and the predictions are informative, they may not be useful unless the predicted risk is used to design the most appropriate follow-up plans for the patients' post-discharge care. Also, the question of identifying the critical variables whose manipulation can be used to reduce readmission probabilities can not be answered by the prediction algorithms.

2.1.2 Causal Bayesian Network

General machine learning prediction models may not be suitable to address the impact of interventions since such models investigate only correlations or associations, rather than causal relationships [71, 72]. For example, suppose that from a clinical history dataset, a hospital recognizes that patients with lower BMI show lower readmission rate. In this case, BMI can be used to predict future readmissions since the correlation can be derived from the dataset for the population. However, we cannot tell that lowering BMI results in readmission reduction. If there exists another factor, such as comorbidity, that affects both BMI and readmission risk, then changing the BMI may have no impact on future readmission. This hypothetical case exemplifies the needs of finding the causal relationships between variables in order to identify risk factors affecting readmissions. For this purpose, a causal Bayesian network can be introduced as it is capable of such causal inferences as well as classic predictions.

A causal Bayesian network (CBN) is a specific type of Bayesian networks where the directions of the arcs should indicate causal relationships between the ending variables. Although the automatic learning of the CBN structure from data is an active challenge, there have been a lot of attempts of algorithm development to learn the structure solely based on data or by adopting domain knowledge. Since a CBN enables causal inferences and suggests how changing one variable is likely to influence the state of another variable, CBNs are widely used in healthcare research where causal discovery can contribute to diagnosis, prognosis, and treatment of diseases. For instance, A study in [73] develops a CBN in order to predict the progression rate of Amyotrophic Lateral Sclerosis (ALS) disease. Paper [74] uses a CBN to support decision making in type 1 diabetes. The authors in [75] developed a mathematical model of chronic myeloid leukemia based on CBNs in order to study possible disease progression mechanisms. Survival prediction and treatment selection in lung cancer care are considered in [76] and CBNs are applied to aid lung cancer experts by providing personalised survival estimates and treatment selection recommendations. Also, a decision support system utilizing CBNs for assisting a management of Warfarin therapy is developed in [77]. However, a CBN has not been developed for targeting COPD readmissions.

Therefore, a causal Bayesian network model analyzing COPD readmission should be developed to identify causal relationships between factors and future readmissions as well as to predict the risk of readmissions. Using this model, the effects on readmission risks of manipulating variables that can be controlled by clinical actions or interventions can be analyzed, which can provide guidelines for intervention planning.

2.1.3 Markov Decision Processes

A Markov decision process (MDP) is a mathematical framework for modeling decision making in problems where outcomes are partly random and partly under the control of a decision maker. MDP models have been used extensively for decision problems in stochastic

domain, including healthcare applications [78,79]. For example, an MDP is utilized in [80] to model kidney transplantation process that allows patients to accept or reject an offered kidney based on its quality. Meanwhile, an MDP model presented in [81] to determine the optimal timing of liver transplantation from a living-donor. In [82], an MDP model is introduced to decide the optimal intervention and timing in mild hereditary spherocytosis to maximize a patient's quality-adjusted life years. In addition, a study in [83] develops an MDP model that determines the optimal time for initiating HIV treatment, and another study in [84] uses an MDP model to evaluate the optimal start time of statin treatment for cardiovascular risk reduction. Furthermore, in [85], a finite-horizon discrete-time MDP model is formulated for the optimal breast cancer biopsy decision.

In spite of the existing work, to the best of our knowledge, there is no analytical model available to dynamically predict patient-specific readmission risks and adjust optimal intervention policies for COPD patients to reduce readmissions. By integrating risk prediction and optimal intervention decision, this dissertation addresses both objectives of developing a causal network and providing a dynamic intervention planning for reducing COPD readmission.

2.2 Opioid Prescriptions to TJR Patients

Numerous studies have been carried out to identify the factors affecting opioid usage after TJR surgery. For example, the risk factors for a higher number of refills or usage in 12 months post-TKR are identified in [39] [40], and [41]. The discovered factors include preoperative opioid use along with younger age, anxiety or depression, female gender, and other intrinsic variables. Paper [42] has shown that more severe pain, worse functioning, symptoms of depression, and a higher preoperative opioid dose are the risk factors. Paper [43] suggests that the type of surgery (TKR/THR), longer hospitalization stay, discharge to a rehabilitation facility, preoperative opioid use, a higher comorbidity score, back pain,

migraine and smoking at baseline are strong predictors for persistent opioid use. In addition, prior opioid use has been recognized in [44] and [45] as a strong risk factor for chronic use of opioids. Paper [46] summarizes that experiencing depression or anxiety and having pain in other places can be significantly related to persistent pain after surgery or postoperative opioid use.

In addition to applying the statistical methods, machine learning techniques have also been used in opioid use studies. Logistic regression models are the prevailing ones to find associations between long-term patterns of opioid usage and the risk factors, see, for instance, [39,41–46]. In addition, papers [86] and [87] intend to predict persistent opioid use or substance dependence. Limited work has been started to use machine learning methods besides logistic regression to classify the patients who experience opioid dependency or abuse (e.g., [87–89]). An opioid classification study usually requires collections of patients' demographic, physiological, and clinical factors related to opioid use. In the above studies, hundreds of thousands of patients' electronic health records (EHR) are used to train the models. However, the aforementioned studies have a common limitation that the EHR data, medical claims data or opioid prescription history are only used for analysis of long-term patterns of opioid usage, such as opioid use after 6 months. In addition, none of those studies suggest further use of their findings in practice, especially in opioid prescriptions.

There are several studies addressing the optimal dose or prescription of medications. For example, [90] and [91] investigate the optimal heparin dosing policy. The former uses multivariate logistic regression to estimate the probability of a given dose of heparin to a patient being sub-therapeutic or supra-therapeutic. The latter develops a reinforcement learning model for optimal dosing policy of heparin to maximize the overall time a patient stays in acceptable therapeutic conditions. In addition, paper [92] introduces LASSO Bandit to solve a restricted version of contextual Markov decision process, and applies the model to identify the correct level of warfarin dose. In the above cases (such as intravenous

unfractionated heparin or warfarin), lab tests of blood clotting are performed for patients who receive heparin after four to six hours, and the test results can work as feedback to provide information of appropriateness of previous heparin dose, which makes it easier to apply data-driven approaches.

Nevertheless, the current literature lacks data-driven studies for optimization of opioid dose or prescription. Among the limited studies, paper [93] collects opioid-use levels through a phone survey to find out potential excess pills. Using this information, a new number of pills to adjust the original prescription has been suggested. However, it recommends a fixed-size prescription approach regardless of different patient characteristics. Another study in [94] emphasizes the urgent need for optimizing opioid prescribing, while the analyses are conceptual rather than mathematical or practical. Moreover, refills are not considered in such studies.

In summary, despite the efforts in existing literature, there is no available study using and predicting the actual level of opioid usage in a short term for patients after TJR surgeries. The goals of the study are first to collect and analyze such information and to develop a classification model through a proposed semi-supervised learning. Also, there is a critical need of a data-driven study for optimizing opioid prescription. By integrating a classification model for anticipating patients' expected opioid use in a short period of time and an optimization model for determining optimal prescription, the ideal amount of opioids to minimize leftover and refills based on the classification results can be obtained. Developing such an integrated framework is the final goal of the study.

2.3 Workflow in Primary Care

The primary care physicians are facing tremendous challenges [49–52]. The physician shortage, overcrowding patient population, and insurmountable EHR tasks make the primary care physicians practicing with increasing workload, which causes the system

unsustainable. For example, the study in [53] indicates that only 27% of a physician's time is on clinical face-time, and 49% on EHR or desk work. Similarly, as shown in [54], primary care physicians spend nearly six hours, almost two thirds of their workday, in EHR related work, such as patient portal and administrative tasks. A time-motion study of primary care physicians' work in 982 visits at 10 clinics is introduced in [95] to investigate the mixed impact of EHRs. It also shows that the physicians spend more time in the EHRs than face-to-face meetings with patients. More clinical and observational studies on EHR workload have been carried out and can be found in, for instance, [96–98].

Many qualitative studies have been conducted to evaluate the primary care physician workflow or tasks. For instance, a comprehensive list of primary care physician tasks performed during a face-to-face patient visit is developed in [99], which consists of 12 major tasks, 189 subtasks, and 191 total tasks. Such a list can help clinics to study workflows and plan changes in EHR implementation. Further evaluation of the variations in primary care physician workflow during patient visits is introduced in [100]. The physician workflow is determined from patient visits to 10 physicians in 10 primary care clinics. It is shown that there exist significant workflow variations during patient visits, thus the healthcare design should support a wide variety of task sequences in primary care. To investigate the impact of medical scribes on workflow, a crossover study of 18 primary care physicians conducted in [101] indicates that using medical scribes is associated with substantial reduction of EHR documentation time and improvement of productivity and job satisfaction.

Besides the qualitative and observational studies, mathematical models to quantify the impact of workflow management with EHR tasks can provide a new perspective. In recent years, operation research methods have been widely used in studying healthcare systems (e.g., monographs [102–106] and reviews [107–109]). Discrete event simulation, Markov chain and queueing models have been the prevailing tools to study patient access and workflow in hospitals and clinics. For example, in ambulatory care, the length of stay of elder patients in long term care facilities is evaluated using Markov models in [110].

The patient flow in outpatient healthcare setting is modeled by a semi-Markov process in [111]. Paper [112] analyzes the workflow and staffing level and identifies the bottleneck stage in a CT test center using a Markov chain model. Similar study is carried out in [113] to investigate the location of follow-up scheduling service for a gastroenterology clinic. In [114], a resource sharing iteration approach is presented to study workflow in a mammography test center. To understand work allocation between physician and medical assistant (MA), the staffing ratio is investigated and a workload balance law is discovered in [115]. Furthermore, joint visits (by physician and MA simultaneously) in primary care clinics are studied in [116] to evaluate various work distribution schemes.

In addition, conceptual single-server queueing models have been introduced to study primary care appointment systems in [117]. Queueing models have also been used to evaluate patient cycle time in outpatient facilities in [118], to analyze nurse practitioner utilization in primary care in [119] and [120], and to study patient no-show in a game theoretic framework in [121]. Recently, electronic visits (message exchange between physicians and patients) and the impact in primary care clinics, such as flow time, panel size, and physician work, are studied in [122] and [123]. Team communication and collaboration are evaluated using queueing models in [124]. However, the above studies do not directly evaluate the physician workflow arrangement issues, such as scheduling face-to-face encounters and EHR tasks and the associated workloads.

The existing literature using Markov chain or queueing models typically focuses on steady state behavior. However, in primary care clinics, finite arrivals and working hours are more realistic. Thus, developing terminating Markov chain models of physician workflow is needed. Moreover, as Markov property may not hold in practice, non-Markovian scenarios should be addressed.

Chapter 3

Prediction of COPD Patient Readmission Risks

3.1 Introduction

In this chapter, we develop a causal Bayesian network that predicts future readmission of COPD patients so that intensive care interventions can be applied to patients identified with high risks for readmission. The causal Bayesian network also can answer questions of finding the variables whose manipulation can be used to reduce readmission probabilities, which cannot be addressed by a typical machine learning prediction model [71,72]. Using the model, the effects of manipulating variables, which can be controlled by clinical actions or interventions, on COPD readmission risks can be analyzed, which can provide guidelines for intervention planning.

To develop such a causal Bayesian network, the data are first collected from electronic health record (EHR) and then preprocessed into appropriate forms for Bayesian network training. Next, the preprocessed data are used to learn Bayesian network structure through a score-based search method. Data-specific domain knowledge is adopted during the training step in order to obtain a causal Bayesian network, which enables us to make causal

inferences about the effect of manipulating variables on future readmissions. Finally, based on the inferences, analysis of post-discharge intervention and treatment plan to reduce the readmission probability can be carried out.

3.2 Methods

The most common and dependable approach to finding a causal relationship between two variables is conducting experiments. In clinical research, randomized controlled trials are used to understand the effect of a treatment on an outcome. However, such experiments can only deal with an univariate relation between a single variable and an outcome, and it is often too expensive, hard or not even ethical to carry out [71, 125]. Thus, in this chapter, instead of finding the exact causal relationships through experiments, we develop a model that is based on the observed data as in general machine learning approaches, and simultaneously employs a limited number of causal relationships known from common knowledge or domain knowledge. For the purpose, we implement a causal Bayesian network because its network structure is capable of representing causal relationships between variables, and the domain knowledge can be easily adopted.

It is worthy to note that the model developed in this chapter is different from general Bayesian networks, in which the directions of the arcs can be meaningless. In the causal Bayesian network, the arcs must have directions that represent causal relationships. Below we introduce both the Bayesian network and the causal Bayesian network, and describe the process of causal inference from the network.

3.2.1 Bayesian Network

A Bayesian network is represented by a probabilistic directed acyclic graph (DAG), in which the vertices are the random variables, and the arcs represent the probabilistic dependencies between the two variables corresponding to the arc's two ending vertices [126, 127]. In

the graph, a vertex v has a conditional probability table $P(v \mid \text{Parent}(v))$, describing the relationship between the variable and its parents. An important characteristic of a Bayesian network is the local directed Markov condition, which implies that, conditioned on its parents, v is independent of variables who are neither descendants nor parents of v in the network. Such a feature allows the joint probability distribution to be described as

$$\begin{aligned} P(v_1, v_2, \dots, v_n) &= \prod_{i=1}^n P(v_i \mid v_{i-1}, \dots, v_1) \\ &= \prod_{i=1}^n P(v_i \mid \text{Parent}(v_i)) \end{aligned} \quad (3.1)$$

In Bayesian networks, the directions of arcs do not need to be meaningful, and some of them may not even follow the chronological order. This is due to the fact that the two networks $A \rightarrow B$ and $A \leftarrow B$ are equivalent in probabilistic models, producing the same marginal distributions and generating the same probabilistic inference to the same queries such as $P(A \mid B)$ [128].

3.2.2 Causal Bayesian Network and Inference

A causal Bayesian network has the same form of Bayesian Networks, except that, in a causal model, a variable A 's parents are A 's direct causes relatively to other variables in the causal network. In other words, an arc $A \rightarrow B$ exists in a causal network if and only if A is a direct cause of B in the specified population, described as the given dataset.

If the causal Bayesian network is known, we can make a causal inference and evaluate the effect of manipulation of a variable v^* . If we change the distribution of v^* to $P'(v^* \mid w)$, where $w \neq v^*$, then we just need to replace the term $P(v^* \mid \text{Parent}(v^*))$ in (6.3) with the newly specified probability $P'(v^* \mid w)$. Then, after the specified manipulation, we obtain

the changed joint distribution $P'(v_1, v_2, \dots, v_n)$ as,

$$P'(v_1, v_2, \dots, v_n) = P'(v^* | w) \prod_{v \in V \setminus \{v^*\}} P(v | \text{Parent}(v)), \quad (3.2)$$

where $V = \{v_1, v_2, \dots, v_n\}$. That is, if the causal network structure is known, and there is an estimated unmanipulated probability $P(v_1, v_2, \dots, v_n)$ from the observed dataset, then the prediction of the effect of an manipulation can be made even if there does not exist observed data regarding the manipulation.

In this study, such an inference rule can be simplified because our objective is to change a variable to a fixed value rather than a specified distribution and evaluate the impact of this change. Using the identity function instead of $P'(v^* | w)$ in (6.4), and using Bayes rule to construct conditional probability, we can obtain the following rule [129]: Let $P(y | x^*, x_i, x_j)$ be the probability that the value of target variable Y is y given that v^* is set to x^* and $v_i = x_i$, $v_j = x_j$ are observed. For any variables v^* , v_i , v_j , and Y , the following rule exists:

$$P(y | x^*, x_i, x_j) = P(y | x^*, x_j),$$

if we delete all arcs pointing to v^* from the original graph and Y is independent of v_i given v^* and v_j in this network.

In this study, such a rule enables us to evaluate the effects of manipulation of controllable variables that are critical to COPD readmissions. Through learning a causal Bayesian network using the dataset in St. Mary's Hospital and applying the inference rule, we intend to evaluate the impact of manipulation of critical controllable variables on readmission, so that appropriate interventions to avoid readmissions can be suggested.

3.3 Data Collection and Processing

3.3.1 Data Collection

St. Mary's Hospital is a community hospital located in Madison, Wisconsin. It is part of SSM Health, which is a Catholic, not-for-profit health system serving the communities across the Midwest. St. Mary's Hospital offers a full range of inpatient and outpatient treatment and diagnostic services in primary care and nearly all specialties in South Central Wisconsin since 1912.

In this study, through extensive literature review and interviews with physicians, nurses, pulmonologists, respiratory therapists, and other staff in St. Mary's Hospital, we identify the risk factors that impact COPD readmissions. These factors are categorized into four groups, describing the social status, physical status, behavior status, and medical information of the patients. A list of these variables is shown in the second column of Table 3.1, where BMI is the body mass index, SpO₂ measures peripheral capillary oxygen saturation, and FEV₁/FVC is defined as the forced vital capacity in 1 second/forced expired volume.

Through accessing to St. Mary's Hospital's data center, we collect a total of 877 COPD patients' records in the last four years. Among them, 556 patients have been admitted mainly due to COPD, and 130 patients have been readmitted within 30 days. Such a dataset is used for model development after preprocessing.

3.3.2 Preprocessing

Data preprocessing is an important step in machine learning studies because the reliability of a predictive model is critically dependent on the quality of its training data. To obtain a reliable model, the following preprocessing steps are carried out.

First, besides removing outliers, by using background knowledge of each variable, we modify or remove the implausible values in the data set, which are considered as mistakenly written values in the data collecting process. For example, the values with different units

Table 3.1: Variables and their discretized values in the COPD patient records

ID	Variable Name	States	Values		
1	Age	3	< 65,	[65, 80),	≥ 80
2	Gender	2	Male,	Female	
3	Marital status	3	Single,	Married,	Widow/ Divorced
4	Employment	3	Not Employed /Retired,	Part /Full time,	Disabled
5	Working environment	2	Polluted /Heavy work,	Not	
6	Ethnicity	3	White,	African American,	Asian
7	Number of children	3	0,	1,	≥ 2
8	Have supporters at home	2	Yes,	No	
9	Height	3	< 1.7,	[1.7, 1.8),	≥ 1.8
10	Weight	3	< 70,	[70, 90),	≥ 90
11	BMI	3	< 25,	[25, 40),	≥ 40
12	Have severe allergies	2	Yes,	No	
13	Have cardiac comorbidity	2	Yes,	No	
14	COPD type	3	Bronchitis,	Emphysema,	Both
15	Tobacco use history	3	< 25,	[25, 50),	≥ 50
16	Use tobacco currently	2	Yes,	No	
17	Use alcohol currently	2	Yes,	No	
18	Use drug currently	2	Yes,	No	
19	Use home equipments	2	Yes,	No	
20	Length of stay	3	< 5,	[5, 8),	≥ 8
21	Number of past admission	3	0,	1,	≥ 2
22	Have been readmitted	2	Yes,	No	
23	Attitude to care	2	Resolved,	Active	
24	Attitude to education	3	Resolved,	Active,	Done
25	Show anxiety during education	2	Yes,	No	
26	Blood pressure	2	< 130,	≥ 130	
27	Pulse	3	<60,	[60, 100),	≥ 100
28	Temperature	3	< 97.7,	[97.7, 99.5],	> 99.5
29	Respiratory rate	2	[10, 26),	≥ 26	
30	Oxygen saturation (SpO ₂)	3	<0.96,	[0.96, 0.99],	> 0.99
31	FEV1/FVC	3	< 0.5,	[0.5, 0.7),	≥ 0.7
32	Readmit in 30 days	2	Yes,	No	

in height are adjusted from centimeter to meter. The SpO₂ values that are larger than 1 are removed. Similarly, the values that are out of the defined range or state are removed. In addition, there exist some variables having too many null values, such as a patient's year of education, whose values are missing in more than 50% of the records. These variables are also removed.

After these preprocessing steps, some variables still may have missing values. In this case, most of the machine learning algorithms require or prefer to impute them in order to deal with the complete data. However, direct imputation may lead to large deviation and produce a bias in the model since the distribution is unknown. In this study, instead of imputation, for each variable, the Bayesian learning algorithm is used to handle the data itself with missing values by estimating the distributions based on the records without missing values.

Next, continuous variables, such as age, height, weight and vital signs, need to be discretized since most Bayesian network learning and inference algorithms are unable to efficiently handle the data that have both discrete and continuous variables with arbitrary distributions. Furthermore, due to the limited number of data records and relatively large number of variables, when there are many states for each variable, the number of observations with each state value will be less than the necessary size to estimate a conditional probability. Thus, we restrict the number of states, i.e., the number of intervals or value ranges, so that each continuous variable is discretized into 2 or 3 states.

In addition, some categorical variables have more than 3 values. For example, "Employment (4)" has five different states. In this case, we merge the states with relatively similar status. Finally, all variables have either 2 or 3 states, which can prevent the model from overfitting.

Note that the discretization task is carried out by following either of the two rules: (1) following a standard guideline [130]; (2) binning the values so that the patient records are well distributed. Vital signs such as body temperature, respiratory rate and BMI follow

the first rule, since they have certain thresholds to classify a person. For example, the body temperature within the range of $97.7 - 99.5^{\circ}\text{F}$ (i.e., $36.5 - 37.5^{\circ}\text{C}$) is considered as a normal temperature, thus it is divided into three categories as shown in Table 3.1. In another example, when the FEV1/FVC ratio is less than 0.7 with a post-bronchodilator, it confirms the presence of persistent airflow limitation [130]. For other variables, such as age or the number of children that do not have specific guidelines to identify a person’s risk to COPD readmission, we set the thresholds to divide the range so that the patient records are well distributed, i.e., each category has a similar number of records. The complete list of variables and their states is presented in Table 3.1.

3.4 Learning a Causal Bayesian Network

3.4.1 Experimental Settings

For general Bayesian networks, there exist various well-known learning methods, such as [131–135]. However, there are no prevailing learning methods for causal Bayesian networks and they usually utilize the learning methods of general Bayesian networks with strong assumptions [136]. In this chapter, in order to obtain a causal Bayesian network, we employ the automatic learning methods from Bayesian networks, which are implemented in *bnlearn R package* [137]. To make it applicable for a causal network, we manually restrict the search space of structures through domain knowledge adoption. In this subsection, the experimental setup for Bayesian network learning methods is described. How to apply this method to generate a causal network is explained in the next subsection.

Typically, the automatic structure learning algorithms for Bayesian networks can be summarized in two categories: (1) constraint-based algorithms that use conditional independence tests, and (2) score-based search algorithms that search for an optimal DAG model that maximizes a goodness-of-fit score in the search space. The score-based search algorithms make use of decomposable scores that allow the total score for a DAG to be

calculated as the sum (or product) of the individual node scores in the network, and evaluate the whole network structure on the set of variables $V = \{v_1, v_2, \dots, v_n\}$. Then, we can ascertain how well a network with a particular structure $P(V) = \prod_i P(v_i | \text{Parent}(v_i))$ fits with the data [138]. To find the best structure, search heuristics based on local addition/removal/reversal of edges that increase the score are the prevailing ones [138–140].

In this study, we use score-based search algorithms. Specifically, hill-climbing search algorithm and tabu-search algorithm [135,141,142] are used as the search heuristic methods. For a score function, we use the negative Bayesian Information Criterion (BIC) score, which is a criterion for model selection among a finite set of models [137]. Since it is based on the likelihood function, which is widely used as a goodness-of-fit score, the model with the lowest BIC (the negative BIC with the largest absolute value) score is preferred. Formally, BIC is defined as [143]:

$$\hat{L} = P(D | \hat{\theta}, G), \quad (3.3)$$

$$BIC = \log(n)k - 2 \log(\hat{L}), \quad (3.4)$$

where \hat{L} is the maximized value of the likelihood function of the network structure G , $\hat{\theta}$ represents the parameter values that maximize the likelihood function, D represents the observed data, n is the number of observations, and k is the number of parameters.

3.4.2 Domain Knowledge Adoption

In order to obtain a causal Bayesian network, the search space of DAGs is manually restricted by adopting domain knowledge related to the data. As explained before, a standard Bayesian network learning algorithm cannot generate a causal network. For example, $A \rightarrow B \leftarrow C \leftarrow D$ and $A \rightarrow B \leftarrow C \rightarrow D$ are Markov equivalent (i.e. the two graph have the same conditional independencies) resulting in the same BIC scores, but they lead to different causal inferences. Thus, the Bayesian network algorithm that yields the DAG with

the highest score does not consider whether the DAG is a causal one or not, and the exact causal DAG cannot be automatically generated.

To solve this issue and obtain a Bayesian network that permits reliable causal inference, we restrict the search space by removing causally impossible arcs, and force the graph to have causally plausible arcs only. Common knowledge or domain knowledge is used in this process. For instance, no other variables can be a cause of “Age (1)” or “Gender (2)”, where the numbers in parentheses are the variable IDs in Table 3.1. Thus, any arcs from a variable to “Age (1)” or “Gender (2)” should be removed. In addition, in the dataset, there is an intrinsic temporal order. For example, the value of variable “Readmit in 30 days (32)” is determined last, while other variables are determined ahead of the discharge time or measured at the time of discharge, which is earlier than “Readmit in 30 days (32)”. Then, the arcs from “Readmit in 30 days (32)” to any of these variables can be removed. In Figure 3.1, these ruled out arcs are illustrated.

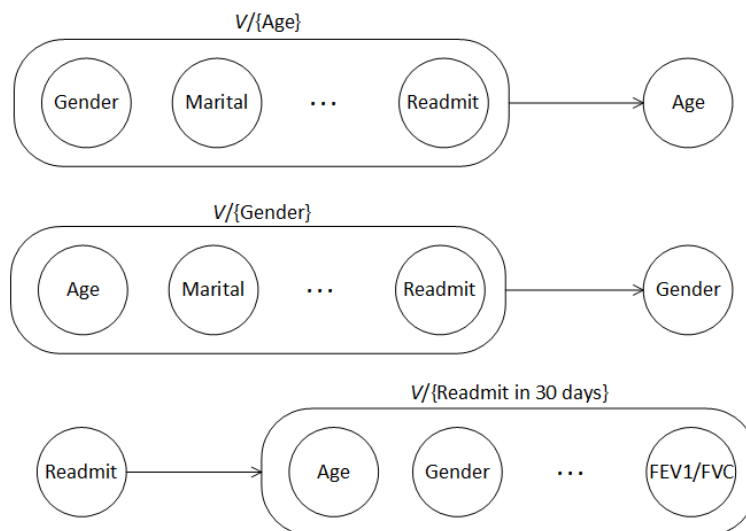


Figure 3.1: Arcs that are ruled out from the search space by data-specific domain knowledge

Remark 3.1. *In addition to removing arcs, new arcs can be added as well. For example, since we focus on finding the relationships between variables and 30-day readmission, using clinical domain knowledge, the arcs from some variables to “Readmit in 30 days (32)” can be manually inserted.*

Additional reasonable arcs can be added as well. For instance, as shown in Figure 3.2, gender may influence weight and height. However, such insertions should be carefully reviewed by healthcare providers to avoid incorrect causality. The resulting models need to be convincing to ensure reliable causal inference.

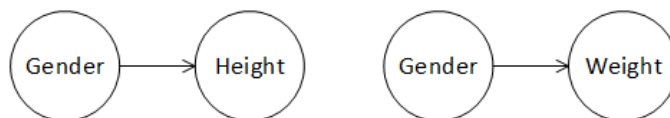


Figure 3.2: Possible arcs that can be forced to be added in candidates of graph structure

3.4.3 Variable Selection

In learning the structure of the network, all variables are used so that any variables that possibly have relationships with readmissions are not ruled out ahead of learning. However, after obtaining the resulting structure, some variables that are not connected to the vertex of “Readmit in 30 days (32)” can be ruled out, and the manipulation of them do not affect the distribution of “Readmit in 30 days (32)”.

3.5 Model Validation and Results

In contrast to prediction models, the validity of a causal network cannot be truly ensured by the dataset, because it involves predictions for the manipulated population that does not exist [71]. However, it is possible to moderately check its causal relationships either by employing background knowledge or by evaluating probabilistic relationships [144].

First, we investigate the resulting structure by checking if the model has an arc that represents an implausible causal relationship between two variables to validate whether the resulting model is sufficient for causal inference. Following the proposed experimental setting, the resulting Bayesian network is obtained and shown in Figure 3.3, where the ancestors of “Readmit in 30 days (32)” are shadowed. It can be observed that most of

the arcs in the graph in Figure 3.3 are causally plausible. For example, “Weight (10)” has “Gender (2)” and “Height (9)” as its parents. While the arc from “Gender (2)” to “Weight (10)” is intentionally added in the learning process, the arc from “Height (9)” to “Weight (10)” is automatically discovered, which is plausible in terms of causality. Other arcs, such as the arc from “Weight (10)” to “BMI (11)”, also causally possible. This fact demonstrates that, to some extent, the resulting causal network is valid.

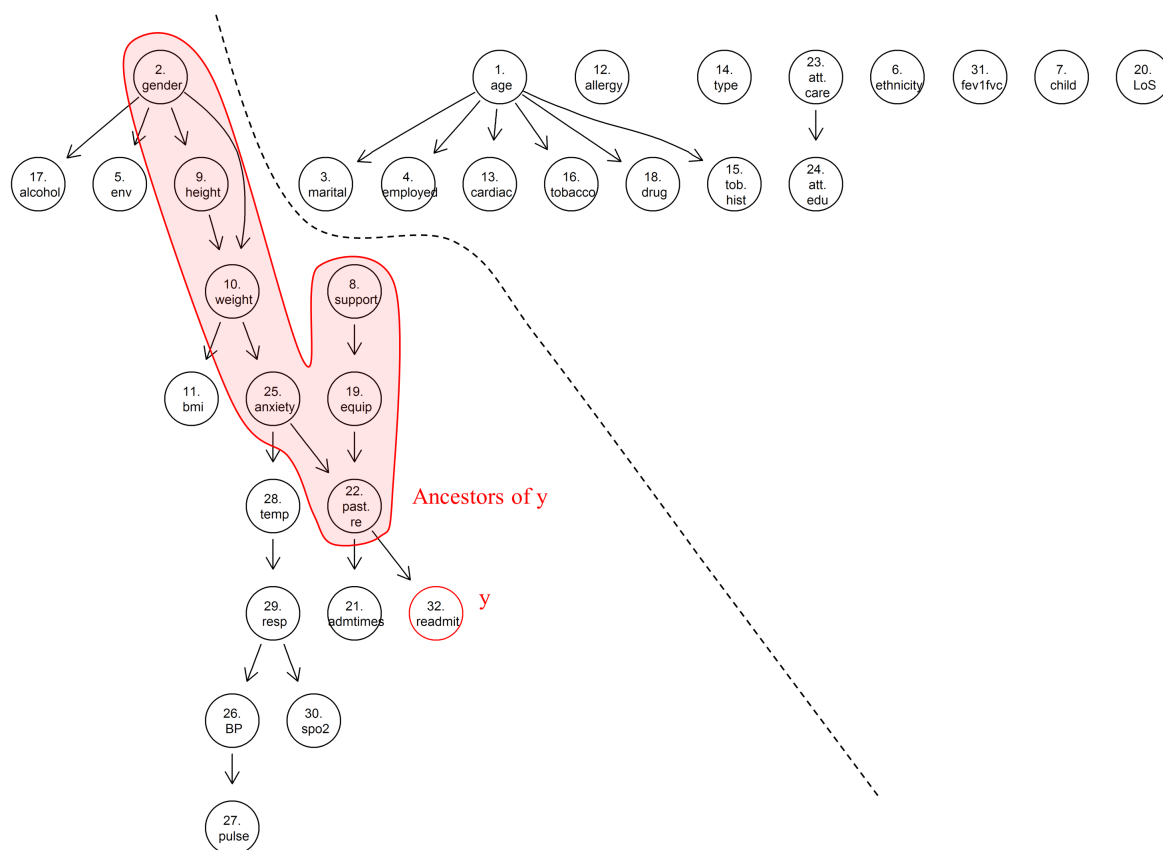


Figure 3.3: The full graph structure of resulting Bayesian network using domain knowledge

Secondly, in order to determine if the resulting model correctly represents the relationships between “Readmit in 30 days (32)” and other variables, we evaluate the predictive accuracy in variable “Readmit in 30 days (32)”, and compare it with other predictive modeling methods. If the resulting Bayesian network achieves a comparable level of prediction accuracy, it implies that the model reflects the relationships between the variables. Note

again that the main purpose is not achieving an accurate predictive inference, but enabling causal inference.

The accuracy of the resulting model for predicting future readmission within 30 days from discharge is shown in Table 3.2 along with the prediction accuracies of other classification models. For the purpose of comparison, the results of Logistic Regression, K-Nearest Neighbors, Naive Bayes, Neural Network, Support Vector Machine and Random Forest are listed. In addition, in order to illustrate the advantage of the proposed model, we add one more result, which is obtained by the Bayesian network automatically generated from the structure learning algorithm without domain knowledge adoption. Note that the models' accuracies in Table 3.2 are evaluated by 10-fold cross validation and they are the best values we can achieve. For each model, different data preprocessing and different model settings are used to obtain the best accuracy.

Table 3.2: Prediction Accuracy of the resulting Bayesian Network and other machine learning methods

Method	Prediction Accuracy
Logistic Regression	0.83
K-Nearest Neighbors	0.77
Naive Bayes	0.79
Neural Network	0.84
Support Vector Machine	0.82
Random Forest	0.88
Bayesian Network (w/o domain knowledge adption)	0.77
Bayesian Network (w domain knowledge adoption)	0.85

As shown in Table 3.2, the prediction accuracies of all classifiers are in the range of 0.77 to 0.88. Random Forest classifier achieves the highest prediction accuracy. However, it does not provide causal inference. The causal Bayesian network with domain knowledge exhibits the second highest prediction accuracy, 0.85, which is high enough to validate that it reflects the variables' probabilistic relationships correctly.

Comparing with the results of a Bayesian network without adopting domain knowledge,

which is depicted in Figure 3.4, there exist remarkable differences. The network in Figure 3.4 is clearly not causal because there are too many arcs representing the reversed temporal order. Although its search space for graph structure is larger and it has a better score, its prediction accuracy is less than that of the causal model. Indeed, its negative BIC is -13426.7 while the causal Bayesian network model has a score of -13464.15.

In summary, the causal Bayesian network with domain knowledge, produced by the above experiment setup, can achieve a high accuracy, and more importantly, it provides causal inferences.

Remark 3.2. *It is of interest to find from Figure 3.3 that in the graphical structure, there is a arc from “Anxiety (25)” to “Have been readmitted (22)”. Since the variable “Anxiety (25)” is collected when a patient shows anxiety when he/she receives education in the hospital before discharge, which typically occurs after “Have been readmitted (22)”, the arc from “Anxiety (25)” to “Have been readmitted (22)” is impossible to be a causal relationship. A possible reason for such an arc is that a multiple-readmitted patient may have shown anxiety during the readmissions. Except for this arc, all other arcs exhibit reasonable causal relationships.*

3.6 Discussion

Using the validated causal network, the impacts of manipulation of variables on the probability of readmission are investigated. Let y be the variable “Readmit in 30 days (32)” and $v_i, i = 1, \dots, 31$, represent the variable with ID i in Table 3.1. Since $P(y | v_1, \dots, v_{31}) = P(y | \text{Parent}(y))$ in a Bayesian network, unless the values of variables in $\text{Parent}(y)$ change, the probability of readmission based on patient status, $P(y | v_1, \dots, v_{31})$, will remain the same. Thus, the probability of readmission changes only if either of the following two cases is satisfied: (1) the parents of y are manipulated, or (2) the ancestors of y are manipulated and thereby the parents of y change.

In the resulting model (see Figure 3.3), the variables that affect the readmission proba-

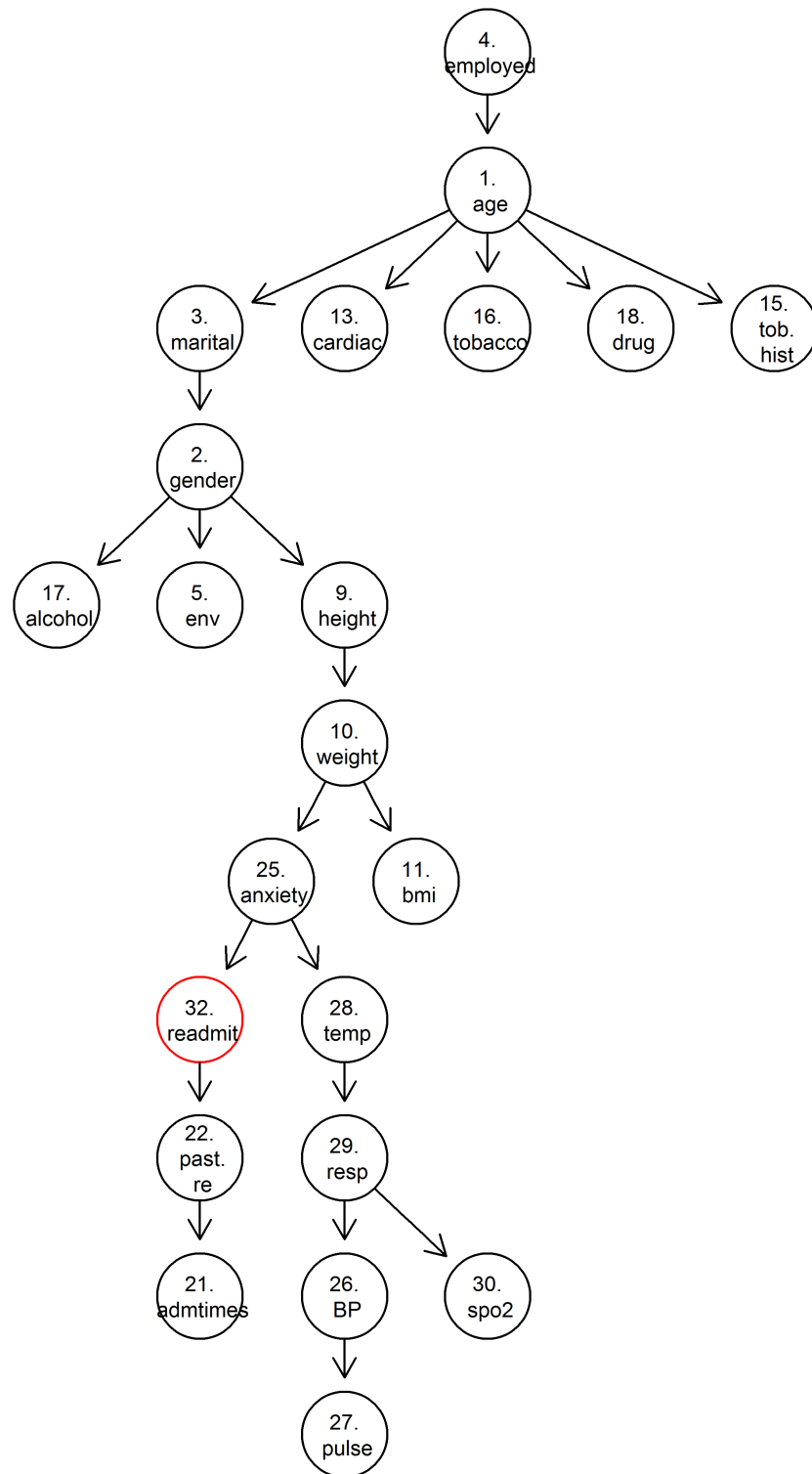


Figure 3.4: The partial graph structure of resulting Bayesian network that is not a causal network

bility are shadowed: “Gender” (v_2), “Height” (v_9), “Weight” (v_{10}), “Show anxiety” (v_{25}), “Have supporters at home” (v_8), “Use home equipment” (v_{19}), and “Past readmission” (v_{22}). However, among them, what can be manipulated are “Weight (10)”, “Show anxiety (25)”, “Have supporters at home (8)”, and “Use home equipment (19)”. If a patient takes an intervention that manipulates one of them, then the values of the descendant vertices alter, so that the effect of interventions can be predicted.

Suppose that a patient has a record x where value of v_i is x_i for all i , and that if the patient takes an intervention that manipulates v_s to make its value change to x'_s , then the descendants of the manipulated variable can be affected. Let v_k be the closest ancestor of y among the variables that are neither an ancestor nor a descendant of v_s . Also, introduce v_j as the parent of y and v_l be the variable that is an ancestor of v_j and a descendant of v_s . Then the updated probability of readmission can be evaluated as follows:

$$\begin{aligned} P(y \mid x_1, \dots, x_{i-1}, x'_s, x_{i+1}, \dots, x_{31}) &= P(y \mid x'_s, x_k) \\ &= \sum_{v_j, v_l} P(y, v_j, v_l \mid x'_s, x_k) \\ &= \sum_{v_j, v_l} P(y \mid v_j)P(v_j \mid v_l, x_k)P(v_l \mid x'_s). \end{aligned}$$

The first equality comes from the rule of do calculus, and the last two equalities are from the structure of the network. The reduction in readmission probability is calculated as:

$$P_{reduced} = \max\{P(y \mid x) - P(y \mid x'_s, x_k), 0\}.$$

If the reduction in readmission probability $P_{reduced}$ is larger than 0, then an intervention can be designed to adjust this variable v_s to reach the state x'_s , and we expect that the patient’s readmission risk within 30 days from discharge can be reduced by $P_{reduced}$.

Using this analysis, we can discover the most effective interventions for a patient with high risk of readmission. Suppose that a patient has ever been readmitted ($x_{22} = 1$),

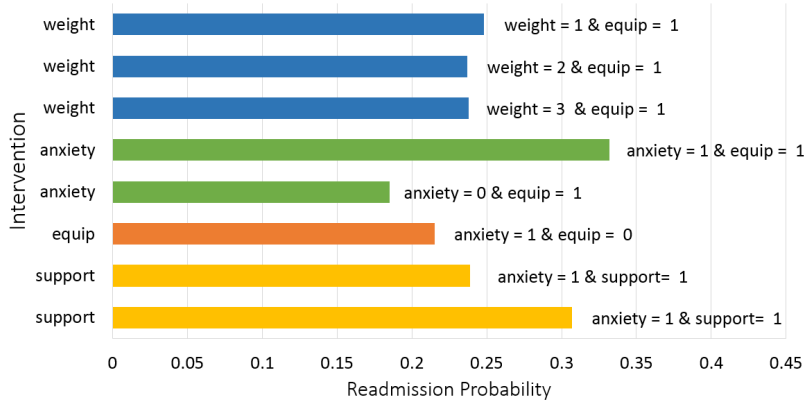


Figure 3.5: The decreased readmission probabilities after manipulating each variable through an intervention

showed anxiety during education in the hospital before the discharge ($x_{25} = 1$), and has a home equipment which helps to breath ($x_{19} = 1$). By using the value of v_{22} , the patient is predicted to have a high probability of readmission: $P(y | x_{25}) > 0.5$. We want to find which variable we need to change in order to most effectively reduce the readmission probability. Using the calculation described above, we predict the updated readmission probabilities after manipulating the controllable variables v_8 , v_{10} , v_{19} and v_{25} . The resulting probabilities are shown in Figure 3.5 and the values of the changed variable and the given status are written next to the bars. From this result, it can be seen that making the patient avoid anxiety ($x'_{25} = 0$) is identified as the most effective way to reduce readmission probability. Therefore, an intervention plan can be designed to let the patient feel less anxiety or do not feel anxiety.

In addition to evaluate intervention plan for high-risk patients to reduce their readmission probabilities, such a framework can also be used to suggest low-risk patients to avoid certain actions that could increase their risks. Then,

$$P_{increased} = \max\{P(y | x'_s, x_k) - P(y | x), 0\}.$$

Analogously, suppose that a patient has never been readmitted ($x_{22} = 1$), and neither

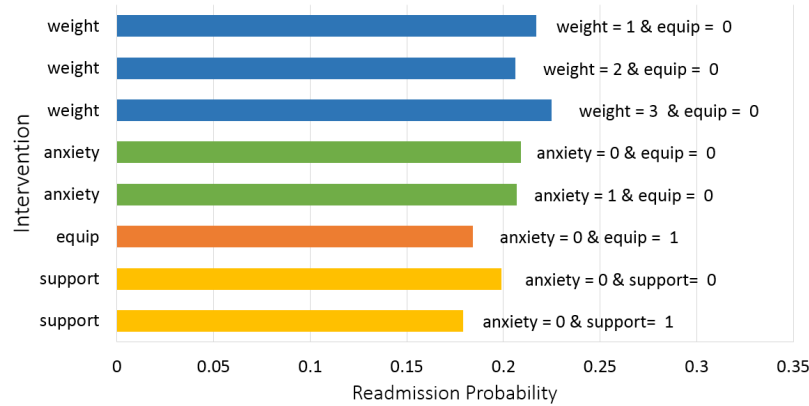


Figure 3.6: The increased readmission probabilities after manipulating each variable through an intervention

feels anxiety during the education ($x_{25} = 0$), nor has any home equipment for COPD palliation ($x_{19} = 0$). Because the past readmission has a large impact on the readmission probability in this model, the predicted future readmission probability for this patient is much less than 0.5, i.e., $P(y | x_{25}) < 0.5$. In order to prevent the patient from getting a higher risk, we calculate new readmission probabilities when the patient's health status is changed. The result is shown in Figure 3.6. As one can see, when the patient's weight state is 3 ($x_{10} = 3$) and $x_{19} = 0$, the readmission probability achieves the highest value. Therefore, the patient should make his/her effort to control the weight.

3.7 Conclustions

In this chapter, a causal Bayesian network that is capable of making causal inference as well as predictive inference is developed. To uncover the most feasible network structure from expert knowledge and data, the search space of network structure is restricted by utilizing domain or expert knowledge of relationship between variables. From the experiment setting, we learn a causal Bayesian network and validate it by observing its structure and comparing its prediction accuracy with other classification algorithms. Using this model, the risk of COPD readmission within 30 days can be predicted, which can support making

decisions for delivery of care for severe COPD patients. Moreover, through the model, the impacts of manipulating critical variables can be analyzed, which can provide guidance to design appropriate post-operative interventions to avoid readmission.

Chapter 4

Prediction of TJR Patients Opioid Consumption

4.1 Introduction

In order to precisely prescribe opioids to TJR patients to meet their need while avoiding the excess opioids, their expected opioids use amounts or use levels should be appropriately predicted by a prediction model. This chapter presents the study of developing such a classification model and conducting a survey study to monitor opioid use of TJR patients and obtain data for the classification modeling.

Since there is no continuous monitoring of actual opioid usage for TJR patients after discharge, a survey study was conducted to collect such information at SSM Health St. Mary's Hospital in Madison, Wisconsin. The result of the survey shows that the amount of opioids used in the same time frame has a large variation across different TJR patients, which signifies the need, and also the possibility, of classification. However, numerous missing responses in the survey are observed, so that the resulting data based on the survey responses can lead to a weak classification model. To effectively handle the missing target values, a tailored semi-supervised learning model is proposed in this study.

4.2 System Description

After receiving TJR surgeries, opioids are often used to manage the post-surgical pain in parenteral or oral forms using injection solution, patches or pills for three or four days during a patient's stay in the hospital. When the patient is discharged from the hospital, additional opioids, typically in form of pills, are usually prescribed, which is referred to as an initial prescription. After consuming all or most of the opioids provided at the initial prescription, a patient can request refills. Also, in many cases, the refill requests are submitted at post-surgical check-ups, which are on two, four, or six weeks after discharge. If the patient keeps requesting more opioids after 90 days from discharge, he/she is directed to pain management and primary care clinics.

The typical types of opioids prescribed in many hospitals may include oxycodone, tramadol, hydrocodone-acetaminophen, acetaminophen-codeine, etc. For example, in SSM Health Saint Mary's Hospital (SMH), oxycodone is prescribed in most cases, and tramadol is mainly used for patients who are older than 80. Other types of opioids are prescribed based on their clinical responses (efficacy, effectiveness, toxicity, and safety). The potency of opioids varies between types, and even within the same type depending on the strength per unit (mg/pill). As the goal of this study is not to consider a specific type of opioids, but to target on the optimal amount (quantity) of opioids, we standardize the opioid content into the unit of Milligram Morphine Equivalent (MME), which is a value assigned to opioids to represent their relative potencies.

The Orthopedic Department in SMH has initiated numerous efforts to reduce opioid over-prescription by calling for physicians to accurately prescribe what the patients would take in the future. To achieve a personalized prescription, identification of patients demanding less or more opioids during a given time period is required. Specifically, opioid demands in a short-term should be considered because a physician usually does not prescribe a TJR patient a very large amount of opioids for long time use, even for the patient with chronic pain. Thus, in order to develop a classification model for patient identification,

obtaining accurate estimates of opioid quantities the patients have taken for a short time period is necessary. However, different from other studies addressing opioid dependency or long term use, such values cannot be inferred only by opioid prescription history information or EHR database because patients may consume less opioids than the prescribed amount. Moreover, the opioid use at later months post-operation cannot identify heavy users of opioids since some patients may only consume opioids intensively within a shorter time period rather than taking them continuously. Thus, the actual amount of opioid usage in a short term after TJR surgery needs to be tracked.

To obtain such information, a patient survey was conducted at SMH to inquire participants about opioid leftover at the time of survey. In addition to the survey, the participants' opioid prescription histories were collected to understand their preoperative opioid usage and gather the postoperative opioid prescription information, and the risk factors were identified and extracted from their electronic health records (EHR). Note that the study was conducted under the approval and management of human subject committee of the hospital and followed the HIPPA protocols and ethical guidelines. All survey and EHR data were de-identified for the analysis.

4.3 Data Collection

4.3.1 Survey Study

A total of 446 patients who received TKR or THR surgery between December 2015 and June 2018 at SMH participated in the survey. The survey was intended to be completed during the patients' post-operative appointments, which occurred typically after two, four, or six weeks post-discharge. In the survey, the patients were asked to disclose the types of opioids being prescribed and the remaining quantities. Note that it is typically easier to ask the leftover quantity rather than counting and keeping track of the number of opioids having been taken for the period between discharge and the survey. Then, based on the

patient's opioid prescription history, the amount of opioid usage after TJR operations can be obtained.

In the survey study, however, some patients participated in the survey too late, thus failing to disclose the short-term usage. Some patients provided invalid answers, such as "unsure", numbers without opioid types, or remaining opioids amount larger than prescription. The flow diagram shown in Figure 4.1 indicates that those late or invalid answers are excluded. In addition, the survey was conducted at two, four, and six weeks post discharge, and the period of opioid consumption for the analysis was set up as two weeks since more survey responses were received at the second week. Figure 4.1 also shows that there are 285 valid values satisfying the criteria while other 161 patients miss the information of opioid consumption within two weeks.

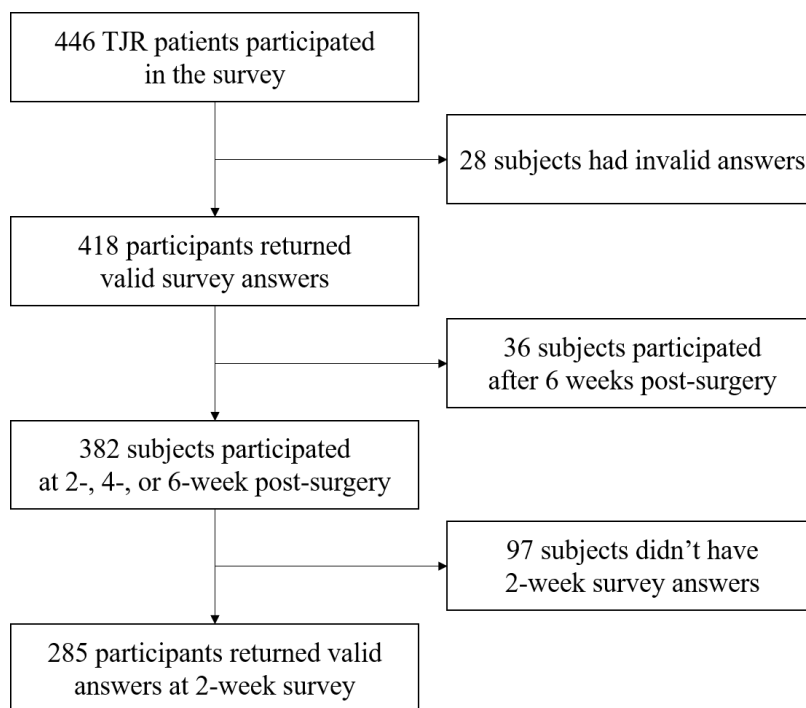


Figure 4.1: Flow diagram of inclusion and exclusion of TJR patient data in the study

4.3.2 Opioid Prescription History

The information of opioid prescription post-hospitalization, such as types and quantities of opioids prescribed at discharge and in subsequent refills within three months, were obtained from opioid prescription history. In addition, the type and quantity of opioids taken by each patient during hospitalization, referred to as facility administered medication (FAM), were also collected. FAM can be utilized as a factor in the classification model because it is known prior to discharge. However, other consumption information related to prescription may not be available for the classification model as the activities occur after discharge. The survey and prescription history are used to calculate the amount of opioid consumption. Using the obtained values, the patients are grouped into multiple classes by setting up the thresholds and spitting the consumption range, and these classes will compose a target variable in the classification model.

4.3.3 Collecting Risk Factors from EHR

In order to construct the feature variables in the model, the characteristics of the participants were extracted from their EHR records. The candidates of features include the risk factors identified in literature review and the variables recommended by the collaborating providers. The list of collected features and their statistics are summarized in Table 4.1. In the table, the continuous variables are presented with their mean values and standard deviations, while the categorical variables are provided by the numbers of records in the categories with percentages.

Table 4.1: Cohort characteristics table

Features		Total (n=446)
Age (years)		65.1 (10.0)
Gender	Male	195 (43%)
	Female	251 (56%)
Procedure	THA	190 (42%)
	TKA	255 (57%)
Height (cm)		170.5 (10.1)
Weight (kg)		94.4 (21.7)
BMI	Extreme	67 (15%)
	Obese	207 (46%)
	Overweight	118 (26%)
	Normal	50 (11%)
	Underweight	1 (0%)
Two-week opioid consumption(MME)		464.7 (403.1)
Facility administered medicine before discharge (MME/day)		44.6 (27.4)
Opium naive	Yes	76 (17%)
	No	369 (82%)
Use opium with in 60 days prior to the surgery	Yes	134 (30%)
	No	311 (69%)
Chronic opium use	Yes	55 (12%)
	No	390 (87%)
Alcohol consumption	Yes	103 (23%)
	No	343 (76%)
Smoking	Yes	46 (10%)
	Former	153 (34%)
	No	247 (55%)
History of depression	Yes	128 (28%)
	No	313 (70%)
History of anxiety	Yes	87 (19%)
	No	359 (80%)
Have hyperlipidemia or migraine	Yes	230 (51%)
	No	214 (47%)
Have other chronic pain	Yes	256 (57%)
	No	190 (42%)
Pulse at surgery		78.6 (13.3)
Blood loss at surgery	Any blood	184 (41%)
	Minimal	259 (58%)
Tourniquet used	Yes	255 (57%)
	No	191 (42%)
Previous surgery experience	Yes	209 (46%)
	No	236 (52%)
Surgical revision	Yes	45 (10%)
	No	399 (89%)

4.4 Data Preprocessing

4.4.1 Constructing a Target Variable

To deal with opioid data, the quantity of opioids should be normalized because different types and potency of prescribed drugs (e.g. oxycodone 5mg, hydrocodone-acetaminophen 5-325mg, and tramadol 50mg) are prescribed to the TJR patients. Morphine Milligram Equivalent (MME) is used to convert different opioids into an equivalent dose of morphine. The daily MME values are calculated as follows:

$$\text{MME/Day} = (\text{Strength per unit}) \cdot (\text{MME conversion factor}) \\ \cdot (\text{Number of units/Days supply}).$$

The target outcomes in the classification are the indicators whether a patient will take less opioids than regular, or more opioids than others. Such separations of patient groups are based on the quantity of consumed opioids, a variable representing the amounts of opioids taken by patients in two weeks needs to be constructed. The value can be obtained by subtracting the remaining opioids at the second week survey from the total amount of opioids prescribed initially and refilled within two weeks.

4.4.2 Preprocessing Feature Values

Data preprocessing or feature engineering is an important step in machine learning studies since the performance and reliability of a resulting model is significantly dependent on the quality of training data. To obtain a better and more robust model, the following steps for preprocessing features have been carried out.

First, features that have unique values or are non-observable up to the time of discharge are dropped. For example, "Patient No." has a unique value in each record, thus it is dropped after integrating opioids demand and EHR data. "Number of opioid refills" and

other features related to post discharge opioid description are dropped because the classification and the subsequent decision of initial prescription should be made at discharge. However, FAM is not excluded even though it is a prescription-related feature because FAM is administered during hospitalization and observed ahead of discharge.

Next, missing values in features are filled out. Since the data used in the study were manually collected and data correctness was pursued, there are few missing values existed in the feature data. Only 3, 5, 2, 3, 1, and 2 missing occur in “weight/BMI”, “History of depression”, “hyperlipidemia or migraine”, “Blood loss”, “Previous surgery experience”, and “Surgical revision”, respectively. Since the numbers are small, null values in the continuous variables are filled with the mean values, and nulls in the rest binary variables are imputed with the most frequent values.

The next step is discretization. The dataset contains continuous variables, such as “Weight” and “BMI”. Transforming the continuous variables into finite numbers of intervals can be beneficial since it often leads to faster and more accurate learning due to data reduction and simplification. However, not all continuous features should be discretized, which depends on the results of testing models using either continuous or discretized forms. When the continuous form works better, still the original form will be used. For example, FAM is not discretized in the final models, while numerical BMI is categorized into standard levels of obesity as shown in Table 4.1. For some other continuous variables, we set the cut-offs and discretize them to maximize the Chi-square test significance so that the variables can help distinguish patients in different classes and result in a better classification results [145]. Specifically, for each variable, we split the intervals using one of the split candidates (e.g., quantiles) and conduct the Chi-square test across patient classes and two split intervals. Then, the best split candidate that maximizes the Chi-square criterion is chosen to make the continuous variable to be binary. Some categories can also be combined into one categorical feature if they describe similar items. For example, “Yes” and “Former” can be combined into one category in “Smoking”.

4.5 Modeling

4.5.1 Missing Data in the Target Variable

Although majority of the patients responded in the survey at their two-week post-surgery timeline, only 285 out of 446 patients have explicitly answered all questions in the survey, which implies that there are only 285 labeled data, and the remaining 161 data are unlabeled.

Since a supervised classification model needs to use labeled data, those missing answers will lead to only 64% of patient data being usable to develop a classification model if we discard the patient records with missing values. Such an abandonment can make the data samples limited to obtain reliable models. In addition, it may not be an appropriate way to deal with MAR data in survey response. To overcome this deficiency, a semi-supervised model employing a tailored pseudo-labeling method that imputes the amount of used opioid and assigns class labels to the missing records based on the predicted values is used in this study.

The proposed pseudo-labeling method is tailored to the data collected in this study, where the data unused in the classification are utilized in the pseudo-labeling step. Specifically, among the three types of available data in the study: survey, prescription history, and EHR, most of the prescription history data, such as the number of refills, when a patient requests refill, and how many opioids are prescribed in the initial prescription or in the refills, can only be observed after the patients are discharged. Thus, such information cannot be used to build the final classification model to predict future opioid usages at the time of discharge. However, they can be used to infer the unanswered questions of opioid usage.

Thus, we first introduce a separate regression model based on prescription history data to predict the unanswered values of opioid usage amount. Then, using the predicted results and considering the uncertainty embedded in the prediction, the unlabeled data will be assigned with pseudo labels, which will be used in the classification model. Such a

framework is illustrated in Figure 4.2.

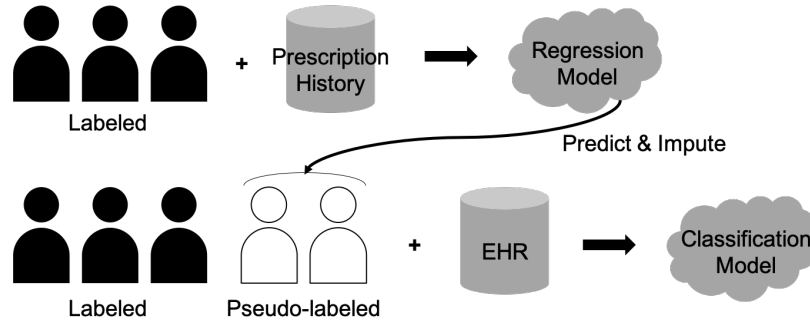


Figure 4.2: Diagram of semi-supervised learning process

The detail steps of pseudo-labeling and classification are presented below.

4.5.2 Pseudo-Labeling through a Probabilistic Model

During the pseudo-labeling step for the unlabeled data, for the patients who didn't answer the two-week survey, the amount of opioid consumption in two weeks post-discharge will be predicted and imputed.

Let matrix $\mathbf{X}_A = [\mathbf{x}_{A,1}, \dots, \mathbf{x}_{A,N_A}]^T$ and vector $\mathbf{y}_A = [y_{A,1}, \dots, y_{A,N_A}]^T$ be the prescription history data and the amount of opioids consumed within two weeks for the patients answering the survey, respectively, where N_A is the number of answered surveys. Analogously, for the patients who do not respond in the two-week survey post-discharge, matrix $\mathbf{X}_U = [\mathbf{x}_{U,1}, \dots, \mathbf{x}_{U,N_U}]^T$ characterizes the set of prescription history records and vector $\mathbf{y}_U = [y_{U,1}, \dots, y_{U,N_U}]^T$ refers to the amount of opioid consumption to be predicted, and N_U is the number of surveys with missing answers.

Different from the known values, the prediction can be uncertain. To reflect the fact that the imputed values are estimations, a probabilistic regression model is used to obtain the estimated values with uncertainties. Specifically, a Bayesian linear regression model is utilized, where the responses are assumed from a probability distribution, such as Gaussian distribution of the form:

$$p(y | \mathbf{X}, \mathbf{w}, \beta) = N(\mathbf{w}^T \mathbf{x}, \beta^{-1}), \quad (4.1)$$

where \mathbf{w} is the vector of weights and β is the precision (inverse of variance) of y .

In addition, parameters \mathbf{w} are also assumed to be sampled from a probability distribution, referred to as a prior probability distribution. Here, a Gaussian distribution is used again as it is the corresponding conjugate prior to (4.1). To simplify the model, we consider a simple form of Gaussian prior, which is a zero-mean isotropic Gaussian with a single precision parameter α :

$$p(\mathbf{w} | \alpha) = N(\mathbf{0}, \alpha^{-1}\mathbf{I}). \quad (4.2)$$

Then, the corresponding posterior distribution over \mathbf{w} is given by:

$$p(\mathbf{w} | \mathbf{y}_A, \mathbf{X}_A, \alpha, \beta) = N(\mathbf{m}_A, \mathbf{S}_A), \quad (4.3)$$

where mean \mathbf{m}_A and covariance \mathbf{S}_A are as follows:

$$\mathbf{m}_A = \beta \mathbf{S}_A \mathbf{X}_A^T \mathbf{y}_A \quad (4.4)$$

$$\mathbf{S}_A^{-1} = \alpha \mathbf{I} + \beta \mathbf{X}_A^T \mathbf{X}_A. \quad (4.5)$$

Note that since the posterior distribution is Gaussian, the maximum of a posteriori probability (MAP) estimate of \mathbf{w} is the mean \mathbf{m}_A [146].

After modeling the posterior and obtaining its MAP estimate using the collected data, the prediction of unknown $y_{U,j}$ in \mathbf{y}_U , $j = 1, \dots, N_U$, will be made, which is the amount of opioid consumption for a patient who has not answered the two-week survey with a record $\mathbf{x}_{U,j}$ in \mathbf{X}_U . To do this, we need to evaluate the predictive distribution defined by:

$$\begin{aligned} & p(y_{U,j} | \mathbf{x}_{U,j}, \mathbf{y}_A, \mathbf{X}_A, \alpha, \beta) \\ &= \int p(y_{U,j} | \mathbf{x}_{U,j}, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{y}_A, \mathbf{X}_A, \alpha, \beta) d\mathbf{w}. \end{aligned} \quad (4.6)$$

As $p(y_{U,j} | \mathbf{x}_{U,j}, \mathbf{y}_A, \mathbf{X}_A, \alpha, \beta)$ is the convolution of two Gaussian distributions, (4.1) and

(4.3), it takes a Gaussian form as well:

$$p(y_{U,j} | \mathbf{x}_{U,j}, \mathbf{y}_A, \mathbf{X}_A, \alpha, \beta) = N(\mathbf{m}_A^T \mathbf{x}_{U,j}, \sigma_A^2(\mathbf{x}_{U,j})), \quad (4.7)$$

where variance $\sigma_A^2(\mathbf{x}_{U,j})$ is given by:

$$\sigma_A^2(\mathbf{x}_{U,j}) = \frac{1}{\beta} + \mathbf{x}_{U,j}^T \mathbf{S}_A \mathbf{x}_{U,j}. \quad (4.8)$$

Note that the first term in (4.8) represents the noise in the dataset and the second one characterizes the uncertainty associated with \mathbf{w} [146].

Now, using the Bayesian linear regression model, not only the opioid consumption of unlabeled data, but also the prediction distributions, can be estimated. Then the predicted values of unlabeled data, which become the point estimate,

$$\hat{y}_{U,j}^{MAP} = \mathbf{m}_A^T \mathbf{x}_{U,j}, \quad j = 1, \dots, N_U,$$

will be used to assign pseudo labels in the dataset. Using them, a classification model can be developed.

However, it should be noticed that if wrong information is obtained by using uncertain pseudo-labeled data, it could degrade the prediction accuracy by misguiding the inference [147]. In other words, if a prediction of the unlabeled data is wrong, the regression step may not yield an improvement (or may even downgrade) over a supervised learning model using the labeled data only. Therefore, we need to exclude the results that can degrade the classification task by investigating the predictive distribution given by (4.7) rather than using the point estimate of $y_{U,j}$ only. In order to exclude the counterproductive pseudo-labels, the thresholds need to be setup to divide and define the patient groups. Then the classification models will be built afterwards.

Let T be a threshold to divide the opioid consumption level. If a predicted $\hat{y}_{U,j}$, $j =$

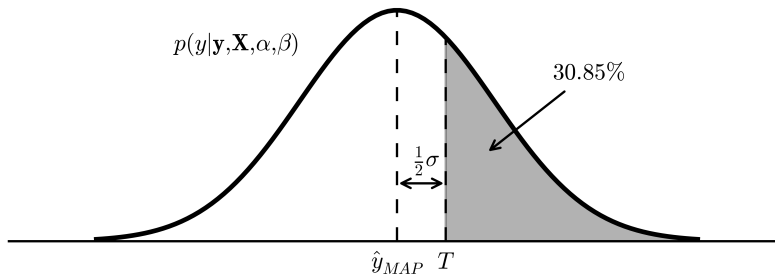


Figure 4.3: Probability of pseudo-labeling failure

$1, \dots, N_U$, is close to T and its corresponding variance $\sigma_A^2(\mathbf{x}_{U,j})$ is not small enough, the record $\mathbf{x}_{U,j}$ is more likely to be assigned to a wrong group. For example, as described in Figure 4.3, if the difference between T and $\hat{y}_{U,j}$ is within one half of the standard deviation, $\frac{1}{2}\sigma_A(\mathbf{x}_{U,j})$, the probability that the prediction gives a wrong label is larger than the shaded area, which is 30.85%. Such a probability of failure in pseudo-labeling can be unacceptably huge. Thus, after predicting the amount of opioid consumption and their corresponding variances, we will exclude the records having the MAP point estimation $\widehat{y}_{U,j}$ with distances from threshold T being less than $\frac{1}{2}\sigma_A(\mathbf{x}_{U,j})$.

The pseudo code of the pseudo-labeling procedure is provided in Algorithm 1, where vector $\mathbf{c} = [c_1, \dots, c_N]^T$ includes both \mathbf{c}_A and \mathbf{c}_U , representing the target values for a classification model in the next step, and $N = N_A + N_U$. If a pseudo-label $c_{U,j}$ is assigned as NULL by the algorithm, the pseudo-label and the corresponding record will be dropped in the following classification model training step.

4.5.3 Building a Classification Model

The first step of the building a classification model is to prepare the preprocessed EHR data \mathbf{Z} and the target values \mathbf{c} , where matrix $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_N]^T$, including both \mathbf{Z}_A and \mathbf{Z}_U . As mentioned before, the known or predicted amounts of consumed opioids \mathbf{y} will be used to construct the labels or pseudo-labels \mathbf{c} . To group the patients and assign labels to them, a classification threshold T needs to be set up, and it should be related to the purpose of

Algorithm 1: Pseudo-labeling using Bayesian linear regression

Input: Prescription history data of answering patients \mathbf{X}_A and non-answering patients \mathbf{X}_U , the amounts of used opioids for answering patients \mathbf{y}_A , hyperparameters α and β , and the classification threshold T_c

Output: Pseudo-labeled dataset $\mathbf{c}_U = [c_{U,1}, \dots, c_{U,N_U}]^T$
 Calculate the covariance of weights $\mathbf{S}_A^{-1} = \alpha \mathbf{I} + \beta \mathbf{X}_A^T \mathbf{X}_A$
 Calculate the mean of weights $\mathbf{w}^{MAP} = \beta \mathbf{S}_A \mathbf{X}_A^T \mathbf{y}_A$

```

for  $j = 1, \dots, N_U$  do
  Predict  $\hat{y}_{U,j} = (\mathbf{w}^{MAP})^T \mathbf{x}_{U,j}$ 
  Calculate variance  $\sigma_A^2(\mathbf{x}_{U,j}) = \frac{1}{\beta} + \mathbf{x}_{U,j}^T \mathbf{S}_A \mathbf{x}_{U,j}$ 
  if  $U_{,j} - T_c > \frac{1}{2} \sigma_A(\mathbf{x}_{U,j})$  then
    if  $\hat{y}_{U,j} < T_c$  then
      |  $c_{U,j} \leftarrow L$ 
    end
    else
      |  $c_{U,j} \leftarrow H$ 
    end
  end
  else
    |  $c_{U,j} \leftarrow \text{NULL}$ 
  end
end

```

classification. The classification results can then provide insights for the optimal opioid dosage and prescription plan to reduce the quantity of unused pills. Specifically, after building the classification model, if a patient is classified to the less opioid-use group, the patient will be recommended to receive less amount of opioids in the initial prescription at discharge. However, since there is no standard criterion or clinical guidance to define “less-opioid-use” or “more-opioid-use” patients, setting up the thresholds to divide the patient groups is critical, but not easy.

By conducting simple classification modeling experiments on the preprocessed data and receiving feedback from the healthcare providers in SMH, different threshold values have been investigated, and 200 MME, denoted as T_L , is selected as the threshold for separating less-opioid-use patients, and 700 MME, denoted as T_H , as the threshold for heavy usage patients. Note that the current practice is to prescribe opioids with 450 MME uniformly to most of patients. Since $T_L = 200$ MME is much smaller than the current practice or

prescription (450 MME), using this threshold for prescriptions determination can achieve substantial left-over reduction. Also, the number of opioid refill requests is expected to decrease by using $T_H = 700$, which is much higher than 450 MME. This will reduce the burden on providers and pharmacists as well.

Using these thresholds, the patients are divided into three groups: using opioids less than 200 MME, between 200 MME and 700 MME, and more than 700 MME within two weeks after discharge. To solve the multi-class classification problem, two independent classification models, where each threshold splits the patient groups for each model, are trained and used. Before training a classification model with $T \in \{T_L, T_H\}$, the unanswered data are assigned with pseudo-labels in accordance with threshold T . Then, classification models, such as Decision Tree and its ensemble methods, including Random Forest and XGBoost, as well as Logistic Regression, Support Vector Machine (SVM), k-Nearest Neighbors (K-NN) and Neural Network models, are trained.

To compare the model performance, we set aside some of the records as a test dataset. Specifically, for further comparisons between the proposed model involving the pseudo-labeled data and the supervised model over the originally labeled data, we select the test data records from the set of originally labeled data. The rest of data will compose the training set. However, the rate of patients using less than 200 MME opioids is relatively small (23%), which makes the data imbalanced and let those patients hard to be detected. Since it is more important to identify less-opioid users considering the purpose of the classification, we over-sample the records of the minority class in training data to balance the classes to overcome the issues of data imbalance. In the case of classification using $T_H = 700$ MME, the original training data are used without over-sampling because, as mentioned above, classifying less-opioid users correctly is more important in reducing opioid over-prescriptions.

With the training data, regardless of whether they are over-sampled or not, we fit them with different classification algorithms and choose the best parameters for each algorithm

using five-fold cross-validations. The parameters for each model tuned in the experiments are as follows: the number of neighbors in K-NN; penalty (L1, L2), regularization strength, and class weight to be balanced or not in Logistic Regression; kernel function, regularization strength, and class weight to be balanced or not in SVM; the number of neurons in the hidden layer in Neural Network; the number of trees, split criterion measure (gini, entropy), and class weight to be balanced or not in Random Forest; the maximum depth of trees and class weight in XGBoost. Other parameters in each model are set through default values or set automatically by the algorithm.

4.6 Results

4.6.1 Pseudo-Labeling with Bayesian Linear Regression

To build the Bayesian linear regression model, the features related to prescription history, such as the amount of initial prescription, the amount of opioid usage in the first, second, and third refills if any, and the number of days passed from initial prescription to refills, are used, along with a few characteristic features, for instance, procedure type and gender. Figure 4.4 illustrates the MAP estimate (blue dots) of the feature weight and its corresponding standard deviation (blue bars) for each feature. Since the posterior distribution of weight \mathbf{w} is Gaussian, the vertical lines with MAPs as centers can depict the posterior distributions while neglecting the correlations between the weights.

From Figure 4.4, we can notice that the weights of initial amount of prescription (“MME initial”), the amount of first short acting opioid refill (“MME SArefill1”), the amount of first long acting opioid refill (“MME LArefill1”), the amount of facility administered medicine before discharge (“FAM”), and the average opioid consumption per day in the first period (“Avg use 1”), are significant and reliable in the resulting model. Note that the “Avg use 1” is obtained by dividing the initially prescribed opioid amount with the number of days passed from initial prescription to first refill. As we can see, the amount

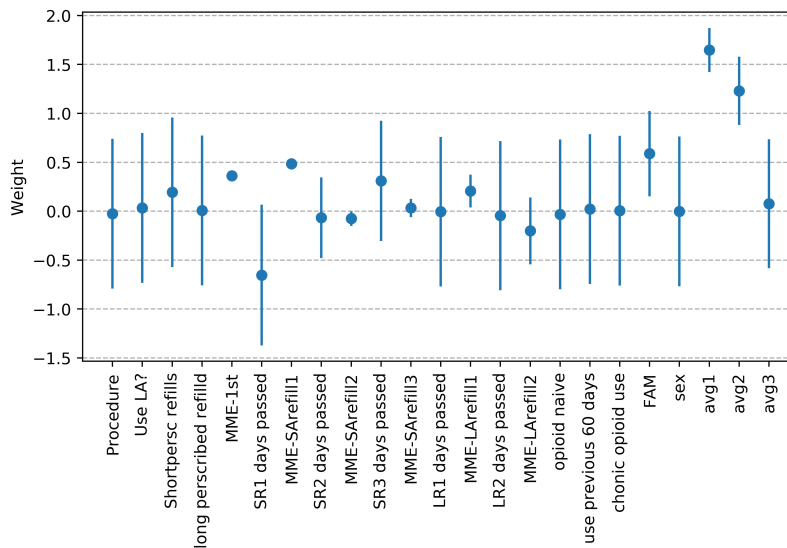


Figure 4.4: Posterior distributions of regression coefficients

of opioid consumption in two weeks are largely related to the initial prescription and the amount of refills, and it is proportional to the estimates of daily usage. Although some features have small weights in the model, we include all the features shown in the figure in the pseudo-labeling step, since mean squared errors on both training and test data and R^2 score show better values when those feature are included than excluded.

Figure 4.5 depicts the distribution of opioid consumption in two weeks obtained from the survey, and the distribution of the predicted opioid consumption of patients who didn't take the survey. The distribution of the predicted results appears to be consistent with the original distribution of opioid consumption, which can infer that the proposed model captures the intrinsic patterns in the data.

4.6.2 Classification Results

The pseudo-labeling step enables us to build a model using more data records, which is expected to help improve the model performance. The accuracy and AUC scores of the resulting models in classifying the patients in the test set are investigated. The results of models with thresholds 200 MME and 700 MME are given in Tables 4.2 and 4.3, respectively.

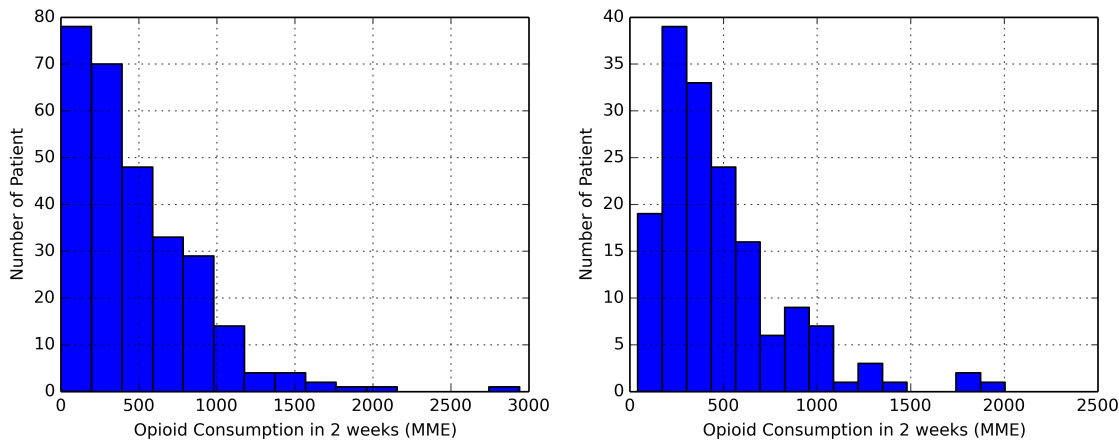


Figure 4.5: Histogram of opioid consumption of originally labeled data (left) and predicted values (right)

For comparison purpose, four additional models with different strategies dealing with missing labels are also evaluated: the model using the complete records only (referred to as *Model 1*), the model imputing the missing amounts of used opioids with mean value and assigning labels (*Model 2*), the model conducting Lasso regression on other features in training data and imputing missing values with predictions and assigning labels (*Model 3*), and the model following the proposed pseudo-labeling without filtering unreliable labels (*Model 4*).

The accuracy is computed based on the same test data that consist of 60 patients randomly selected from the labeled data for all models. All hyperparameters and the best subset of features for each model are tuned by five-fold cross-validations based on the training set, including the rest of labeled data (225 records) and the pseudo-labeled data (where the number of records depends on which pseudo-labeling strategy is applied), and are chosen when the corresponding model achieves the highest cross-validation accuracy among all the models whose resulting AUC scores are higher than 0.75. When a method cannot achieve the 0.75 AUC in 5-fold cross-validation, the method should be ignored in further testing. Nevertheless, for reference purpose, the results of testing accuracy for the models showing low cross-validation AUCs are also included (with parenthesis in the

Table 4.2: Classification Results with threshold 200 MME

Method	Proposed		Model 1		Model 2		Model 3		Model 4	
	accuracy	AUC	accuracy	AUC	accuracy	AUC	accuracy	AUC	accuracy	AUC
K-NN	0.70	0.71	0.68	0.63	0.62	0.50	0.70	0.72	0.65	0.66
Logistic Regression	0.70	0.68	0.67	0.62	0.75	0.73	0.70	0.70	0.67	0.69
SVM	0.82	0.75	0.73	0.66	0.73	0.64	0.73	0.68	0.70	0.64
Neural Network	0.80	0.70	0.73	0.66	0.73	0.68	0.73	0.74	0.75	0.71
Random Forest	0.83	0.65	0.80	0.63	0.78	0.59	0.78	0.67	0.77	0.66
XGBoost	0.83	0.72	0.80	0.74	0.78	0.71	0.78	0.67	0.78	0.73

Table 4.3: Classification Results with threshold 700 MME

Method	Proposed		Model 1		Model 2		Model 3		Model 4	
	accuracy	AUC	accuracy	AUC	accuracy	AUC	accuracy	AUC	accuracy	AUC
K-NN	0.73	0.60	0.73	0.62	(0.78)	0.61	0.77	0.62	(0.72)	0.59
Logistic Regression	0.80	0.66	0.75	0.59	(0.80)	0.63	0.78	0.67	0.77	0.62
SVM	0.77	0.62	(0.75)	0.58	(0.73)	0.50	0.75	0.58	(0.75)	0.58
Neural Network	0.82	0.70	(0.78)	0.69	(0.80)	0.63	0.82	0.69	(0.72)	0.50
Random Forest	0.82	0.68	0.75	0.59	(0.73)	0.50	0.78	0.64	0.75	0.59
XGBoost	0.82	0.69	0.80	0.68	(0.80)	0.63	0.82	0.69	0.80	0.66

tables).

As shown in Table 4.2, XGBoost and Random Forest models achieve the best test accuracy of 83% among all the models trained through a proposed framework when the threshold is $T_L = 200$ MME. When identifying patients who use more than $T_H = 700$ MME of opioids in two weeks, the Random Forest model using the proposed pseudo-labeling shows the best performance with 82% accuracy, which is presented in Table 4.3. Most of the models using the proposed pseudo-labeling surpass the performance of those trained using the labeled data only or those using other imputation schemes.

To facilitate a comparison between the proposed scheme and standard imputation approaches, the proposed pseudo-labeling, no imputation, single imputation, regression imputation, and labeling without filtering (i.e., *Models 1-4*, respectively), are compared, and the best accuracy in each machine learning algorithm is colored in both tables. In both tables, the models using the proposed pseudo-labeling method (*Proposed*) performs better than all others in general except in a few models resulting deficient performances. Only in one case, logistic regression in $T = 200$ classification shows higher accuracy when mean value imputation (*Model 2*) is used, and the result of K-NN in $T = 700$ classification

is better when regression imputation (*Model 3*) is applied. However, the performances of other pseudo-labeling methods are showing significantly impaired. Furthermore, the comparisons between the models using the proposed pseudo-labeling with filtering and without filtering underline the superiority of the proposed model over the classic approaches.

Remark 4.1. *Tables 4.2 and 4.3 show the frequent dominance of semi-supervised learning with the proposed pseudo-labeling method over other labeling methods and pure supervised learning, which is an indication of superiority of the proposed method. A rigorous proof of the soundness and superior power of the proposed algorithm can be included in future study focusing on methodology.*

4.6.3 Important Risk Factors

The permutation importance of features in the top five models (XGBoost, Random Forest, SVM, Neural Network, and Logistic Regression) learned in the proposed framework are investigated in order to identify the features playing more important roles in the classification models. The permutation importance of a feature is obtained by measuring the decreasing pattern in model accuracy when the feature is unavailable or omitted [148]. Tables 4.4 and 4.5 provide the list of top seven features arranged in the order of highest permutation feature importance score for each of the five models. The blank cells in both tables indicate that the rest of features do not show significant importance scores.

Table 4.4: Top 7 important features in the resulting models with threshold 200 MME

	XGBoost	Random Forest	SVM	Neural Network	Logistic Regression
F1	FAM	FAM	FAM	FAM	FAM
F2	smoker	other chronic pain	pulse	pulse	pulse
F3	depression	opioid naïve	age	smoker	BMI-obesity
F4	BMI (binary)	smoker	weight	sex	depression
F5	surgical revision	preoperative use	height	BMI-obesity	preoperative use
F6	other chronic pain	HL or migraine	smoker	other chronic pain	other chronic pain
F7	HL or migraine	depression	other chronic pain		HL or migraine

As shown in both tables, the amount of “FAM” before discharge is identified as the most important feature across all the classifiers and thresholds. This fact implies that the

Table 4.5: Top 7 important features in the resulting models with threshold 700 MME

	Random Forest	XGBoost	Neural Network	Logistic Regression	SVM
F1	FAM	FAM	FAM	FAM	FAM
F2	pulse	sex	chronic use	chronic use	sex
F3	age	chronic use	pulse	age	anxiety
F4	weight	alcohol	weight		BMI (binary)
F5	preoperative use	smoker	surgical experience		pulse (binary)
F6	height	depression	opioid naïve		chronic use
F7	chronic use	BMI (binary)			surgical experience

amount of facility administered medicine before discharge, which is the amount of opioids that a patient was received while being hospitalized, highly affects the total amount of opioid taken within two weeks after discharge. In addition, there are features that are common in one of the tables but not shown in others. Those features can be interpreted as the important features for a given classification corresponding to a certain threshold. For example, in Table 4.4, “smoker” is identified as one of the most important features in all but Logistic Regression for identifying less-opioid-use patients, while it is not influential when classifying more-opioid-use patients (see Table 4.5). It may infer that the patients who do not smoke tend to use less opioids or avoid using opioids. However, it does not distinguish between moderate and heavy opioid users. This result coincides with the conclusion in [43] that smoking is a risk factor of persistent opioid use.

It is also observed in Table 4.4 that whether a patient has other chronic pain, such as back pain, in addition to post-surgical pain, is a crucial risk factor to identify patients who will use less than 200 MME opioids. Meanwhile, Table 4.5 shows that whether a patient has long history of frequent opioid prescription prior to the surgery (“chronic use”) is related to the identification of patients who are more likely to use more than 700 MME opioids. From this fact, it can be inferred that patients with other types of chronic pain are more likely to easily consume opioids than who do not have them so that having other chronic pain can be an indicator of opioid less users. On the other hand, chronic opioid use can be a pointer of the expected overdose rather than a sign of using or not. Therefore, for those who have such signs, higher amount of opioid prescription can be available with

ensuring more careful management.

4.7 Conclusions

In this Chapter, we introduce a case study to classify patients' opioid usage at SSM Health St. Mary's Hospital. Specifically, a machine learning classification model is developed to identify patients expected to use less opioids and to detect opioid over-users within two weeks after undergoing total joint replacement surgeries. To obtain the data regarding their opioid usage patterns in a short period of time, a survey study was first conducted. The resulting values of the amounts of used opioids work as a target variable. Variables related to patient characteristics are extracted from the patient's electronic health records and used as features in the classification models.

To deal with the missing target values, a Bayesian linear regression based pseudo-labeling procedure is proposed. The pseudo-labeling first predicts the expected amount of opioid usage and assigns pseudo labels to the patients based on the predictions and their uncertainties. Then, the pseudo labeled data are mixed with the originally labeled data to build a classification model. Such an approach is not only useful for prediction of opioid usage for TJR patients, but also applicable to many other scenarios where a target variable in a classification model has an underlying continuous variable with abundant missing values.

Chapter 5

Intervention Planning for COPD Patients Post Discharge

5.1 Introduction

This chapter is devoted to modeling and analysis of intervention process to reduce hospital readmissions for patients with chronic obstructive pulmonary disease (COPD). As compliance is a major issue among COPD patients, and economic burden and social support can be important factors affecting compliance and readmission, we propose to hospital management to reimburse out-of-pocket and transportation costs for COPD patients visiting primary care physicians (PCPs) and respiratory rehabilitation centers, which are used as incentives to encourage them complying with patient-specific intervention plan. It is expected that the increased compliance level can lead to reduced readmissions, whose savings in readmission cost and penalties will exceed the incentive cost. Then we introduce an optimization model to minimize COPD readmission rate under incentive budget constraint and patients' readmission risks. Solving the problem, the minimal readmission rates are evaluated and the conditions to achieve the optimal solution are derived, which can provide a guideline for hospital management to plan appropriate incentive budget

and benchmark the desired readmission rate. A case study at a community hospital is presented to illustrate the method. Finally, cost-effective analysis, sensitivity studies, and implementation discussions are carried out.

5.2 COPD Intervention Process

5.2.1 Intervention Process

To reduce COPD readmission, effective patient specific interventions are needed. Thus, analysis of the intervention process is carried out first. In this chapter, using the post-discharge interventions at St. Mary's Hospital as an example, we introduce an optimization framework.

Specifically, at St. Mary's Hospital, when a COPD patient is discharged, he/she is instructed to take medication, use home equipment (e.g., inhalers, nebulizers), and quit smoking. At discharge, a COPD patient will be scheduled to have a visit to the primary care physician (PCP) and multiple services at rehabilitation center (usually twice a week). Within 72 hours after discharge, a nurse phone call is expected to check the patient's status. A pulmonary clinic visit should also be scheduled, but usually needs to wait for more than a month. Similar intervention processes can be observed in many other hospitals.

However, some patients may not follow these instructions so that compliance becomes an important issue [15]. Reasons for non-compliance can be due to economic burden and lack of social support [16], such as patient's inability to finding a PCP, lacking transportation, or being unable to afford to the rehab and office visit co-pays. As a result, some patients may skip one or all the intervention services. Therefore, four possible scenarios exist in the intervention process (as shown in Figure 5.1): no compliance to intervention, visiting PCP only, visiting rehab center only, and visiting both rehab center and PCP. As one can understand, the ideal scenario is that all patients comply to interventions to visit both rehab centers and PCP offices.

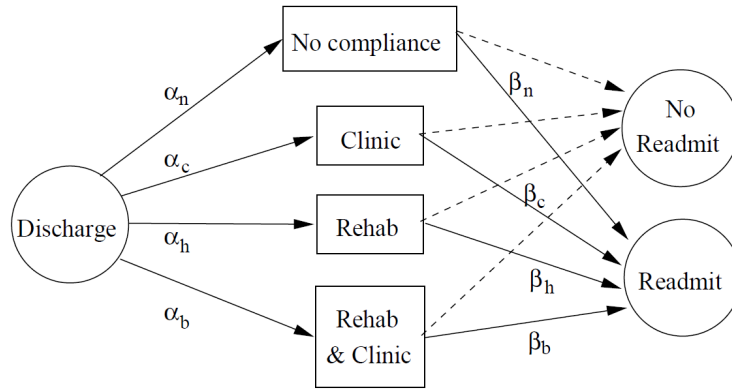


Figure 5.1: Intervention flow model

In Figure 5.1, α_i^k represents the percentage of patients with readmission risk level k following scenario i , where $i = n, c, h, b$, representing no compliance, PCP visit only, rehab visit only, and both rehab and PCP visits, respectively, and $k = H, L$, indicating high and low risks for readmission, respectively; while β_i^k characterizes the probability to readmission for patients with risk level k in scenario i .

In order to encourage patients to comply with these interventions, the hospital has been proposed to provide services or incentives for patients to improve their compliance to provider's instructions, for example, reimbursing out-of-pocket expenses for interventions such as co-pays or supporting transportation cost. Through such intervention encouragement, it is expected that patients' percentages following the intervention instructions will be improved and the resulting readmission rate will be reduced. Existing literature has suggested that PCP visits and rehab services can lead to lowered readmissions [149, 150].

Then questions arise: What will be the maximal benefit, i.e., the smallest readmission rate, under an incentive budget constraint? What will be the conditions or configurations of patient compliance to interventions when the optimal readmission rate is achieved? Answers to the first question enable us to understand the best outcome we can attain given a certain budget level. Answers to the second question could make us realize where are the differences comparing to the ideal scenario. Then, the management can use these answers to determine how much incentive to provide in order to reach a desired readmission rate,

and plan additional efforts targeting to the specific patients to comply to interventions.

5.2.2 Optimal Intervention Model

In this chapter, we focus on optimizing the benefit of incentives, i.e., finding the minimal readmission rate for a given incentive budget.

Clearly, the readmission probabilities will be different for patients with different risk levels taking different intervention activities. In other words, both the COPD patients' risk levels for readmission and the associate possibility in compliance to physician's instruction will affect the readmission rate. Therefore, based on COPD patients' records extracted from hospital database, we first identify the significant factors associated with COPD readmission. Then, using these factors, a prediction model is introduced to categorize COPD patients into high- ($k = H$) and low-risk ($k = L$) groups.

For each group, the readmission probabilities in each intervention scenario (i.e., β_i^k , $i = n, c, h, b$) will be estimated. In other words, a COPD patient with readmission risk k taking intervention i has probability β_i^k to be readmitted into the hospital within 30 days, and $1 - \beta_i^k$ not to be readmitted. Based on observations and data analysis we assume that $\beta_b^k < \beta_h^k < \beta_c^k < \beta_n^k$ for both risk groups.

The cost to encourage a patient to take intervention i is denoted as C_i . Particularly, the average co-pay and/or transportation cost for a patient to visit PCP is denoted as C_c ; while the average cost for taking rehab services during the 30-day period is C_h . Then the cost of taking both interventions is defined by C_b . Due to frequent visits to rehab centers, we assume $C_c < C_h < C_b$. Note that these costs are mainly due to economic burden so that they are indifferent to patient's health status. Thus, the two risk groups share the same cost values.

Now, define the number of COPD patients within the system by N . Assume the proportion of high-risk patients among all patients is $\rho < 1$. Then the expectation of COPD

patients' readmission rate P_r can be calculated as

$$P_r = \sum_{i \in \{n, c, h, b\}} \left(\rho \alpha_i^H \beta_i^H + (1 - \rho) \alpha_i^L \beta_i^L \right), \quad (5.1)$$

where, as explained before, α_i^k is the proportion of patients in group k who take the i -th intervention through the hospital incentive program. The expected incentive cost of the hospital to reimburse COPD patients is given by

$$C = N \sum_{i \in \{c, h, b\}} \left(\rho C_i \alpha_i^H + (1 - \rho) C_i \alpha_i^L \right). \quad (5.2)$$

Denote the hospital's available budget to pay incentives to cover the out-of-pocket and transportation costs of COPD patients who take the interventions as C_d . Then the following problem can be formulated:

$$\begin{aligned} & \text{minimize} && P_r = \sum_{i \in \{n, c, h, b\}} \left(\rho \alpha_i^H \beta_i^H + (1 - \rho) \alpha_i^L \beta_i^L \right) \\ & \text{such that} && N \sum_{i \in \{c, h, b\}} \left(\rho C_i \alpha_i^H + (1 - \rho) C_i \alpha_i^L \right) \leq C_d, \\ & && \sum_{i \in \{n, c, h, b\}} \alpha_i^j = 1, \quad j \in \{H, L\} \\ & && 0 \leq \alpha_i^j \leq 1, \quad 0 \leq \beta_i^j \leq 1, \\ & && \text{for all } i \in \{n, c, h, b\}, \quad j \in \{H, L\}. \end{aligned} \quad (5.3)$$

The solution of this problem can provide benchmarks of the best outcome a hospital can achieve in reducing readmission rates under given incentives. In addition, understanding the conditions to reach such an outcome can indicate how much differences in terms of patient compliance comparing with the desired (full compliance) case. Then, by evaluating the results due to different incentive budgets, hospital management can determine the appropriate incentive level based on their target for readmission rate.

5.3 Optimality Analysis

In this section, optimality analysis is carried out by solving the preceding linear programming problem. To find optimal solutions, we introduce a graphical framework which solves the problem and graphically represents the impacts of parameters and budget value on the readmission results so that it facilitates sensitivity analysis.

Specifically, in this framework, we project the problem into a 2 dimensional space. We first describe the problem in a $2D$ space when $\rho = 1$, i.e., when there is only one risk group (e.g., high-risk patients). Let the x - and the y -axes represent the cost and readmission probability, respectively. Each intervention is represented by its corresponding point that consists of its cost and the resulting readmission probability. Considering that the sum of the proportion variables is exactly 1 and those variables are the coefficients of the equations, the formulas for P_r and C are convex combination of the four points representing the four intervention types. We refer to such a region of convex combinations as a feasible region of the problem (see either of the two pink quadrangles in Figure 5.2). Thus, the original problem is transformed into a problem of identifying the minimum value of y in this region while satisfying the constraint on x value.

Next, when $\rho > 0$, the feasible region can be projected in the $2D$ space and this feasible region forms an octagon (the middle green shaded area in Figure 5.2), of which eight sides are respectively parallel to one of the sides of two quadrangles (the upper and lower pink shaded areas in Figure 5.2) corresponding to two risk groups with high-risk (the upper quadrangle) and low-risk (the lower quadrangle) patients.

5.3.1 Piecewise Linear Function

The minimized readmission probabilities of the possible budget values per patient (denoted as $x = \frac{C_d}{N}$) are represented as a piecewise linear function in terms of x , whose graph is constructed by the lower sides of the octagon. Using Algorithm 2, we introduce a piecewise

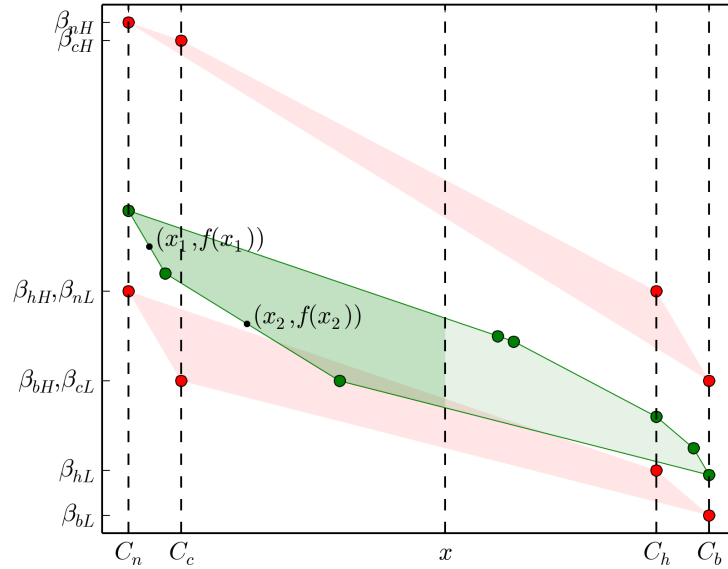


Figure 5.2: Projection of the problem in a 2D space

linear function $f(x)$.

5.3.2 Optimal Solutions

Using the above algorithm, the optimal readmission rate function can be constructed and a specific value of readmission rate under a given budget value can be obtained. When the budget value is set and the corresponding readmission rate is obtained, an optimal solution, which characterizes the optimal incentive condition, can be generated. Such a condition is derived according to which region the unit budget (x) belongs to. If the point of cost x and its optimal readmission rate $f(x)$ (i.e. point $(x, f(x))$) is the point where the two edges intersect, the optimal solution will consist of allocating one intervention for each risk group. If it lies on an edge, the optimal solution includes one intervention for one risk group and two interventions for another risk group, where the ratio between the two interventions is decided based on the position of the point.

In the resulting piece-wise linear function, each piece (edge) corresponds to an intervention for either high-risk patient group or low-risk patient group. Tracing the resulting

Algorithm 2: Readmission Rate Function Constructing Algorithm

Step 1: For each patient group, calculate the magnitudes of slopes of all possible edges between two points within the group. Specifically,

1a) Let $S_{i,j}^k$, $i, j \in \{n, c, h, b\}$, $k \in \{H, L\}$, be the magnitude of slope of the edge between (C_i, β_i^k) and (C_j, β_j^k) , which are the points of group k .

1b) Calculate $S_{i,j}^k$ as follows:

$$S_{i,j}^k = \left| \frac{\beta_i^k - \beta_j^k}{C_i - C_j} \right|. \quad (5.4)$$

Step 2: Find all lower part edges of the feasible regions for high-risk and low-risk groups using the following conditional statements.

2a) If $S_{n,b}^k = \max\{S_{n,i}^k : i \in \{c, h, b\}\}$, then $(C_n, \beta_n^k) - (C_b, \beta_b^k)$ is the only bottom edge of group k .

2b) If $S_{n,h}^k = \max\{S_{n,i}^k : i \in \{c, h, b\}\}$, then $(C_n, \beta_n^k) - (C_h, \beta_h^k)$ and $(C_h, \beta_h^k) - (C_b, \beta_b^k)$ are the lower part edges of group k .

2c) Otherwise,

- if $S_{c,b}^k = \max\{S_{c,i}^k : i \in \{h, b\}\}$, then $(C_n, \beta_n^k) - (C_c, \beta_c^k)$ and $(C_c, \beta_c^k) - (C_b, \beta_b^k)$ are the lower part edges of group k .
- else, $(C_n, \beta_n^k) - (C_c, \beta_c^k)$, $(C_c, \beta_c^k) - (C_h, \beta_h^k)$ and $(C_h, \beta_h^k) - (C_b, \beta_b^k)$ are the lower part edges of group k .

Step 3: Sort the edges according to the steepness. Specifically, combine the lists of edges in two groups into one list and sort the edges in the list in descending order of the magnitudes of slopes.

Step 4: Iteratively construct the function equation.

4a) Initialization: Set the initial interventions of each group as $g_H = g_L = n$ and starting point $(x_{\text{curr}}, y_{\text{curr}})$ as $(C_n, \rho\beta_{g_H}^H + (1 - \rho)\beta_{g_L}^L)$, which is the left-most point and set the initial function as $f(x) = 0$.

4b) Take out the steepest edge from the list and let this edge be $(C_i, \beta_i^k) - (C_j, \beta_j^k)$. Then, β_i^k must be either $\beta_{g_H}^H$ or $\beta_{g_L}^L$.

4c) If $k = H$, then change the value of g_H to the value of j . If $k = L$, then change the value of g_L to the value of j .

4d) Set the next point $(x_{\text{next}}, y_{\text{next}})$ as $(\rho C_{g_H} + (1 - \rho)C_{g_L}, \rho\beta_{g_H}^H + (1 - \rho)\beta_{g_L}^L)$

4e) Add the following term into $f(x)$:

$$\left(y_{\text{curr}} + \frac{\beta_i^k - \beta_j^k}{C_i - C_j} (x - x_{\text{curr}}) \right) I(x_{\text{curr}} \leq x < x_{\text{next}}),$$

where $I(\cdot) = 1$ when $x_{\text{curr}} \leq x < x_{\text{next}}$, otherwise $I(\cdot) = 0$.

4f) Set $(x_{\text{curr}}, y_{\text{curr}}) = (x_{\text{next}}, y_{\text{next}})$ and repeat this step from 4b) until the list becomes empty. \square

function from the left-most point until approaching the budget value x , we can update the optimal solution. If budget x and the corresponding optimal readmission rate $f(x)$ are on the left-most point, then the optimal solution consists of no intervention for both risk groups. If the point $(x, f(x))$ is on the first edge that corresponds to intervention i of risk group k , then the optimal solution implies that some patients in risk group k take intervention i while the others take no intervention (e.g., $(x_1, f(x_1))$ in Figure 5.2, where $k = L, i = c$).

When another end point of the edge is reached, the optimal solution changes to that all patients in risk group k take intervention i , and patients in the other risk group takes no intervention. In the next edge, the optimal solution suggests that some patients in risk group k' go for intervention i' while the others remain unchanged, if the edge corresponds to intervention i' of risk group k' (such as $(x_2, f(x_2))$ in Figure 5.2, where $k' = H, i' = b$).

The optimal result is updated by repeating this procedure following the result function until reaching the budget value x .

5.3.3 Special Cases with Closed Form Solutions

In general, the optimal readmission rate and incentive allocation can be found by constructing the piecewise linear function. However, in some special cases, there exist closed-form solutions. Below, we introduce such a special case, which occurs when both high-risk and low-risk groups have the same lowest edges, i.e. when the interventions in both groups have the same order of cost-effectiveness.

Specifically, consider a case where the lowest edge of both groups is $(C_n, \beta_n) - (C_b, \beta_b)$. This implies that taking both interventions is the most cost-effective choice for all patients. Generally, rehabilitation is practiced periodically in order to manage a patient's health status, so that it reduces readmission rate more than that for PCP visit. On the other hand, PCP visit is usually a one-time intervention, so that it costs less than rehab services. However, if a patient take both interventions, there can be a synergy effect, thereby it can

bring the biggest cost-effective readmission reduction. To enumerate the candidates for optimal intervention combinations according to the budget value, we can divide the overall cost range into 4 parts. The optimal solution and its practical meaning for each case are explained below:

- (i) $x \geq C_b$. In this case, visiting both rehab center and PCP will minimize the readmission rate, since both β_b^H and β_b^L are the smallest readmission probabilities in each groups. In other words, the optimal readmission rate is

$$P_r^* = \rho\beta_b^H + (1 - \rho)\beta_b^L,$$

where

$$\alpha_b^H = \alpha_b^L = 1.$$

- (ii) $(1 - \rho)C_b \leq x \leq C_b$. If the budget is not enough to cover all patient taking both interventions, the optimal readmission rate will be higher since it implies that some patient are not taking any intervention. Since reimbursement to all patients taking both interventions is not an available option, the best possible result is better to split the budget to cover some patients taking one of the two interventions. To minimize readmission rate, it should cover the patients who have larger difference in readmission rates between taking both interventions and taking nothing. In other words, the optimal intervention solution will lead to more cost-effective readmission reduction. Then, we obtain

$$P_r^* = \begin{cases} \rho\beta_b^H + (1 - \rho)(\beta_b^L\alpha_b^L + \beta_n^L\alpha_n^L), & \text{if } \frac{\beta_n^H - \beta_b^H}{C_b} \geq \frac{\beta_n^L - \beta_b^L}{C_b}, \\ \rho(\beta_b^H\alpha_b^H + \beta_n^H\alpha_n^H) + (1 - \rho)\beta_b^L, & \text{otherwise,} \end{cases}$$

where

- if $\frac{\beta_n^H - \beta_b^H}{C_b} \geq \frac{\beta_n^L - \beta_b^L}{C_b}$,

$$\alpha_b^H = 1, \quad \alpha_b^L = \frac{\tilde{C}_d - \rho C_b}{(1 - \rho)C_b}, \quad \alpha_n^L = 1 - \alpha_b^L,$$

- otherwise,

$$\alpha_b^L = 1, \quad \alpha_b^H = \frac{\tilde{C}_d - (1 - \rho)C_b}{\rho C_b}, \quad \alpha_n^H = 1 - \alpha_b^H.$$

(iii) $\rho C_b \leq x < (1 - \rho)C_b$. In this budget range, it is even not enough to cover all low-risk patients taking both interventions. Thus, the optimal solution includes some low-risk patients taking nothing. In fact, when all high-risk patients take both interventions and less low-risk patients take both, or when as many as possible low-risk patients take both while high-risk patients take nothing, the optimal solution can be obtained. For example, if high-risk patients have smaller difference in readmission rate between the two interventions, the second scenario leads to the smallest readmission rate. Therefore, the optimal solution is

$$P_r^* = \begin{cases} \rho(\beta_b^H \alpha_b^H + \beta_n^H \alpha_n^H) + (1 - \rho)\beta_n^L, & \text{if } \frac{\beta_n^H - \beta_b^H}{C_b} \geq \frac{\beta_n^L - \beta_b^L}{C_b}, \\ \rho\beta_n^H + (1 - \rho)(\beta_b^L \alpha_b^L + \beta_n^L \alpha_n^L), & \text{otherwise,} \end{cases}$$

where

- if $\frac{\beta_n^H - \beta_b^H}{C_b} \geq \frac{\beta_n^L - \beta_b^L}{C_b}$,

$$\alpha_b^H = 1, \quad \alpha_b^L = \frac{\tilde{C}_d - \rho C_b}{(1 - \rho)C_b}, \quad \alpha_n^L = 1 - \alpha_b^L,$$

- otherwise,

$$\alpha_b^H = 1, \quad \alpha_b^L = \frac{\tilde{C}_d}{(1 - \rho)C_b}, \quad \alpha_n^L = 1 - \alpha_b^L.$$

(iv) $x < \rho C_b$. In this range, for both groups, it is not possible to cover all patients taking

both. Thus, the optimal solution consists of some patients of one of the two groups taking both interventions and the rest taking nothing. When patients who experience more difference in readmission rates between the two options take both interventions, the optimal solution can be obtained.

$$P_r^* = \begin{cases} \rho(\beta_b^H \alpha_b^H + \beta_n^H \alpha_n^H) + (1 - \rho)\beta_n^L, & \text{if } \frac{\beta_n^H - \beta_b^H}{C_b} \geq \frac{\beta_n^L - \beta_b^L}{C_b}, \\ \rho\beta_n^H + (1 - \rho)(\beta_b^L \alpha_b^L + \beta_n^L \alpha_n^L), & \text{otherwise,} \end{cases}$$

where

- if $\frac{\beta_n^H - \beta_b^H}{C_b} \geq \frac{\beta_n^L - \beta_b^L}{C_b}$,

$$\alpha_n^L = 1, \quad \alpha_b^H = \frac{\tilde{C}_d}{\rho C_b}, \quad \alpha_n^H = 1 - \alpha_b^H,$$

- otherwise,

$$\alpha_n^H = 1, \quad \alpha_b^L = \frac{\tilde{C}_d}{(1 - \rho)C_b}, \quad \alpha_n^L = 1 - \alpha_b^L.$$

Figure 5.3 presents a summary of this analysis, where the candidates of possible solutions according to each range of budget value are illustrated. The texts above the vertical lines represent the cost value of the case that all high-risk patients get the intervention on the upper one and all low-risk patients get the intervention on the lower one. The texts along with the arrow indicate intervention combinations, left one for high-risk, and right one for low-risk patients. For example, “Both, None & Both” implies that all high-risk patients take both interventions, while for low-risk patients, some take both and some take no intervention. The double directional arrows indicate that the changes can be made in either direction. Note that the cost is increasing from left to right, but the location of each point is not related to readmission risk.

Similar analyses for other cases are also illustrated. Figure 5.4 represents the candidates

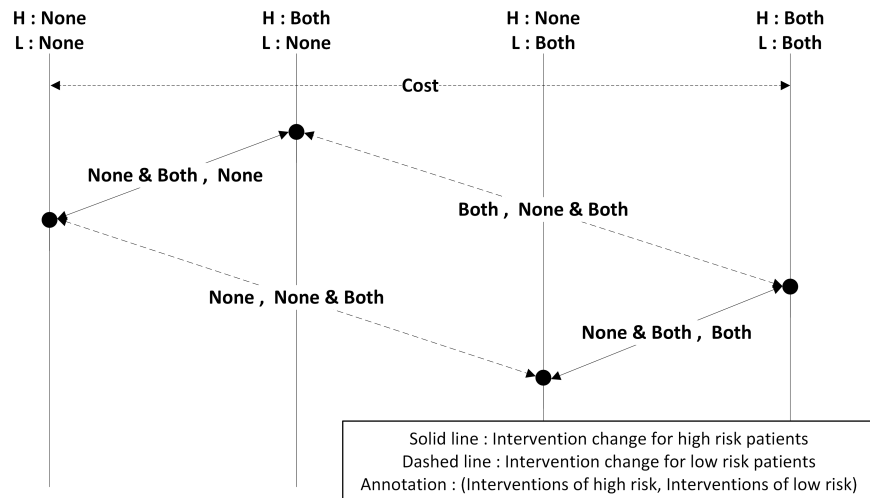


Figure 5.3: Illustration of candidates for optimal combination of interventions

of possible solutions in each budget range when rehabilitation is the most cost-effective intervention. In this case, as the budget value increases, incentives can be provided to some patients for rehab services, and the choice is dependent on the effectiveness of the intervention in each group. If the reduced readmission rate through rehab is higher for the high-risk group than that for the low-risk group, encouraging the high-risk patients to visit rehab center is more beneficial.

When PCP visit is the most cost-effective intervention, the candidates of possible solutions are illustrated in Figures 5.5 and 5.6. In the case of Figure 5.5, taking both is the second most cost-effective intervention, while in the case of Figure 5.6 taking rehab service is the second. In both cases, similar to the scenarios in Figure 5.4, as the budget value increases, encouraging some patients to visit PCP is beneficial. In Figure 5.5, as the budget value increases further, providing incentives to patients taking both interventions will be the next optimal solution. While in Figure 5.6, covering rehab visits is the next optimal solution and taking both interventions will be the last one.

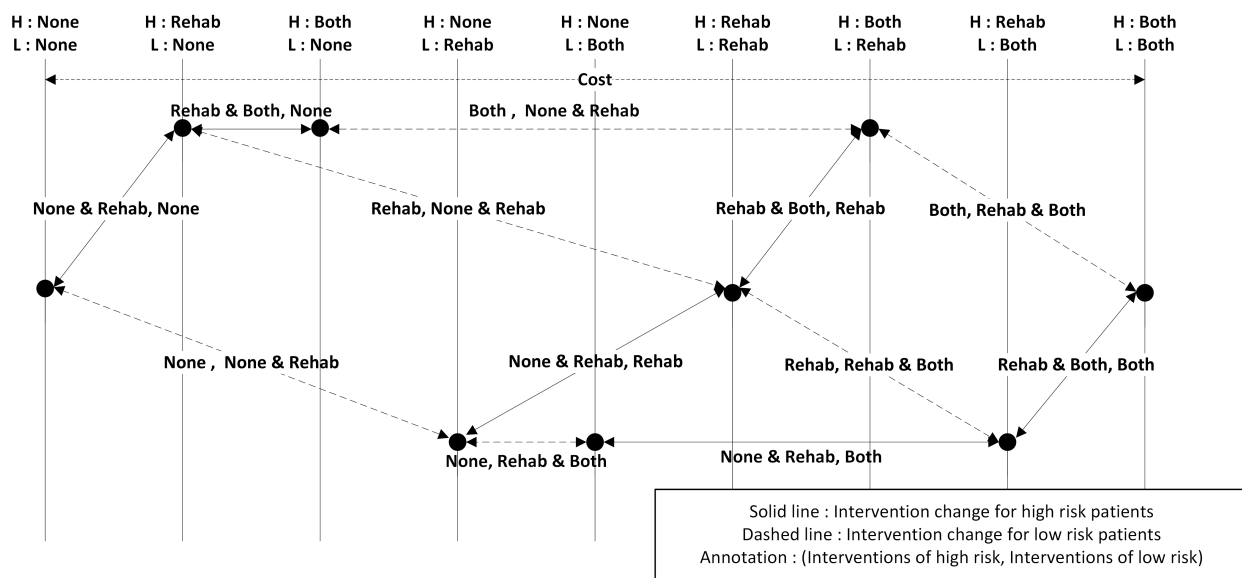


Figure 5.4: The candidates for optimal combinations of interventions, when rehab is the most cost-effective intervention

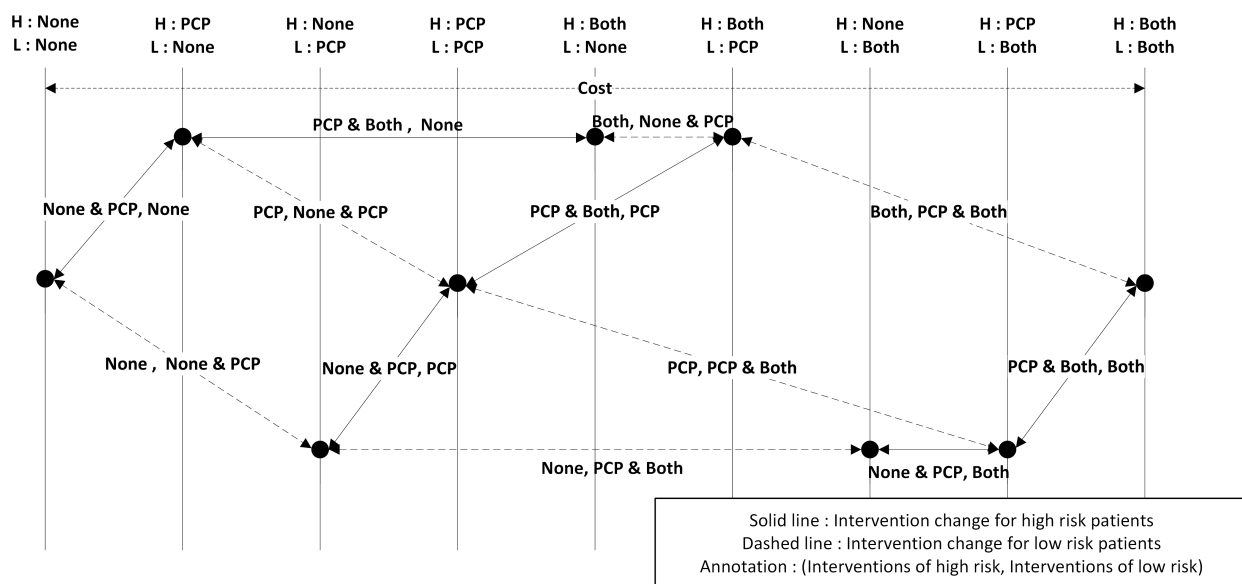


Figure 5.5: The candidates for optimal combinations of interventions, when PCP is the most cost-effective intervention and taking both is next

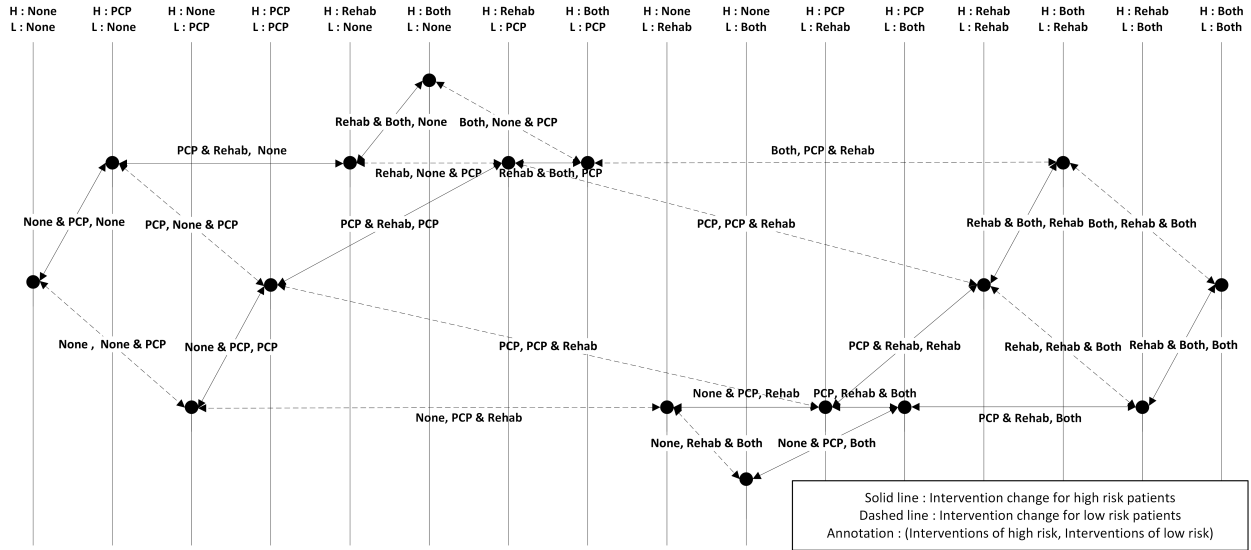


Figure 5.6: The candidates for optimal combinations of interventions, when PCP is the most cost-effective intervention and rehab is next

5.4 Case Study

Collaborating with physicians, nurses, pulmonologists, respiratory therapists, and other staff in SSM Health St. Mary's Hospital, and based on close to 900 patients' records, we identify the factors that impact COPD patients' readmissions.

Specifically, based on the collected data in St. Mary's Hospital's database and the estimates provided by the staff at St. Mary's Hospital (summarized in Table 5.1), we carry out a case study to investigate the impact of incentives on readmission probability for high- and low-risk patients, where the risks are evaluated based on a random forest prediction model with 88% accuracy.

Table 5.1: Case study parameters ($N = 50$, $\rho = 0.3$, $C_d = 5000$)

	Both PCP and Rehab	Rehab	PCP	None
β_i^H	0.20	0.30	0.58	0.60
β_i^L	0.05	0.10	0.20	0.30
C_i	220	200	20	0

In Step 1 of Algorithm 2, we calculate the magnitudes of slopes of all possible edges.

These values are:

$$\begin{aligned} S_{n,c}^H &= 0.001, & S_{n,r}^H &= 0.0015, & S_{n,b}^H &= 0.0018, \\ S_{c,r}^H &= 0.0016, & S_{c,b}^H &= 0.0019, & S_{r,b}^H &= 0.005, \end{aligned}$$

for high-risk group and

$$\begin{aligned} S_{n,c}^L &= 0.005, & S_{n,r}^L &= 0.001, & S_{n,b}^L &= 0.0011, \\ S_{c,r}^L &= 0.0006, & S_{c,b}^L &= 0.008, & S_{r,b}^L &= 0.0025, \end{aligned}$$

for low-risk group.

In Step 2, for the high-risk group, the lower part of the polyhedron consists of one edge of $(C_n, \beta_n^H) - (C_b, \beta_b^H)$ (see Figure 5.7), since the edge between the points for no intervention and both interventions has the largest slope size, among all the edges that consist of the point of no intervention as one of the end points. For the low-risk group, the lower part edges are $(C_n, \beta_n^L) - (C_c, \beta_c^L)$ and $(C_c, \beta_c^L) - (C_b, \beta_b^L)$, since among all the edges connected with the point of no intervention, the edge between no intervention and PCP has the largest magnitude of slope, and the edge between PCP and both interventions has larger slope than that of the edge between PCP and rehab.

In Step 3, the edges are sorted by their slopes in a descending order. There are 3 lower part edges resulting from Step 2 with slope sizes 0.005, 0.0018, 0.0008 in the following order: $(C_n, \beta_n^L) - (C_c, \beta_c^L)$, $(C_n, \beta_n^H) - (C_b, \beta_b^H)$, and $(C_c, \beta_c^L) - (C_b, \beta_b^L)$.

In Step 4, the optimal readmission rate function is constructed as follows:

$$\begin{aligned} f(x) &= I(0 \leq x < 14)(-0.005x + 0.39) \\ &\quad + I(14 \leq x < 80)(-0.0018(x - 14) + 0.32) \\ &\quad + I(80 \leq x < 220)(-0.0008(x - 80) + 0.2012). \end{aligned}$$

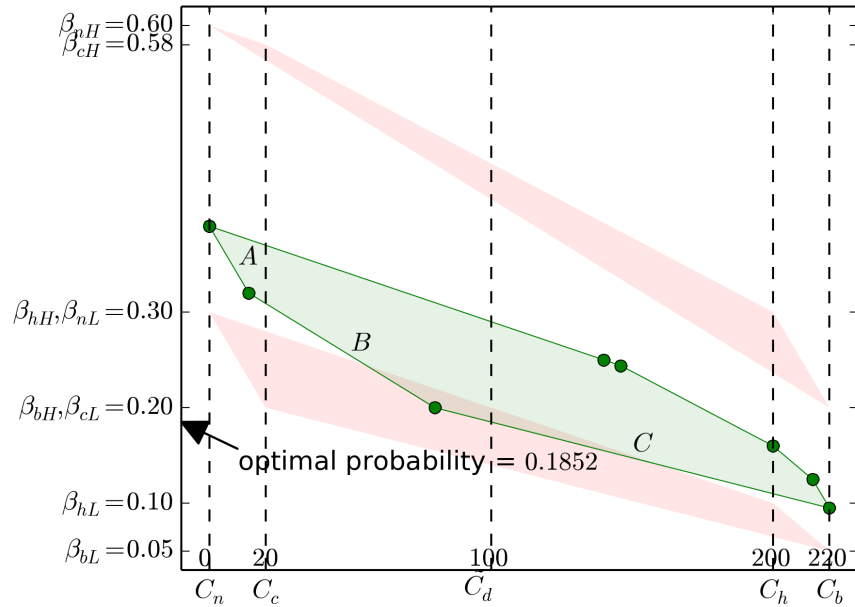


Figure 5.7: Case study illustration

From this function, conditions to achieve an optimal intervention for minimal readmission rate can be generated. The optimal intervention plans are then determined according to which region the incentive budget per patient, x , belongs to. In this case study, from the left-most edge, if the budget value belongs to the x -range of the first edge (“A” in Figure 5.7), the optimal choice is to encourage low-risk patients to visit PCPs as much as possible. If the budget value belongs to the x -range of the second edge (“B” in Figure 5.7), the optimal choice is to encourage all low-risk patients to visit PCPs and some high-risk patients to visit both rehab and PCPs as much as possible. If the budget value belongs to the x -range of the last edge (“C” in Figure 5.7), the optimal choice is to encourage all high-risk patients to get both while some low-risk patients to get both interventions as much as possible and others to visit PCPs.

In this case study, the incentive budget per patient is $x = \frac{B}{N} = \frac{5000}{50} = 100$ and it belongs to the last edge. Thus, with this incentive budget, the optimal readmission rate is $f(100) = -0.0008(100 - 80) + 0.2012 = 0.1852$. Such a result implies that all the high-risk patients ($0.3 \cdot 50 = 15$) and some low-risk patients ($\frac{100-80}{220-80} \cdot 0.7 \cdot 50 = 5$) take both

interventions, and the rest of low-risk patients ($50 - 15 - 5 = 30$) visit PCPs. As shown in Figure 5.7, when the vertical line of $x = 100$, which indicates the given unit budget value, is drawn, the intersection of the line with the third lower edge is corresponding to the edge between the points for PCP and both in the low-risk group $((C_c, \beta_c^L) - (C_b, \beta_b^L))$. Therefore, the optimal solution can be achieved when all high-risk patients and several low-risk patients take both interventions, while the rest visit PCP.

5.5 Discussions

The above study presents a graphical framework to derive optimal solutions for post-discharge interventions for COPD patients. Such a framework can help hospital management to design its incentive budget, understand the best outcome associate with the budget, and encourage patients to take necessary interventions to reduce readmission rate. In addition to derive the optimal solutions, the proposed framework has more advantages. First, it provides an equivalence to the incremental cost-effectiveness ratio (ICER) method, which is a widely used method for evaluating clinical interventions. Secondly, the framework facilitates sensitivity analysis.

5.5.1 Cost-effectiveness

Using this framework, cost-effective analysis can be carried out implicitly without evaluating the ICER because the concept of ICER is replaced by the slope of edge in the framework. Define the effect of interventions as the amount of reduction in COPD readmission rate. Then, the ICER between two interventions, A and B, can be obtained by

$$ICER = \frac{\text{Cost}(A) - \text{Cost}(B)}{\Delta P_r(A) - \Delta P_r(B)}, \quad (5.5)$$

where $\text{Cost}(k)$ and $\Delta P_r(k)$, $k = A, B$, are the cost and the reduction in readmission probability of intervention k , respectively.

With Algorithm 2, the steepest slope of the sides will be obtained. The absolute value of the slope between the two points of interventions A and B is the reciprocal of ICER between interventions A and B.

$$\begin{aligned} \left| \frac{P_r(A) - P_r(B)}{\text{Cost}(A) - \text{Cost}(B)} \right| &= \left| \frac{(P_r(A) - P_r(\text{Base})) - (P_r(B) - P_r(\text{Base}))}{\text{Cost}(A) - \text{Cost}(B)} \right| \\ &= \frac{(P_r(\text{Base}) - P_r(A)) - (P_r(\text{Base}) - P_r(B))}{\text{Cost}(A) - \text{Cost}(B)} \\ &= \frac{\Delta P_r(A) - \Delta P_r(B)}{\text{Cost}(A) - \text{Cost}(B)} = \frac{1}{ICER}, \end{aligned}$$

where $P_r(\text{Base})$ is the baseline readmission probability under current intervention. Thus, the steepest slope value implies the smallest ICER. Therefore, by following the lower sides of the octagon, we greedily select interventions that have smaller ICER values between itself and the current intervention in each group.

To elaborate the relationship between the suggested framework and the ICER analysis, consider the case study in Section 8.6. We start from the baseline which has no budget ($C_d = 0$) and no interventions are taken for both risk groups. When the budget increases, the intervention that has the lowest ICER value is selected. At the baseline, the candidates for intervention include:

- PCP visit for high-risk group (I_c^H);
- rehab visit for high-risk group (I_r^H);
- both rehab and PCP visits for high-risk group (I_b^H);
- PCP visit for low-risk group (I_c^L);
- rehab visit for low-risk group (I_r^L),
- and both PCP and rehab visits for low-risk group (I_b^L).

The ICER values between the baseline and those interventions are 1000, 667, 550, 200, 1000 and 880, respectively (see Figure 5.8). Using the ICER criteria (i.e., the smallest ICER

value), the choice of PCP visit for low-risk group is selected. As the budget increases, more low-risk patients can take this intervention, until the budget per person, $\tilde{C}_d = \frac{C_d}{N}$, reaches the value of 14 (i.e., $(1 - \rho)C_c$), which implies all low-risk patients visit PCPs.

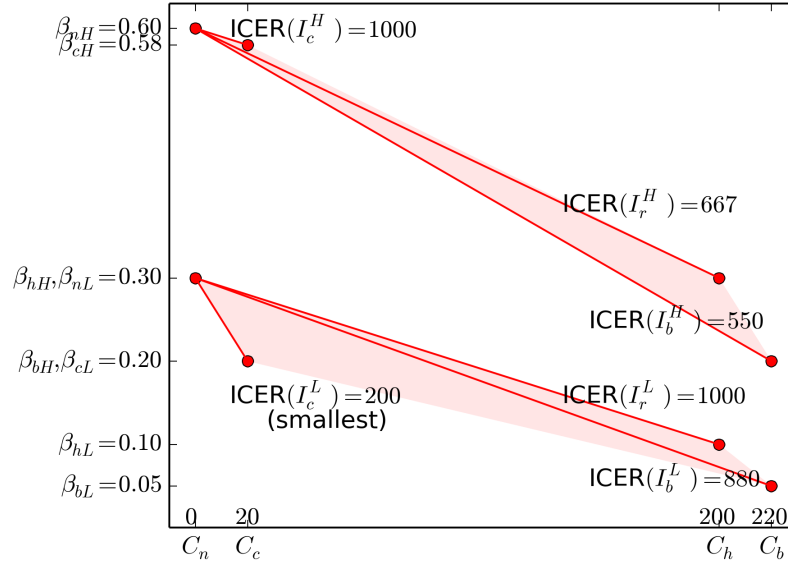


Figure 5.8: ICER illustration

When the budget increases to more than 14, the rest of budget can encourage patients taking another intervention in order to reduce the readmission rate. At this point, the possible intervention candidates are I_c^H , I_r^H , I_b^H , I_r^L and I_b^L (since I_c^L is already used). The ICER values of the interventions for high-risk patients are the same as before: 1000, 667, 550, while the ICER values of the interventions for low-risk patients are changed since the baseline intervention is changed to PCP visit already. Thus, following the same procedure, the ICER values between new interventions and PCP visit are 1800 and 1333, for I_r^L and I_b^L , respectively. Among these five candidates (three “old” two “new”), taking both interventions for high-risk patients has the smallest ICER value, which is then selected as the next intervention. Therefore, the number of high-risk patients who take both interventions increases until all high-risk patients take this option. Now the level of budget per patient is $\tilde{C}_d = 14 + \rho C_b = 80$.

Furthermore, when \tilde{C}_d increases from 80, the next intervention will be selected and there are only two possible candidates left: I_r^L and I_b^L since the high-risk patients already take both interventions. Between these two interventions, taking both interventions for low-risk patients has a smaller ICER value, so that it is selected for low-risk patients.

With all these selections, the optimal readmission rates can be calculated. As one can see, by greedily selecting the intervention that has the smallest ICER value among all the listed possible intervention candidates, we obtain the same optimal solution to achieve the minimal readmission probability. However, such a procedure is relatively more complicated. In addition, the optimization framework has an analytical significance from the point of view that it not only provides an optimal solution and a simple and straightforward way to accomplish the goal, but also can graphically illustrate how the readmission rate is decreased.

5.5.2 Sensitivity Analysis

To study the effect of parameter variation (e.g., cost, readmission probability), first we consider the costs of visits to PCP and rehab center. With the same setting in Table 5.1 except the cost of PCP visit, the optimal readmission rate increases as the cost increases. As shown in Figure 5.9, the resulting optimal readmission rate is monotonically increasing and the growth rate is changed in various ranges. These slope variations come from either a change of optimal policy or an update of the ratio between the difference of readmission rates due to interventions and the difference of the costs. For example, if the cost of PCP visit changes from 20 to 50, then the optimal intervention will change so that the slope of the resulting readmission rate changes (see Figure 5.9).

Figure 5.10 shows how the resulting intervention solution changes. With $C_c = 20$, the optimal solution corresponds to the third lower edge (“C” in Figure 5.10) on lower octagon, which has green vertices in Figure 5.10. However, when the cost changes to 50, the resulting intervention changes to the second lower edge (“B^{new}” in Figure 5.10) of the upper octagon

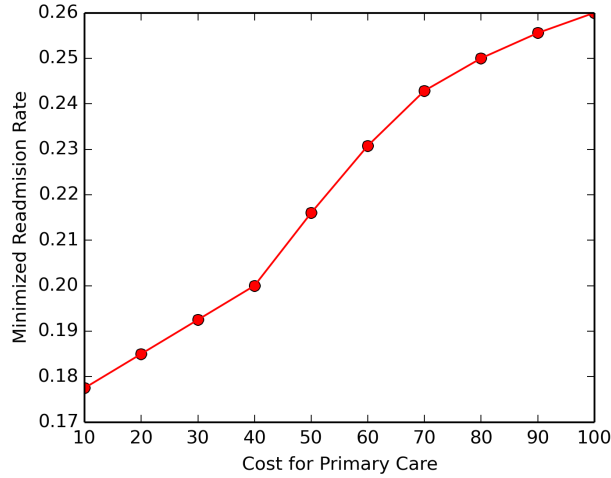


Figure 5.9: The resulting optimal readmission rate with the change of PCP cost

whose vertices are yellow.

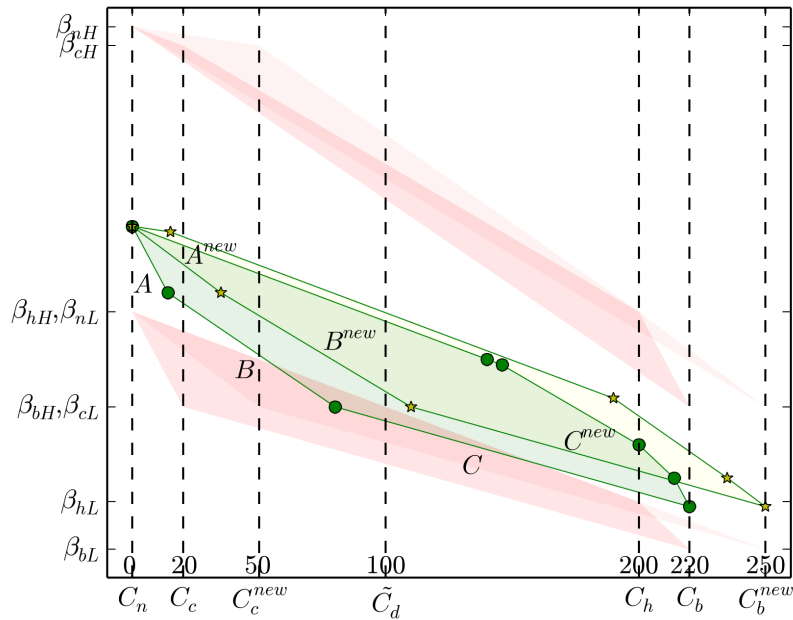


Figure 5.10: An optimal policy transition by the PCP cost shift

Similarly, the optimal readmission rate increases as the cost of rehap visit increases, or decreases if such a cost decreases. As illustrated in Figure 5.11, the resulting optimal readmission rate is monotonically increasing. However, its growth rate changes, which

decreases until the cost is between 220 and 240, and starts to increase if the cost is beyond it.

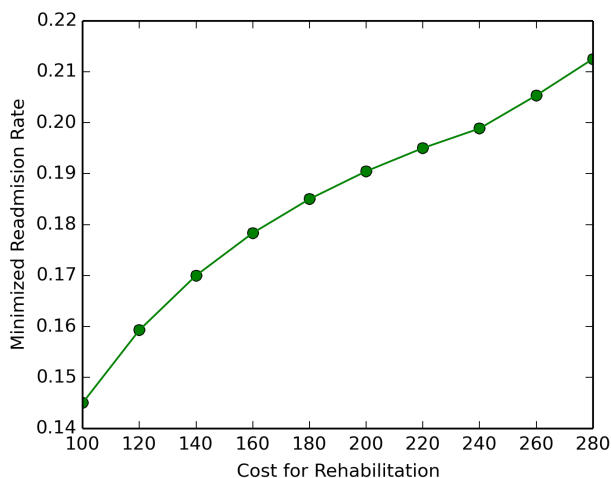


Figure 5.11: The resulting optimal readmission rate with the change of rehab cost

The cause of this growth rate change can be explained using Figure 5.12, by comparing the costs of rehab visit between 200 and 150. The resulting optimal solutions still correspond to the third lower edge (“C” in upper octagon and “C^{new}” in lower octagon in Figure 5.12) even if the cost is changed. However, the slope of the edge is changed to be steeper, so that a decreasing speed of readmission rate is observed.

Figures 5.13 and 5.14 show the result of the readmission probability change of each intervention. In Figure 5.13, the red graph (i.e., the line with solid dots) illustrates that the resulting readmission rate does not change even if the readmission probability of no intervention for high-risk patients decreases, until it reaches a certain point. This is because until this point, the change of readmission probability does not alter the optimal solution at the given budget value. As mentioned earlier, the optimal solution corresponds to the third edge and this fact or the position of the third edge does not change with the decreasing readmission probability of no intervention for high-risk patients.

However, when this readmission probability decreases further beyond a certain point, the magnitude of the second lower edge’s slope decreases so that it is less than that of the

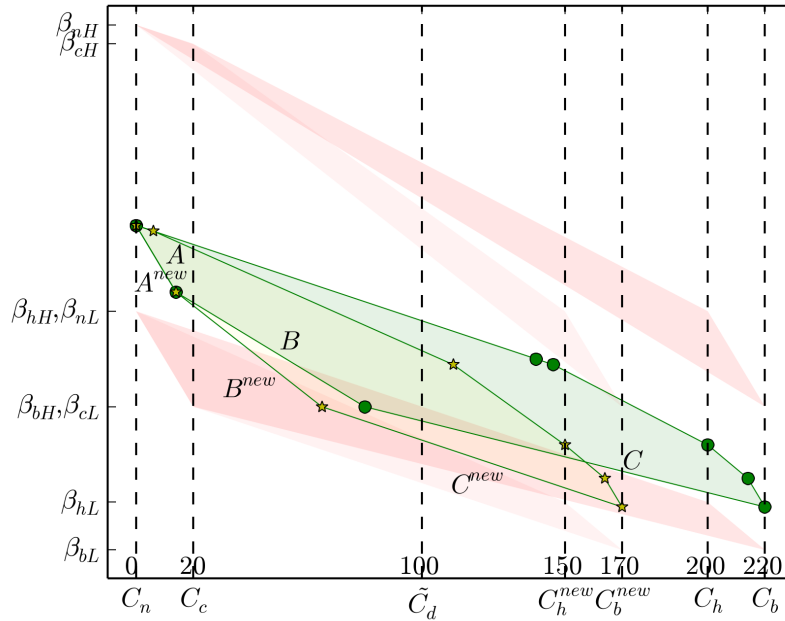


Figure 5.12: An optimal policy transition by the rehab cost shift

third lower edge. Then there is a switch from the third edge to the second, and the optimal solution changes to correspond to the second edge. Its position descends as the readmission probability of no intervention further decreases. The other curves in Figures 5.13 and 5.14 illustrate the similar movements. The resulting optimal readmission rate is maintained at the same level until the decrease of each readmission probability of interventions makes a change in the optimal solution.

5.5.3 Variation Analysis

There are a few implementation issues regarding the incentives. First, in the intervention model, the optimal solution may include patients within the same group taking different incentives. Then, if such a solution is to be implemented in practice, additional efforts need to be devoted to further segmenting patients with different groups, for instance, based on social-economic conditions.

Secondly, even though the incentives can be used to encourage patients to take inter-

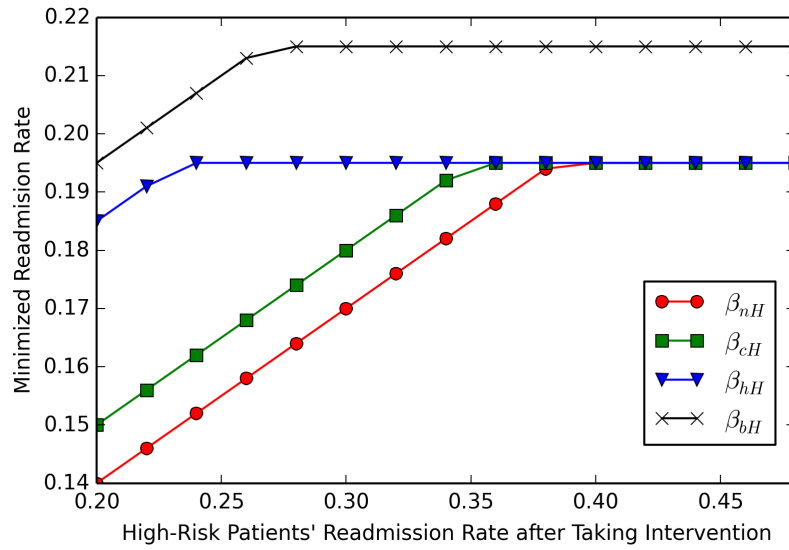


Figure 5.13: The resulting optimal readmission rate changes by the change of each intervention's readmission rate

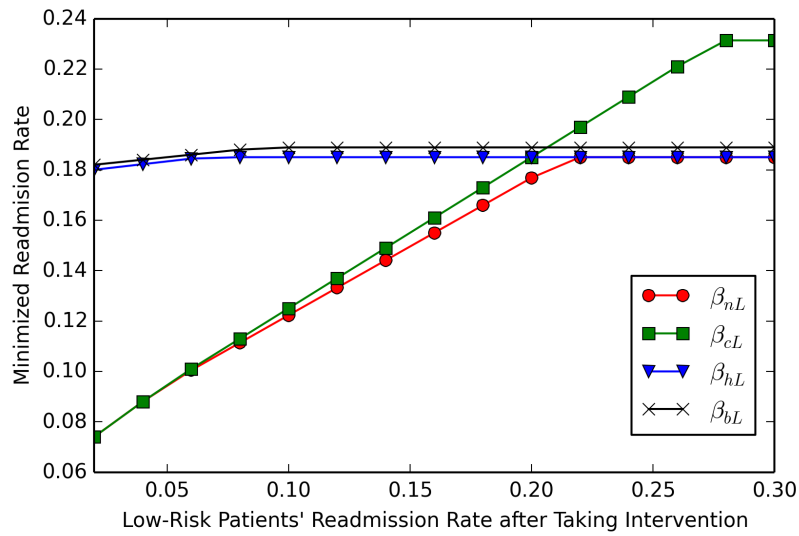


Figure 5.14: The resulting optimal readmission rate changes by the change of each intervention's readmission rate

ventions, it is still possible that some patients may not comply with the instructions. Below we investigate this issue and show that the framework generates robust solutions under certain conditions. Let ϵ be the percentage of patients who do not follow the instruction despite of the incentives, and P_1 and P_2 be the resulting readmission probabilities with intervention Scenarios 1 and 2, respectively. Assume Scenario 1 is the best, thus $P_1 < P_2$. Due to $\epsilon > 0$, the actual readmission probabilities realized for Scenarios 1 and 2 will be P_1^{realized} and P_2^{realized} . We would like to find out under which condition Scenario 1 remains as the best (i.e., $P_1^{\text{realized}} > P_2^{\text{realized}}$).

To find this condition, assume that the base scenario before applying the new interventions is the same for both cases. Denote β_0 , β_1 and β_2 as the readmission probabilities under base intervention, and interventions of Scenarios 1 and 2, respectively. Also define α_1 and α_2 as the desired intervention rates that patients follow Scenarios 1 and 2, respectively. Since Scenario 1 is the optimal one, the reduced readmission probability of Scenario 1 will be greater than or equal to that of Scenario 2, i.e., the following statement is true:

$$\alpha_1(\beta_0 - \beta_1) \geq \alpha_2(\beta_0 - \beta_2). \quad (5.6)$$

Thus, we would like to find out the condition of ϵ that Scenario 1 remains optimal even though the realized readmission probability is less than the desired outcome. In other words,

$$(\alpha_1 - \epsilon)(\beta_0 - \beta_1) \geq (\alpha_2 - \epsilon)(\beta_0 - \beta_2). \quad (5.7)$$

Define

$$f(\epsilon) = (\alpha_2 - \epsilon)/(\alpha_1 - \epsilon).$$

Then equation (5.7) is equivalent to

$$f(\epsilon) \leq (\beta_0 - \beta_2)/(\beta_0 - \beta_1).$$

Consider the derivative of $f(\epsilon)$

$$f'(\epsilon) = \frac{\alpha_2 - \alpha_1}{(\alpha_1 - \epsilon)^2}.$$

In addition, from (5.6), we have

$$f(0) \leq \frac{\beta_0 - \beta_2}{\beta_0 - \beta_1}.$$

Then if at least one of the following two conditions is satisfied, equation (5.7) is satisfied and the optimality of Scenario 1 is retained.

$$\text{Condition 1: } \alpha_2 \geq \alpha_1,$$

$$\text{Condition 2: } x \leq \frac{\alpha_1(\beta_0 - \beta_1) - \alpha_2(\beta_0 - \beta_2)}{\beta_2 - \beta_1}. \quad (5.8)$$

If any of these two conditions is not satisfied due to patients' incompliance, the optimal solution would not be optimal anymore. Thus, additional efforts to ensuring these conditions to be satisfied need to be pursued.

5.6 Conclusions

In this chapter, an intervention framework is proposed to reduce COPD readmissions. Through targeted incentives, COPD patients are encouraged to comply with the intervention plan. Based on each patient's specific readmission risk, the co-pay and transportation costs for patients visiting PCPs and rehab centers can be reimbursed. An optimization model is then introduced to minimize COPD readmission rate under a given incentive budget. Solving the optimization problem, we obtain the minimal readmission rate and the conditions to achieve it. These results can provide a guideline or quantitative tool for hospital management to design incentive budget by comparing the results of optimal readmission rates under different incentive options. Through a case study at St. Mary's Hospital, we illustrate the applicability of the method. Moreover, the sensitivity, variation and cost-effective analyses of the model are also carried out. Such a work provides a

quantitative tool for hospital management to design appropriate intervention policies and plan necessary budget levels to reduce COPD readmissions. Such a methodology can be transformed to study readmissions for patients with other symptoms.

Chapter 6

Dynamic Intervention Decision for Reducing COPD Readmissions

6.1 Introduction

In clinical environment, medical interventions encompass procedures, processes, equipment, and medications used during healthcare encounters (both face-to-face or electronically, i.e., virtual). The intervention decisions are often made sequentially during the course of disease progression and treatment process. New interventions should be designed based on the patient's illness and medical conditions, as well as the effect of prior ones. Thus, the interventions should be tailored to each individual patient by considering his/her health condition and various characteristics in order to achieve better outcome. However, the existing analytical work related to COPD readmission mainly focuses on data analysis to predict the risk of readmission. The limited studies devoted to intervention policies have not addressed the dynamic nature of post discharge care and not tailored to control or manipulate the specific factors for each individual patient. Therefore, there is a need to develop personalized intervention plans and dynamically update intervention decisions. As causal Bayesian network models are capable of finding the causal relationships between

variables, while Markov decision process (MDP) models can be used to find optimal solutions to solve stochastic and dynamic decisions, this chapter integrates the two models and proposes a causal Bayesian network based Markov decision process (CNMDP) model to provide personalized dynamic decision support for intervention planning.

6.2 Causal Network Markov Decision Process

In this section, we first describe the COPD intervention decision problem and discuss the limitations using a standard MDP model. Then, we propose a CNMDP model to overcome the difficulties.

6.2.1 COPD Intervention Decision Problem

Suppose that the information of current state of a COPD patient's readmission risk at every decision point is known. The goal is to find the most effective intervention based on the patient's risk state and dynamically update the optimal intervention in every specified period. Since the maximum length of monitoring period of readmission for COPD patients is 30 days, a 4-week monitoring period is selected and the best intervention is updated at the beginning of each week based on the updated risk status.

Such a decision problem may be formulated as a naive finite-horizon Markov decision process by defining the system state as the readmission risk (e.g., high and low) and the readmission status (readmitted or not), and the actions by the available interventions and a non-intervene action. The problem can be solved when the transition probabilities between the states are known. However, a patient's health status, such as the risk of COPD readmission, is conceptual rather than objective. Thus, the effectiveness of interventions at each state and following transitions between the states are not easy to measure or detect. Therefore, a direct application of MDP model becomes hardly possible and an alternative way needs to be identified.

Although the direct transition probabilities between risk states for a given intervention may not be available, if the impact of the intervention on a specific variable is observable and measurable, then the state transition can be obtained by incorporating such a variable into the system and investigating the relationship between the variable and the targeting health status. Thus, all variables related to the risk of readmission can be engaged into an MDP model and the variables and the associated relationships can be represented in a graphical form. For this purpose, in the graph, the effect of controlling variables via an intervention should be appropriately propagated to the subsequent variables in order to represent the intervention's resulting effect. Since a causal Bayesian network (CBN) can make such a causal inference, which analyzes the response of the effect variable when the cause is changed, a CBN is introduced to the MDP model to analyze the risk changes resulting from the interventions.

6.2.2 Markov Decision Process Using Causal Bayesian Networks

In this study, we extend the MDP model by using a causal Bayesian network, in which each state consists of the current values of a patient's characteristics as well as the indication of readmission. The network determines the relationships between the characteristics and readmission. Then, the problem is formulated by defining the state space and transition probabilities, from which the definitions of actions and reward function are obtained.

Suppose a patient's characteristics affecting the readmission risk are represented by a set of variables $\mathcal{X} = \{x_1, \dots, x_n\}$. Let the status of a patient at time t be denoted as (\mathcal{X}_t, y_t) , where y_t is 1 if the patient is readmitted at time t and 0 otherwise, and \mathcal{X}_t defines the values of variables in \mathcal{X} at time t . In addition, a variable is referred to as *controllable* if it can be manipulated and its value can be changed through an intervention. Assume that \mathcal{X} contains at least one controllable variable.

- Time Horizon: $T = \{1, 2, 3, 4\}$. As we assume intervention decisions are updated weekly, there are four decision epochs after discharge, where $t = 1$ represents the

time of discharge, and each t of the remaining epochs indicates the beginning of t -th week from discharge. For a general scenario, if the period of intervention is assumed as n days, then the set of decision epochs can be generalized as $T = \{1, 2, \dots, \lfloor \frac{m}{n} \rfloor\}$, where m is the number of days in total monitoring period. For instance, $m = 30$ and $n = 7$ in the COPD readmission problem.

- **States:** At each decision epoch, a patient's health state is represented by the vector of health characteristic values \mathcal{X} , and the readmission state is characterized by the value of y . Let S_x be the set of values that $x \in \mathcal{X}$ can have, and $S_{\mathcal{X}}$ be the Cartesian product of all $S_x, x \in \mathcal{X}$. Then, the state space S can be defined as follows,

$$S = \bigcup_{\forall \mathcal{X} \in S_{\mathcal{X}}} \bigcup_{\forall y \in \{0,1\}} \{(\mathcal{X}, y)\}, \quad (6.1)$$

where the number of states $|S| = 2 \cdot |S_{\mathcal{X}}|$ as $y \in \{0, 1\}$.

- **Action taken at status s :** $a(s) \in A = \{a_0, a_1, \dots, a_n\}$, where a_0 indicates no intervention, and a_1 to a_n represent n types of interventions. Here we only consider effective interventions so that each intervention will affect at least one of the variables in \mathcal{X} . For example, if there exist a variable x in \mathcal{X} , representing the smoking status of a patient, which is considered as a controllable variable, then action set A contains an intervention that can alter the smoking status (from currently smoking to non smoking).
- **Transition Probability:** $p_t(s' | s, a)$ defines the transition probability from state s to state s' by action a at time t , and can be derived from the network structure and corresponding probability tables. The detailed explanations are provided in the next subsection.
- **Rewards:** When action a is taken at state s at time t and a transition from $s = (\mathcal{X}, y)$

to $s' = (\mathcal{X}', y')$ occurs, the reward $r_t(s, s', a)$ is defined as follows:

$$r_t((\mathcal{X}, y), (\mathcal{X}', y'), a) = \begin{cases} -c_r^t - c_a^t & \text{if } y \neq 1, y' = 1, \\ -c_a^t & \text{otherwise,} \end{cases} \quad (6.2)$$

where c_a^t is the intervention cost of action a and c_r^t is the readmission cost occurring due to the readmitted patient at time t .

- The goal is to choose a policy π^* that will maximize the cumulative function of random rewards, i.e., the expected discounted sum of rewards over a potential horizon.

Remark 6.1. *Although the above formulation focuses on the readmission issue for COPD patients, the problem can be generalized to abate their health deterioration for longer follow-up period by extending the time horizon of decisions and adjusting the reward and state definition by substituting the readmission-related elements with other components.*

Under this formulation, the next step is to estimate the transition probabilities between states following the interventions. In order to obtain such probabilities, only the correlation information is not enough. The impact of variable manipulation by an intervention on the subsequent variables should be reflected in the network, i.e., the relationships needs to be causal. Thus, a causal Bayesian network is utilized to represent the relationships between variables and calculate the probabilities of state transitions. Below the structure of the network, transition probabilities, assumptions, and conditional probabilities are introduced.

6.2.2.1 Dynamic Causal Bayesian Network

Assume that there exist a set of variables $\mathcal{X} = \{x_1, \dots, x_n\}$ affecting the future readmission risk represented by y . When the relationships between variables are stationary, they can be represented by a graphical structure through learning a causal Bayesian network.

A causal Bayesian network has the same form of a Bayesian network, but a parent node of a variable x in a causal network should be x 's direct cause. In other words, an arc $x \rightarrow x'$ exists in a causal network if and only if x is a direct cause of x' in the specified population described by the given dataset. A Bayesian network can be represented by a graphical structure $G = (V, A)$ where $V = \{x_1, \dots, x_n\}$ and A are the sets of nodes and arcs, respectively, and the joint probability distribution in G is represented by the following form.

$$p_G(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p_G(x_i | \text{Parent}_G(x_i)). \quad (6.3)$$

If a causal Bayesian network is known, the causal inference can be drawn and the effect of manipulation on a variable x can be evaluated. When the distribution of x is replaced with $p'_G(x | z)$ where $z \neq x$, we just need to replace the term $p_G(x | \text{Parent}_G(x))$ in (6.3) with a newly specified density $p'_G(x | z)$. Then, we obtain the following rule,

$$p'_G(x_1, x_2, \dots, x_n) = p'_G(x | z) \prod_{x' \in V \setminus \{x\}} p_G(x' | \text{Parent}_G(x')). \quad (6.4)$$

which implies that, if the causal network structure is known and an unmanipulated density $p_G(v_1, v_2, \dots, v_n)$ is estimated based on the observed dataset, the prediction of the effect of an intervention manipulating values of variables can be made even if the manipulation is not observed in the given data.

In addition, as the notion of time T is involved in the MDP model, temporal evolution of variables should be considered in the model. A dynamic BN is an extension of BNs that represents the temporal transition of variables over time, where nodes and arcs in the graph represent random variables and probabilistic dependencies between variables across time and those at the same time, respectively. As intervention decisions are made in discrete time, we construct a dynamic BN using a causal BN over a discrete number of time steps, which is the same with the decision epochs. Both the probabilistic dependencies between variables and the probability distributions describing the temporal dependencies

are assumed to be time invariant, so that the dynamic BN is composed of the same causal BN at each epoch and the same temporal transition probabilities that can be utilized between the epochs.

Now consider a causal Bayesian network model of COPD readmission represented by the graphical structure depicted as the first graph in Fig. 6.1. As one can see, there are four epochs in the time horizon. The states of a patient are characterized by a set of health status variables $\mathcal{X} = \{x_a, \dots, x_f\}$ and readmission status y . The links between variables in \mathcal{X} in each network represent the causal relationships, while the links across the epochs illustrate both the causalities and the transitions, characterized by transition probabilities. The arrows pointing toward readmission states indicate that the causal Bayesian network model predicts the readmission risks. As the causal Bayesian network prediction is updated in every epoch, it captures the dynamic nature of post discharge care and intervention.

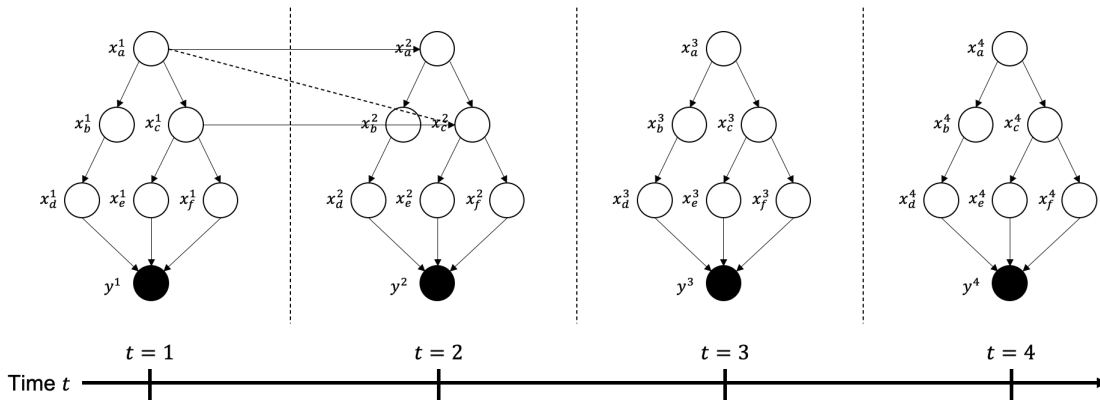


Figure 6.1: Structure of CNMDP

In COPD and many other post-discharge care processes, the state of a variable at time $t + 1$ largely depends on the state of the variable at time t . Fig. 6.1 provides an illustration of a dynamic causal Bayesian network. For example, x_c^{t+1} largely depends on x_c^t due to the stronger persistence of variable x_c , but also is affected by x_a^{t+1} because of the causal relationship depicted in the Bayesian network. If there is no intervention on x_c at time t and x_a^{t+1} is not different from its previous state x_a^t , then x_c^{t+1} will be stable as x_c^t (except in the case of its natural degradation). However, if x_a^{t+1} is different from x_a^t , x_c^{t+1} will be affected

by the manipulation in its parent node because of the causal relationship between x_a and x_c . Therefore, in this case, the state of x_c^{t+1} depends on its previous state x_c^t , its parent node x_a^{t+1} and parent's past state x_a^t as depicted in Fig. 6.1. For other nodes, the dependency rule is the same, but for simplicity, these relationships are not illustrated in Fig. 6.1.

Remark 6.2. Fig. 6.1 implies that the node corresponding to y , which indicates readmission, does not have an out-going arc and always locates at the last layer in the effective causal Bayesian network. If there exist descendant nodes of y or nodes that are not connected to y in a causal Bayesian network G , such nodes can be removed from G in the analysis because a manipulation on these nodes does not affect y .

6.2.2.2 Transition Probability

Let $s \in S$ and $s = (s_{\mathcal{X}}, s_y)$. When $p(s' | s, a)$ is the probability of transition $s \rightarrow s'$ given action a , the transition probability in the proposed model can be decomposed by the conditional probabilities in the dynamic causal Bayesian network. Let y be a binary variable that has an absorbing state, for example $y = 1$ represents the readmitted status. As shown in Figure 6.1, when $y_t \neq 1$, the transition probability from a state (\mathcal{X}_t, y_t) to a state $(\mathcal{X}_{t+1}, y_{t+1})$ can be obtained by the following equation.

$$p((\mathcal{X}_{t+1}, y_{t+1}) | (\mathcal{X}_t, y_t), a) = p(y^{t+1}, x_a^{t+1}, \dots, x_f^{t+1} | x_a^t, x_b^t, \dots, x_f^t, a) \quad (6.5)$$

$$= p(y^{t+1} | \text{Parent}_G(y)^{t+1}) \prod_{x \in \mathcal{X}} p(x^{t+1} | x^t, \text{Parent}_G(x)^{t+1}, \text{Parent}_G(x)^t, a), \quad (6.6)$$

where $\text{Parent}_G(x)^t$ represents the parent nodes of x in the network G at time t . The first equality in (6.5) comes from the fact that patient characteristics \mathcal{X}_t consist of variables $x_a^t, x_b^t, \dots, x_f^t$ in the graph, while the second equality (6.6) specifies the conditional probabilities in the network. Thus, for the system in Figure 6.1, it can be further expanded

as:

$$\begin{aligned}
& p((\mathcal{X}_{t+1}, y_{t+1}) \mid (\mathcal{X}_t, y_t), a) \\
&= p(y_{t+1} \mid x_f^{t+1}, x_e^{t+1}, x_d^{t+1}) p(x_f^{t+1} \mid x_f^t, x_c^{t+1}, x_c^t, a) p(x_e^{t+1} \mid x_e^t, x_c^{t+1}, x_c^t, a) \\
&\quad \cdot p(x_d^{t+1} \mid x_d^t, x_b^{t+1}, x_b^t, a) p(x_c^{t+1} \mid x_c^t, x_a^{t+1}, x_a^t, a) p(x_b^{t+1} \mid x_b^t, x_a^{t+1}, x_a^t, a) p(x_a^{t+1} \mid x_a^t, a).
\end{aligned} \tag{6.7}$$

6.2.2.3 Assumptions

With the above description, to facilitate estimating the parameters of transition probabilities, the following assumptions are introduced.

Assumption 6.1. (A Natural Order of States) *All variables in \mathcal{X} are categorically represented using ordered numbers $0, 1, 2, \dots, n$, where there exists a nature order of states such that a higher number indicates a better status.*

First, we assume that the states of each variable $x \in \mathcal{X}$ have a natural ordering with respect to hospital readmission. For example, it is commonly acknowledged that smoking contributes to COPD exacerbation. If there exist three states, currently-smoking, smoking-in-the-past, and never-smoking, for the variables corresponding to smoking, then these states are naturally ordered associated with readmission, i.e., having values 2, 1, and 0, respectively. Using this assumption, we can also obtain a partial order on S , which can be used for structural properties and solution algorithms.

Lemma 6.1. *State space S has a partial order, denoted as \preceq , such that $s \preceq s'$ if and only if $s_x \leq s'_x$ for all $x \in \mathcal{X}$ and $s_y \leq s'_y$, where $s = (s_{\mathcal{X}}, s_y)$, $s' = (s'_{\mathcal{X}}, s'_y)$.*

Proof: See Appendix A. ■

Secondly, we assume there is no effect on health status if interventions are not applied.

Assumption 6.2. (Control Effect Only) *If no intervention affecting x or its ancestors is taken at epoch t , the value of x does not change at the next epoch $t + 1$.*

This assumption infers that if no intervention is taken, then all $x \in \mathcal{X}$ representing the current health status will not vary and the readmission probability based on the values of $x \in \mathcal{X}$ remains the same. It implies that there is no natural deterioration or improvement of health condition without an intervention. Such an assumption is reasonable since the period between decision epochs is too short to make a natural change.

Thirdly, it is assumed that an intervention will not harm the patient status.

Assumption 6.3. (No Worsening) *A state of x never becomes worse due to an intervention.*

Note that after discharge, the risk of readmission in the next epoch $P(y^{t+1} | \mathcal{X}^{t+1}) = p \neq 0$, and interventions will be assigned to the patient in order to reduce readmission probability from p to $p' < p$. Such efforts will be continued in the subsequent epochs. Thus, in a similar sense to Assumption 2, we assume a state does not become worse and an intervention can only improve or maintain a status. In the extreme case that the risk of readmission is zero, the patient won't need any intervention due to this assumption. Note that in this context, a patient's non-compliance to physician's instruction (such as resuming smoking) is not considered. Such issues can be studied in future work.

Under this assumption and Lemma 6.1, we obtain

Lemma 6.2. *Given an initial state s^1 at the first decision epoch, a state transition can occur only to a subset of space S , $S' = \{s \mid s \succeq s^1, s \in S\}$, throughout the entire decision process.*

Proof: See Appendix A. ■

Lemma 6.2 restricts the state space of a problem instance when the first health condition is given, which enables more efficient computation. In addition, from the proof of Lemma 6.2 (see Appendix A), the following result arises.

Corollary 6.3. *A transition from s to $s' \preceq s$ never occurs in the entire MDP process.*

6.2.2.4 Conditional Probability Tables

Finally, based on the above assumptions, the conditional density for all variables can be obtained by defining the conditional probabilities.

Lemma 6.4. *If intervention a does not affect variable $x \in \mathcal{X}$, then the conditional probabilities regarding x are independent of a , and*

$$p(x^{t+1} = s'_x \mid x^t = s_x, \text{Parent}_G(x)^{t+1} = \text{Parent}_G(x)^t) = \begin{cases} 1 & \text{if } s'_x = s_x \\ 0 & \text{otherwise} \end{cases} \quad (6.8)$$

$$p(x^{t+1} = s'_x \mid x^t = s_x, \text{Parent}_G(x)^{t+1} \neq \text{Parent}_G(x)^t) = \begin{cases} \frac{p_G(x^{t+1}=s'_x \mid \text{Parent}_G(x)^{t+1})}{\sum_{s'_x \geq s_x} p_G(x^{t+1}=s'_x \mid \text{Parent}_G(x)^{t+1})} & \text{if } s'_x \geq s_x \\ 0 & \text{otherwise} \end{cases} \quad (6.9)$$

Proof: See Appendix A. ■

Lemma 6.5. *When an intervention a affects variable $x \in \mathcal{X}$, the conditional probabilities regarding x and a are independent of the parent nodes of x , and*

$$p(x^{t+1} = s'_x \mid x^t = s_x, a_x) = \begin{cases} p_{a_x, (s'_x, s_x)} & s'_x \geq s_x \\ 0 & \text{if otherwise} \end{cases} \quad (6.10)$$

where $p_{a_x, (s'_x, s_x)}$ is the transition probability of the states involving x , and represents the efficacy of intervention a_x .

Proof: See Appendix A. ■

The above two lemmas substantially reduce the number of parameters composing the transition probability tables (See Appendix A) in a given CNMDP network structure to the number of parameters regarding intervention effects.

6.3 Solving CNMDP Problems

Solving the proposed CNMDP problem, we obtain a policy that represents the best action for each state in the MDP at each time, referred to as the optimal policy $\pi = (\pi_1, \dots, \pi_T)$. For most MDPs, learning algorithms are used to derive the optimal policies through learning value functions. A value function $v : s \rightarrow \mathbb{R}$ describes the expected objective value obtained following policy π at current state $s \in S$.

6.3.1 Value Function

By the definition of value functions, the optimal policy and value function at current state s can be represented in recursive forms through the following equations:

$$\pi_t(s) = \operatorname{argmax}_{a \in A} \left\{ \sum_{s' \in S} p_t(s' | s, a) (r_t(s, s', a) + \gamma v_{t+1}(s')) \right\}, \quad (6.11)$$

$$v_t(s) = \sum_{s' \in S} p_t(s' | s, \pi(s)) (r_t(s, s', \pi(s)) + \gamma v_{t+1}(s')). \quad (6.12)$$

For the CNMDP model in this study, to simplify notations, let S_0 denote the set of states $\{(\mathcal{X}, y) \mid (\mathcal{X}, y) \in S, y = 0\}$, and $S_1 = \{(\mathcal{X}, y) \mid (\mathcal{X}, y) \in S, y = 1\}$. If a patient is readmitted, the value function generates zero, thereby no intervention to prevent a readmission is needed and state transitions are not considered because further rewards will be zero, i.e., $v(s) = 0$ for all $s \in S_1$. If a patient has not been readmitted at current time t , the value function can be represented as follows: For all $s \in S_0$,

$$v_t(s) = \sum_{s' \in S} p_t(s' | s, \pi) [r_t(s, s', \pi) + \gamma v_{t+1}(s')] \quad (6.13)$$

$$= \sum_{s' \in S_0} p_t(s' | s, \pi) \cdot [r_t(s, s', \pi) + \gamma v_{t+1}(s')] + \sum_{s' \in S_1} p_t(s' | s, \pi) \cdot [r_t(s, s', \pi) + \gamma v_{t+1}(s')] \quad (6.14)$$

$$= -c_\pi - c_r \cdot \sum_{s' \in S_1} p_t(s' | s, \pi) + \gamma \cdot \sum_{s' \in S_0} p_t(s' | s, \pi) \cdot v_{t+1}(s'), \quad (6.15)$$

where c_π and c_r are intervention cost of π and readmission cost, respectively. Equality (6.13) comes from the definition of value function (6.12). In (6.14), a set of next states $s' \in S$, is partitioned into S_0 and S_1 depending on y in the next state. For transitions into both types of states along the optimal intervention policy π , the negative intervention cost $-c_\pi$ composes the rewards, while the negative readmission cost $-c_r$ is added into the rewards when transitioning into the readmission state, and finally resulting in (6.15).

Let $P(x)$ be a directed path from node x to target variable y in a given network G , and $L(x)$ be the set of variables in any possible paths from x to y in G , so that $L(x) = \cup P(x)$. Then, the following result is obtained:

Lemma 6.6. *If intervention a_x affecting x is taken, transition probabilities for a given action a_x do not involve other variables $x' \in \mathcal{X} \setminus L(x)$. Then, if $s'_{x'} = s_{x'}$, $\forall x' \in \mathcal{X} \setminus L(x)$, we obtain*

$$\begin{aligned} p_t(s' | s, a_x) &= \prod_{z \in L(x)} p(s'_z | s_z, s'_{\text{Parent}_G(z)}, s_{\text{Parent}_G(z)}, a_x) \\ &= p_t(s'_y | s'_{\text{Parent}_G(y)}) \prod_{z \in L(x) \setminus \{x, y\}} p(s'_z | s_z, s'_{\text{Parent}_G(z)}, s_{\text{Parent}_G(z)}) \cdot p(s'_x | s_x, a_x). \end{aligned} \tag{6.16}$$

Otherwise, $p_t(s' | s, a_x) = 0$.

Proof: See Appendix A. ■

6.3.2 Solution Algorithm

In the COPD readmission problem, there exists a specific length of monitoring period to consider, i.e., preventing a patient's readmission within 30 days. Thus, when the time interval between decision points is set as one week, there will be four periods and intervention decisions occur at $t \in \{1, 2, 3, 4\}$. Therefore, we obtain a finite-horizon Markov decision process. In this case, the problem can be solved by Backward induction using the Algorithm 3 below.

Algorithm 3: Backward Induction

Require: A causal Bayesian network on patient characteristic factors \mathcal{X} with corresponding conditional probability tables, a patient's initial state $s_{\mathcal{X}}^1$, cost values of readmission and interventions c_r and c_a for all $a \in A$, a discount factor γ and time horizon T ;

Ensure: The optimal policy π that maximize value function and the optimal value v ;

Set $v_{T+1}(s) = 0$ for all $s \in S$;

for $t = T, T - 1, \dots, 1$ **do**

for $s_{\mathcal{X}} \in S_{\mathcal{X}}$ **do**

for $a \in A$ **do**

$p_r(s_{\mathcal{X}}, a) = \sum_{s'_{\mathcal{X}} \in S_{\mathcal{X}}} p((s'_{\mathcal{X}}, 1) \mid (s_{\mathcal{X}}, 0), a)$;

$pv_n(s_{\mathcal{X}}, a) = \sum_{s'_{\mathcal{X}} \in S_{\mathcal{X}}} p((s'_{\mathcal{X}}, 0) \mid (s_{\mathcal{X}}, 0), a) \cdot v_{t+1}((s'_{\mathcal{X}}, 0))$;

end

$\pi_t((s_{\mathcal{X}}, 0)) = \operatorname{argmax}_{a \in A} \{-c_a - c_r \cdot p_r(s_{\mathcal{X}}, a) + \gamma \cdot pv_n(s_{\mathcal{X}}, a)\}$;

$v_t((s_{\mathcal{X}}, 0)) = \max_{a \in A} \{-c_a - c_r \cdot p_r(s_{\mathcal{X}}, a) + \gamma \cdot pv_n(s_{\mathcal{X}}, a)\}$;

end

end

However, since the state space S consists of all combination of patients' characteristics \mathcal{X} , a computation issue arises due to extremely large state space to calculate all the transition probabilities and functions values at each time t . This leads to a huge computation intensity, particularly at lines 5 and 6 in Algorithm 3, even though the model only considers a relatively short time horizon.

To solve this issue, Lemmas 6.2, 6.4 and 6.5 are applied by utilizing the concept of dynamic programming. First, Lemma 6.2 reduces the size of state space from $|S|$ to $|S'|$, where $S' = \{s \mid s \succeq (s_{\mathcal{X}}^1, 0), s \in S\} \subseteq S$, and the size of the state space of characteristic factors from $|S_{\mathcal{X}}|$ to $|S'_{\mathcal{X}}|$, where $S'_{\mathcal{X}} = \{s_{\mathcal{X}} \mid s_{\mathcal{X}} \succeq s_{\mathcal{X}}^1, s_{\mathcal{X}} \in S_{\mathcal{X}}\} \subseteq S_{\mathcal{X}}$. Next, if an intervention a affects x (i.e., $a = a_x$), only the states of x and its descendants can vary. Thus, from Lemmas 6.4 and 6.5, the conditional probabilities for other factors are identities. Then, Lines 5 and 6

with the associated transition probabilities in (6.6) can be refined as:

$$\begin{aligned} p_r(s_{\mathcal{X}}, a) &= \sum_{s'_{\mathcal{X}(a)} \in S'_{\mathcal{X}(a)}} p(((s'_{\mathcal{X}(a)}, s_{\mathcal{X} \setminus \mathcal{X}(a)}), 1) | ((s_{\mathcal{X}(a)}, s_{\mathcal{X} \setminus \mathcal{X}(a)}), 0), a) \\ &= \sum_{s'_{\mathcal{X}(a)} \in S'_{\mathcal{X}(a)}} p(1 | s'_{\text{Parent}_G(y)}) \prod_{x \in \mathcal{X}(a)} p(s'_x | s_x, s'_{\text{Parent}_G(x)}, s_{\text{Parent}_G(x)}, a), \end{aligned} \quad (6.17)$$

$$\begin{aligned} p v_n(s_{\mathcal{X}}, a) &= \sum_{s'_{\mathcal{X}(a)} \in S'_{\mathcal{X}(a)}} p(((s'_{\mathcal{X}(a)}, s_{\mathcal{X} \setminus \mathcal{X}(a)}), 0) | ((s_{\mathcal{X}(a)}, s_{\mathcal{X} \setminus \mathcal{X}(a)}), 0), a) \cdot v_{t+1}(((s'_{\mathcal{X}(a)}, s_{\mathcal{X} \setminus \mathcal{X}(a)}), 0)) \\ &= \sum_{s'_{\mathcal{X}(a)} \in S'_{\mathcal{X}(a)}} [p(0 | s'_{\text{Parent}_G(y)}) \prod_{x \in \mathcal{X}(a)} p(s'_x | s_x, s'_{\text{Parent}_G(x)}, s, a)] \cdot v_{t+1}(((s'_{\mathcal{X}(a)}, s_{\mathcal{X} \setminus \mathcal{X}(a)}), 0)), \end{aligned} \quad (6.18)$$

where $\mathcal{X}(a) \subset \mathcal{X}$ consists of the characteristic factor that a affects (i.e., factor x if $a = a_x$) and its descendant factors, and $S'_{\mathcal{X}(a)}$ is the restricted space space of factors in $\mathcal{X}(a)$. Furthermore, as the transition probabilities between the states are stationary (time-invariant), the results of transition probabilities $p(s' | s, a)$ learned by (6.6), rather than $v_t(s)$, are stored to reduce computation time.

6.4 Structural Properties

With the MDP model introduced above, the following properties can be obtained.

Lemma 6.7. *For all initial state $s \in S_0$, any intervention $a \in A$ decreases the expected readmission probability, i.e.,*

$$\sum_{s'_{\mathcal{X}} \in S_{\mathcal{X}}} p_G(y = 1 | s'_{\mathcal{X}}) \cdot p((s'_{\mathcal{X}}, 1) | (s_{\mathcal{X}}, 0), a) \leq p_G(y = 1 | s_{\mathcal{X}}). \quad (6.19)$$

Proof: See Appendix A. ■

In order to facilitate proofs of the following propositions, we assume that variables $x \in \mathcal{X}$ are binary and efficacy of intervention a_x for all $x \in \mathcal{X}$ are represented by a single

parameter $p_{a_x} \in \mathbb{R}$. However, the binary assumption can be released and the followings still hold in categorical variables as well with appropriate settings in the propositions.

Assumption 6.4. (Binary Variable) *All variables $x \in \mathcal{X}$ are binary and have values 0 or 1, and there is a nature order of states where 1 indicates a better status. The efficacy of intervention a_x is $p_{a_x} \in \mathbb{R}$, where $p_{a_x,(v',v)} = p_{a_x}$ when $v' > v$.*

Proposition 6.1. *Let $P(x)$ be a path from $x \in \mathcal{X}$ to y in the network G . Suppose that there is only one path from x to y for a $x \in vX$. If all interventions corresponding to $z \in P(x) \setminus \{y\}$ have the same efficacy p_a and cost c_a , then the intervention associated with the closest z to y is the most preferred decision for all t .*

Proof: See Appendix A. ■

Proposition 6.1 shows that an intervention on a variable closer to y is preferred when intervention efficacy is the same and there is only one path connecting a variable and y . However, not in all cases interventions in the last layer is the optimal. For example, when the network is given as $G = (V, A)$ where $V = \{x_1, x_2, x_3, x_4, y\}$, $E = \{(x_1, x_2), (x_2, x_3), (x_2, x_4), (x_3, y), (x_4, y)\}$, an intervention on x_2 can be more beneficial than on x_3 or x_4 . Especially in the case that $P_G(y = 0 \mid x_3 = x_4 = 1)$ is much larger than $P_G(y = 0 \mid x_3 = 0, x_4 = 1)$ or $P_G(y = 0 \mid x_3 = 1, x_4 = 0)$ and $P_G(x_3 = 1 \mid x_2 = 1)$ and $P_G(x_4 = 1 \mid x_2 = 1)$ are close to 1, intervention on x_2 may result in larger increment on $p(y = 0)$ than other interventions. However, the following observation shows there exists more general preference of interventions on descendants.

Observation 1. *For any x , suppose that there exists z such that $z \in \cap P(x) \setminus \{x, y\}$. When there are two different interventions affecting on x and z , respectively, and if their efficacy p_a and costs c_a are the same and both of them are applicable at t , then the intervention on z is preferred to the intervention on the parent node for any t .*

In addition, the proof of Proposition 6.1, specifically the comparison between $u_{t,a_{z_{n-1}}}(s)$ and $u_{t,a_{z_{n-2}}}(s)$, implies that a small perturbation on intervention efficacy p_a associated with

an ancestral variable won't change the result given that $P_{G'}(z_{n-1} = 0 \mid z_{n-2} = 1)$ is big enough. It proves the following corollary.

Corollary 6.8. *The intervention with the highest efficacy is not always the optimal decision.*

The next proposition provides an insight of applicable intervention cost.

Proposition 6.2. *There exists a threshold of the ratio between the intervention cost and the readmission cost where if the ratio is larger than the threshold, the intervention is ignored.*

Proof: See Appendix A. ■

Proposition 6.2 and its proof show that in order for an intervention to be useful and comparable, the intervention cost should be less than a certain value. By using the threshold, the health care providers will be able to cut out some of interventions without solving the MDP problem. Moreover, the threshold will be beneficial when designing affordable and effective interventions by providing the range of proper intervention cost depending of its effect.

Proposition 6.3. *For every pair of interventions, if intervention costs are the same, there exist a certain threshold $R_{s,t}$ of the ratio between intervention efficacy so that an intervention dominates another intervention if the ratio is larger/smaller than the threshold at state s and time t .*

Proof: See Appendix A. ■

Proposition 6.3 shows that there is a linear boundary between areas of p_{a_1} and p_{a_2} where one intervention dominates another. It is useful in the sensitivity analysis of intervention efficacy. It shows that if p_{a_1} is larger then $R_{s,t}p_{a_2}, p_{a_2}$ won't be chosen no matter what efficacy values other available interventions have.

6.5 Case Study

6.5.1 Study Description

To demonstrate the applicability of the proposed model, a case study of designing an optimal intervention plan for COPD patients at SSM Health St. Mary's Hospital (SMH) is introduced. In the case study, electronic health records of COPD patients are collected and a causal Bayesian network targeting COPD readmission prediction is developed. Then, based on the learned causal network, a CNMDP problem is formulated and an intervention plan depending on the patient's status is provided.

6.5.2 Learning a Causal Bayesian Network

When a causal Bayesian network consisting of all possible variables affecting future readmission probability is given, the structure of CNMDP can be formulated as shown in Figure 6.1. If such a graphical structure is not provided, the causal network can be learned from the dataset describing patients status and readmission events. In Chapter 3, a causal Bayesian network (given in Figure 6.2) is learned based on the data set presented in Table 3.1, as shown (see Chapter 3 for details).

In such a graphic network structure, the subset of variables affecting readmission is the set of ancestors of readmission variable in the graph. The variations of other variables not connecting to readmission variables on one direction in the original network do not affect the readmission probability. Thus, the variables considered in the CNMDP model include: "gender (x_s)", "height (x_h)", "weight (x_w)", "anxiety (x_a)", "support (x_s)", "equipment (x_e)", "past readmission (x_p)", and "readmission (y)". The subgraph marked in Figure 6.2 is induced by leaving those variables and arcs between the remaining variables, which will be utilized in the CNMDP model. By duplicating the network in every epoch and connecting them, the dynamic causal network in Figure 6.3 is composed.

Remark 6.3. *Generally speaking, if a given dataset is not obtained from causal experiments such as*

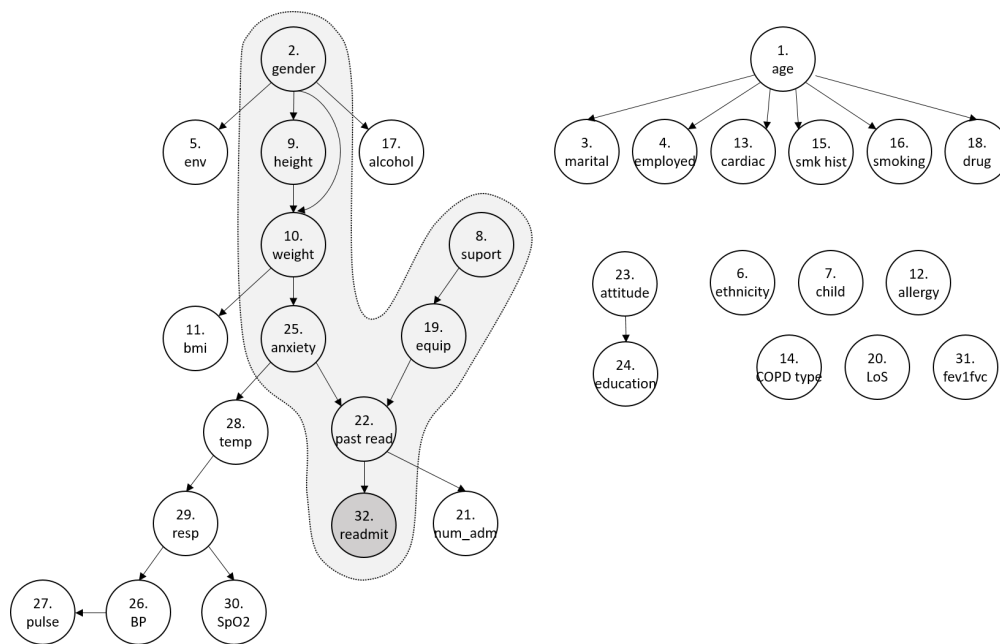


Figure 6.2: A full diagram of the learned causal network for COPD readmission (Figure 3.3 in Chapter 3)

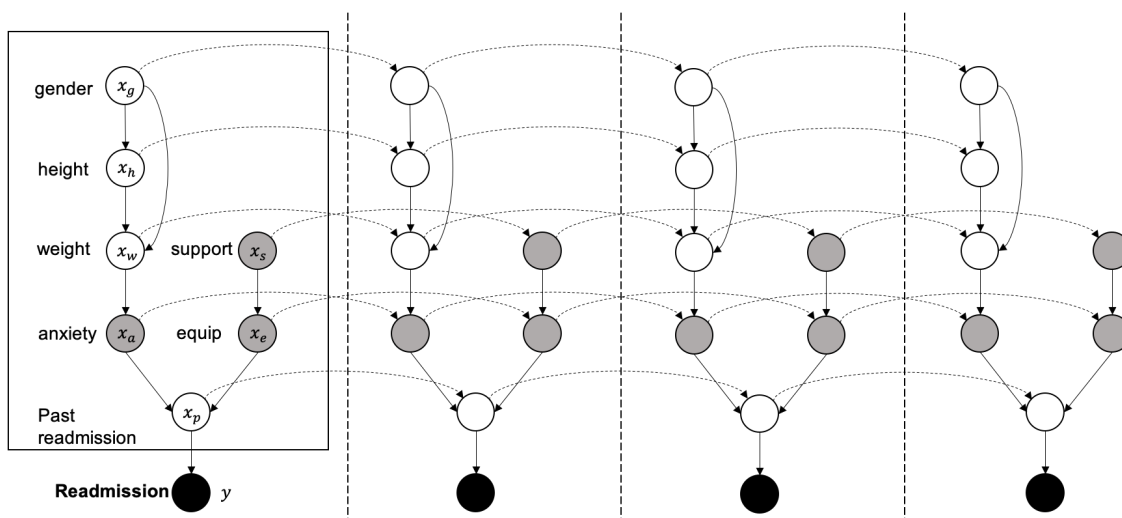


Figure 6.3: CNMDP model for COPD readmission

randomized controlled trials, a machine learning model learned from the dataset cannot guarantee that it represents the causal relationships [138]. Nevertheless, many heuristics have been suggested to identify causal relationships from observation data with background knowledge [151, 152] or without external information [153].

6.5.3 Constructing and Solving MDP Model

As mentioned before, a controllable variable can be manipulated by an intervention, and a variable that does not change over time is referred to as *static*. Among the remaining variables, “anxiety”, “support” and “equip” are classified as controllable, while “gender” and “height” are static since they cannot be manipulated by an intervention. In addition, variable “past readmission” is determined by clinical history, thus it cannot be controlled through any intervention. Furthermore, even though “weight” can be controllable, it is challenging to observe weight change during a short time period (such as four weeks). Therefore, we only consider interventions that can influence x_a (anxiety), x_s (support), and x_e (equip), and denote the interventions as a_a , a_s , and a_e , respectively. These variables with relevant interventions are shaded with grey color in Figure 6.3.

From Lemmas 6.4 and 6.5, the transition probabilities in the defined CNMDP model with the network structure G shown in Figure 6.3 can be obtained from the conditional probabilities in the causal Bayesian network, $p_G(x_h | x_g)$, $p_G(x_w | x_g, x_h)$, $p_G(x_a | x_w)$, $p_G(x_e | x_s)$, $p_G(x_p | x_a, x_e)$, and $p_G(y | x_p)$, and the probabilities regarding the effects of interventions, p_{a_a} , p_{a_s} , and p_{a_e} . Appendix B provides the details of obtaining the first type conditional

probabilities from the learned causal network G . For the second type probabilities, we have

$$\begin{aligned}
& p((\mathcal{X}', y') \mid (\mathcal{X}, 0), a) \\
&= p(y' \mid x'_p, a) p(x'_p \mid x_p, x'_a, x'_e, x_a, x_e, a) p(x'_a \mid x_a, x'_w, x_w, a) p(x'_e \mid x_e, x'_s, x_s, a) \\
&\quad \cdot p(x'_w \mid x_w, x'_h, x'_g, x_h, x_g, a) p(x'_s \mid x_s, a) p(x'_h \mid x_h, x'_g, x_g, a) p(x'_g \mid x_g, a) \tag{6.20}
\end{aligned}$$

$$\begin{aligned}
&= p(y' \mid x'_p) p(x'_p \mid x_p, x'_a, x'_e, x_a, x_e) p(x'_a \mid x_a, a) p(x'_e \mid x_e, x'_s, x_s, a) p(x'_s \mid x_s, a) \\
&\quad \cdot I\{x'_w = x_w, x'_h = x_h, x'_g = x_g\}, \tag{6.21}
\end{aligned}$$

where $I\{\cdot\}$ equals to 1 if the condition in the parenthesis is true and 0 otherwise. Note that x_g , x_h , and x_w have no corresponding interventions and they are assumed invariant over the decision time. Thus,

$$\begin{aligned}
p(x'_g \mid x_g, a) &= p(x'_g \mid x_g) = I\{x'_g = x_g\}, \\
p(x'_h \mid x_h, x'_g, x_g, a) &= p(x'_h \mid x_h, x'_g, x_g) = I\{x'_h = x_h, x'_g = x_g\}, \\
p(x'_w \mid x_w, x'_h, x'_g, x_h, x_g, a) &= p(x'_w \mid x_w, x'_h, x'_g, x_h, x_g) = I\{x'_w = x_w, x'_h = x_h, x'_g = x_g\}.
\end{aligned}$$

In addition, as x_w does not vary, we have

$$p(x'_a \mid x_a, x'_w, x_w, a) = p(x'_a \mid x_a, a).$$

Suppose $p_{a_a} = 0.5$, $p_{a_s} = 0.8$, and $p_{a_e} = 0.6$. For a COPD patient who falls into the states where $x_a = 0$, the optimal decision of intervention at the first decision epoch is a_a . In the case that a patient does not feel anxiety during his/her hospitalization ($x_a = 1$) such that intervention a_a does not affect the patient's readmission risk, the optimal decisions are shown as a_e when $x_e = 0$ and a_e is applicable, otherwise no intervention needs to be taken.

When the intervention efficacy p_{a_a} , p_{a_s} , and p_{a_e} vary, the optimal intervention differs. Figure 6.4 shows how the optimal intervention alternates depending on the probabilities

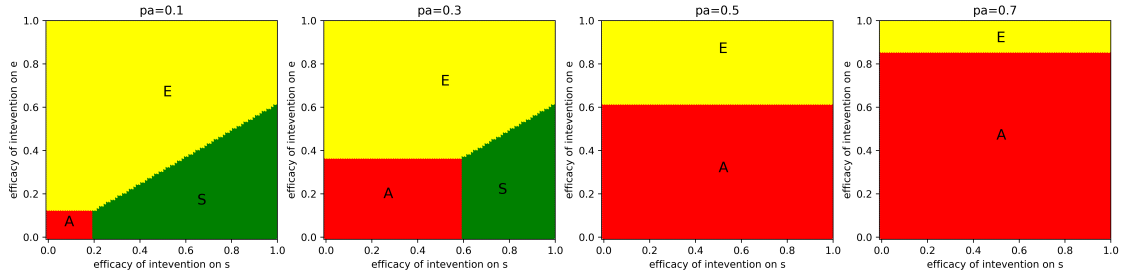


Figure 6.4: Resulting optimal policy depending on different efficacy

for intervention efficacy for patients whose states correspond to $x_a = x_s = x_e = 0$. In every subplot, p_{a_s} and p_{a_e} values are indicated in x-axis and y-axis, respectively. The corresponding p_{a_a} is written in each subplot.

As Proposition 6.1 states, if $p_{a_s} = p_{a_e}$, then a_e is preferred since x_s is a parent node of x_e in the network. From the first plot with $p_{a_a} = 0.1$, it is shown that when p_{a_e} becomes larger, intervention a_e dominates the others. Similar patterns are shown in other subplots as well.

Comparing the different subplots we can observe the effect of p_{a_a} on optimal policies. When p_{a_a} increases, the area where a_a dominates gets larger. The first two subplots show that Proposition 6.3 holds. It is shown that there is a linear boundary between the areas where a_e and a_s dominate, and this boundary doesn't change no matter how other intervention efficacy vary. When p_{a_s} increases, the range of p_{a_e} where a_s can dominate a_e gets larger, which implies that if p_{a_s} is significantly higher than p_{a_e} , the intervention on x_s acts more effectively than the one on its child node x_e , otherwise, a_e is more favorable than a_s in most cases.

6.6 Conclusions

This chapter proposes a causal Bayesian network based Markov decision process (CN-MDP) model to evaluate the effectiveness of interventions for each patient and develop corresponding intervention strategy. In the model, a causal Bayesian network using patient records is employed to predict how much an intervention can affect each patient's health

status. Then, the network is integrated into a Markov decision process to make optimal decision on intervention planning to minimize the risk of health deterioration or readmission. Different from a standard MDP model the proposed CNMDP represents states as vectors of patient features (e.g. risk factors regarding COPD readmission and patient characteristics), which enables personalized and realistic intervention decision. Also, by incorporating the causal Bayesian network, prediction of intervention effects and the associated transition probability setup are facilitated. To illustrate the applicability of the proposed model, a case study at SSM Health St. Mary's Hospital is presented. Furthermore, in order to deliver managerial insights, related propositions are provided and proved. Such a model can be applied to COPD and many other diseases or healthcare problems to design personalized, prompt, and effective intervention plan.

Chapter 7

Optimization of Opioid Prescription Post TJR Surgery

7.1 Introduction

In Chapter 4, a classification model for identifying patients who are expected to use less, moderate, or higher amount of opioids after discharge from TJR surgery. Based on the classification model, in this chapter, we aim to optimize opioid prescription to reduce both leftovers and number of refills based on classification of opioid demands for surgical patients in post-surgery recovery period. Specifically, we propose an analytical framework that integrates classification and optimization. In this framework, a machine learning based classification model first identifies each patient's expected demand level of opioids. Then, based on the prediction, optimal prescriptions derived from a stochastic programming model that minimizes both leftovers and number of refills will be prescribed to the TJR patients. Such a framework can provide support for clinicians to prescribe the optimal quantity of opioids that meets the patient's need for pain control, and in turn help reduce the overall opioid over-prescription and overdose. To illustrate the applicability and benefits of the framework, a case study at a community hospital for TJR patients is introduced.

7.2 An Integrated Opioid Prescription Optimization Framework

7.2.1 Problem Description

When opioids are prescribed, a one-size-fits-all pattern is typically followed in practice. For example, in SSM Health Saint Mary's Hospital, most of the patients are prescribed 60 pills of 5mg oxycodone, which is equivalent to 450 mg oral Morphine, as a standard. Such a universalized practice can lead to a large amount of unused medicines, which may contribute as a source for potential non-medical use. Figure 7.1 shows the amounts of opioids being prescribed at discharge and used in a two-week post-discharge period. The values are derived from the collected data in the preceding survey study (see details in Chapter 4). A point in the blue (dark) line indicates the normalized amount of opioids that a patient has used in two weeks after discharge, while the indices for the corresponding patients are sorted by the used opioid amounts. The yellow (light) bar in the background illustrates the amount of opioids prescribed at discharge (almost uniformly with 450 MME as the current prescription standard) for the corresponding patient, which suggests that the prescribed amount is not closely correlated to the actual usage. Patients who use more opioids than the prescribed (patients whose yellow bars are under the blue line) will request refills, and the green cross marker in the figure indicates the number of refills requested by a patient in the same period.

As one can see, for patients who show patterns of less usage than others, fewer opioids should be prescribed to reduce excessive opioids. Likewise, for patients who may require more medicine, more opioids should be prescribed so that the number of refills can be reduced. In other words, the precise amounts of opioids should be prescribed to meet the actual demands of TJR patients. To achieve this, an integrated analytical framework for opioid prescription optimization is proposed. Figure 7.2 depicts the outline of the entire framework.

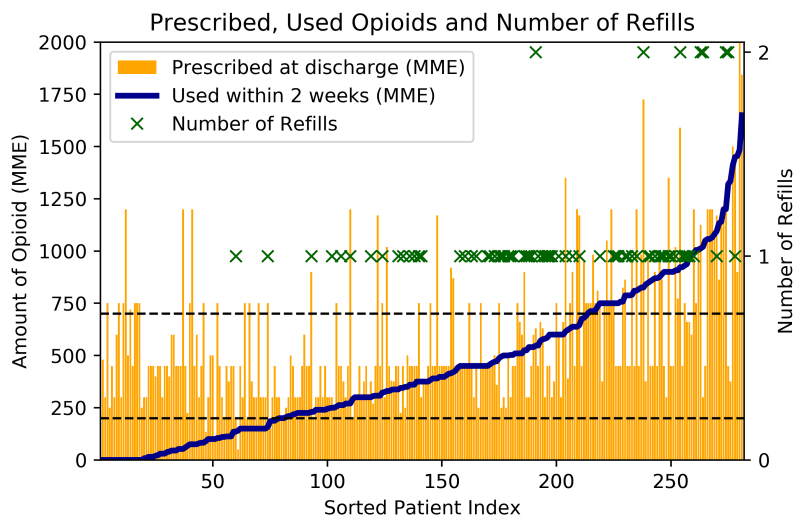


Figure 7.1: Patterns of opioid prescribed amount at discharge and usage within 14 days

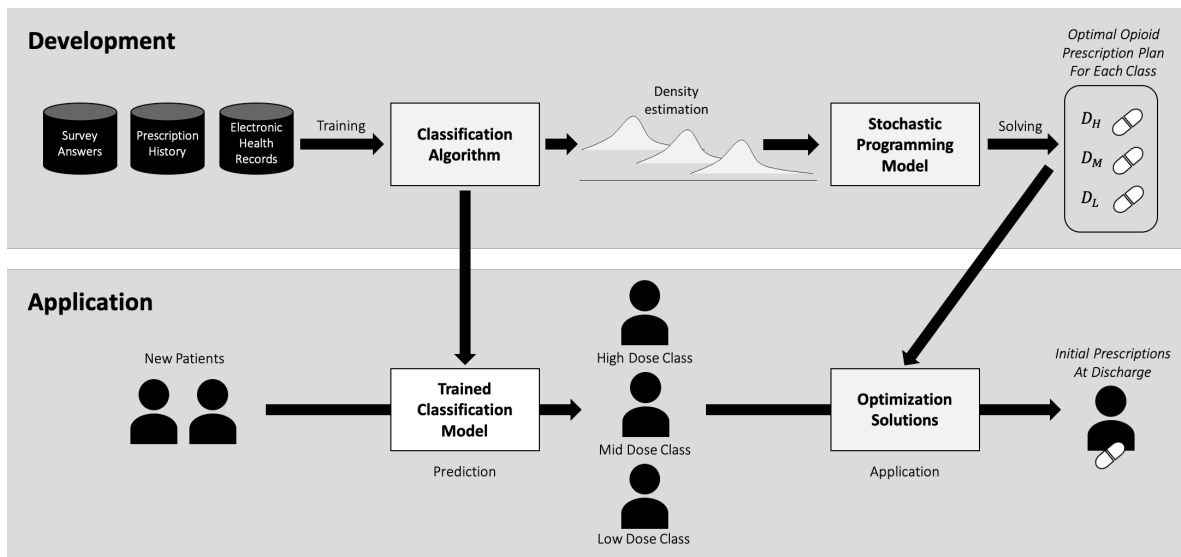


Figure 7.2: An analytical framework of classification and stochastic programming models

As shown in Figure 7.2, the optimal decision of opioid prescription is based on each patient's expected opioid usage, which are predicted through a machine learning model developed in Chapter 4. Thus, the integrated framework consists of two main phases: classification and optimization. As the first phase, data collection and classification model development are conducted and presented in Chapter 4. Then, in the second phase, which this chapter focuses on, we develop a stochastic programming model involving each patient's opioid demand as uncertain parameters to minimize both leftovers and number of refills. The parameters of the optimization model, including the distributions of random opioid demands, are estimated based on the trained classification model, and the optimal opioid prescription plans are obtained by solving the stochastic program. After the classification model is trained and the stochastic model is solved once, the trained model and prescription solution will be used in practice to classify new TJR patients and apply the obtained optimal prescription plans to them.

Remark 7.1. *Note that a classification model rather than a regression is employed in this study for two reasons. Firstly, the collaborating healthcare practitioners prefer classifications that group the patients into specific categories to the diverse numerical results of regression. Secondly, if a regression model is employed, the regression error will directly impact the result of opioid leftover as the predicted demand will determine the prescribed amount. Whereas, by classifications, the difference between the actual opioid demand and the prescribed amount can be mitigated as the classified demand level will cover a range of demand values.*

Below the formulation of the proposed optimization model is introduced.

7.3 Stochastic Programming for Opioid Prescription

7.3.1 Model Formulation

The main goal in this study is to reduce opioid leftovers. Clearly, a simple approach is to prescribe a minimal amount of opioids to all patients and provide refills whenever they request. However, frequent refills can lead to substantial inefficiency, higher workload, and inconvenience, from all perspectives of physicians, pharmacists, and patients, thus are not favorable. In fact, hospitals intend to prescribe enough amount in a moderate level to most patients in order to avoid recurrent refill requests which increases the administrative burden on providers and their staff [34,154]. Therefore, the objective should be set to avoid frequent refills as well as to reduce leftovers.

In this research, different opioid demand levels of new TKR/THR patients are predicted by the preceding classification model. As mentioned above, the patients are classified into three groups, denoted by L , M , and H , to distinguish and identify them as demanding less, moderate, or more opioids, respectively. Then, the optimal prescription should be determined for each group based on their needs. To formulate the model, the following assumptions are introduced:

Assumption 7.1. *The amount of opioids in subsequent refills is the same as in the initial prescription.*

Assumption 7.2. *The actual opioid use of each patient is random, and each opioid-use level class (i.e., Groups L , M , and H) has its own opioid demand distribution.*

Assumption 7.1 represents the typical refill pattern in current practice in many hospitals, where the amount of refilled opioids during a short time period after discharge is usually the same as that in initial prescription [154]. Assumption 7.2 emphasizes that the demands in each group are uncertain before filling prescriptions, and indicates that they may follow

different distributions, which is in accordance with intuition and the data obtained in this study (see, for instance, the pattern in Figure 7.1).

In order to handle the case of random opioid demands, a stochastic programming model is developed. Let G and K be the set of opioid-use level groups and the set of patients who receive TJR surgeries, respectively. For each patient $k \in K$, introduce $x_I^{(k)}$, $x_A^{(k)}$, and $y^{(k)}$ as the amount of initially prescribed opioids, the total amount of opioids prescribed within a given time period (e.g., two weeks) after discharge, and the number of refills requested within the same period, respectively. Let $R^{(k)}$ be the uncertain opioid demand required by patient k , and $w^{(k)}$ be the remaining amount of opioids, namely, leftover, at the end of the time period. In addition, for each patient $k \in K$ and each group $g \in G$, denote $R_g^{(k)}$ as the uncertain parameter representing opioid demands required by patient k in group g . Also, introduce $z_g^{(k)}$ as the binary parameter indicating patient k 's group classification, i.e., $z_g^{(k)} = 1$ if patient k is classified to opioid-use group g , and 0 otherwise. The *decision variables* in the model are $D_g, g \in G$, which are the amounts of opioids that should be prescribed to a patient in group g at the time of discharge.

Using these definitions, a two-stage nonlinear stochastic programming model for optimal opioid prescriptions is formulated as Problem 7.1 below:

Problem 7.1.

$$\text{minimize} \quad \mathbb{E}_R \left[\sum_{k \in K} w^{(k)} + \lambda \sum_{k \in K} y^{(k)} \right] \quad (7.1a)$$

$$\text{subject to} \quad x_I^{(k)} = \sum_{g \in G} D_g z_g^{(k)}, \quad \forall k \in K \quad (7.1b)$$

$$R^{(k)} = \sum_{g \in G} R_g^{(k)} z_g^{(k)}, \quad \forall k \in K \quad (7.1c)$$

$$x_A^{(k)} = x_I^{(k)} (y^{(k)} + 1), \quad \forall k \in K \quad (7.1d)$$

$$x_A^{(k)} \geq R^{(k)}, \quad \forall k \in K \quad (7.1e)$$

$$w^{(k)} = x_A^{(k)} - R^{(k)}, \quad \forall k \in K \quad (7.1f)$$

$$D_g \geq 0 \quad \forall g \in G \quad (7.1g)$$

$$y^{(k)} \geq 0, y^{(k)} \in \mathbb{Z}, \quad \forall k \in K \quad (7.1h)$$

Objective function (7.1a) combines the two proposed objectives to minimize the expected sum of all opioid leftovers while minimizing the expected total number of refills. Parameters $z_g^{(k)}$ of patient group g are assigned by a classification model. Due to the classification results, the values of $z_g^{(k)}$ satisfy $\sum_{g \in G} z_g^{(k)} = 1$ for all k . Group $g \in G$, where patient $k \in K$ belongs to, determines the initial opioid prescription for the patient using the optimal prescription decision in (7.1b), also it determines the distribution of opioid demand for patient $k \in K$ in (7.1c). In (7.1d), the total amount of opioids prescribed to patient $k \in K$ is obtained by the sum of initially prescribed amount $x_I^{(k)}$ and the amounts in subsequent refills $x_I^{(k)} \cdot y^{(k)}$, which should satisfy the patient's demand described in (7.1e). The wasted amount of opioids is the difference between prescribed opioids and demand during the given time period, as shown in (7.1f).

As the proposed problem is nonlinear and mixed-integer, to effectively solve the problem without relaxation, first, model decomposition and associated propositions are presented below as viable solution approaches. Then, distribution estimation is described to determine the distribution of random opioids demands of patients in each group.

7.3.2 Model decomposition

Note that Problem 7.1 involves a nonlinear constraint in (7.1d) and integer variables $y^{(k)}$, which makes the problem arduous to solve in general. However, as we observe, the decision variable D_g in Problem 7.1 is associated with the corresponding group $g \in G$. When the decision variable is fixed and group values $z_g^{(k)}$ are given, terms $w^{(\cdot)}$ and $y^{(\cdot)}$ and constraints for different patients k and k' are not interacting. Due to this fact, Problem 7.1 can be decomposed into simpler subproblems, which will facilitate solving the original problem. Therefore, we introduce Propositions 7.1 and 7.2 below to decompose the model and solve it efficiently.

First, we introduce decomposition at the class level.

Proposition 7.1 (Class-level decomposition). *Each decision variable D_g , $g \in G$ in the original formulation can be obtained by solving Problem 7.2 below separately: For each group $g \in G$,*

Problem 7.2.

$$\text{minimize} \quad \mathbb{E}_{R_g} \left[\sum_{k \in K_g} w^{(k)} + \lambda \sum_{k \in K_g} y^{(k)} \right] \quad (7.2a)$$

$$\text{subject to} \quad x_A^{(k)} = D_g (y^{(k)} + 1), \quad \forall k \in K_g, \quad (7.2b)$$

$$x_A^{(k)} \geq R_g^{(k)}, \quad \forall k \in K_g, \quad (7.2c)$$

$$w^{(k)} = x_A^{(k)} - R_g^{(k)}, \quad \forall k \in K_g, \quad (7.2d)$$

$$D_g \in \mathbb{R}_+, \quad y^{(k)} \in \mathbb{Z}_+, \quad \forall k \in K_g, \quad (7.2e)$$

where $K_g \subset K$ is a set of patients classified in group $g \in G$.

Proof: See Appendix C. ■

Proposition 7.1 decomposes original Problem 7.1 into $|G|$ subproblems in the form of Problem 7.2. It results in one decision variable in each subproblem with less constraints.

Note that (7.1b) and (7.1c) are not required in Problem 7.2 since all patients in K_g are in the same class g and a single demand distribution is involved in each subproblem. Also, the numbers of variables and constraints are substantially reduced because K_g is a subset of K . Therefore, Proposition 7.1 simplifies the original problem significantly.

In addition, since the constraints for different k 's are independent and patients in the same g share the same demand distribution, Problem 2 can be further simplified. Thus, we consider individual level decomposition next.

Proposition 7.2 (Individual-level decomposition). *The optimal decision for D_g can be obtained by solving Problem 7.3 below: For each $g \in G$,*

Problem 7.3.

$$\min_{D_g} \mathbb{E}_{R_g}[D^g(y+1) + \lambda y] \quad (7.3a)$$

$$s.t. \quad D^g(y+1) \geq R_g, \quad (7.3b)$$

$$D_g \in \mathbb{R}_+, y \in \mathbb{Z}_+. \quad (7.3c)$$

Proof: See Appendix C. ■

Due to Proposition 7.2, we can obtain the optimal value of the same decision variable D^g with $\frac{1}{|K_g|}$ times less constraints. Then, from Propositions 7.1 and 7.2, we can solve Problem 7.1 that has $|G|$ decision variables with $6 \times |K| + |G|$ constraints, $|K|$ integer variables, and $|K|$ random variables, by solving Problem 7.3 for $|G|$ times, which only involves a single decision variable with two constraints, one integer variable, and one random variable.

In the two-stage formulation of Problem 3, D^g is the first-stage decision variable and y is the second-stage decision variable. To solve Problem 7.3 for $g \in G$, suppose that there are a set of samples, $\Xi = \{\xi_1, \dots, \xi_S\}$, where $S = |\Xi|$, generated from distribution of R_g with the same probabilities. For any decision D , we can estimate the expected value of the

objective function by averaging its values related to the realized samples, which leads to the so-called sample average approximation (SAA). Then, since variable y depends on a realized value ξ of R_g , we can represent the problem as follows:

Problem 7.4.

$$\min_{D \geq 0} \frac{1}{S} \sum_{s=1}^S [D(y_s + 1) + \lambda y_s] \quad (7.4a)$$

$$\text{s.t. } D(y_s + 1) \geq \xi_s, \quad s = 1, \dots, S, \quad (7.4b)$$

$$y_s \geq 0, y_s \in \mathbb{Z} \quad s = 1, \dots, S. \quad (7.4c)$$

From constraints (7.4b) and (7.4c), we can obtain $y_s \geq \frac{(\xi_s - D)_+}{D}$. Since the model will minimize y , which is an integer, we solve the second-stage mixed-integer problem as

$$y_s = \left\lceil \frac{(\xi_s - D)_+}{D} \right\rceil, \quad s = 1, \dots, S.$$

Therefore, we can finally obtain the following 1-dimensional optimization problem (Problem 7.5) without any constraints except non-negativity.

Problem 7.5.

$$\min_{D \geq 0} \frac{1}{S} \left\{ D \cdot \left(\sum_s \left\lceil \frac{(\xi_s - D)_+}{D} \right\rceil + S \right) + \lambda \sum_s \left\lceil \frac{(\xi_s - D)_+}{D} \right\rceil \right\}. \quad (7.5a)$$

Comparing to the large-scale mixed integer programming using a large number of samples, the final sub-problem is easier to solve because it only involves one directional line search and the possible solution region is bounded due to the following theorem.

Theorem 7.1. *The optimal solution of Problem 7.5, D , is bounded to the maximum sample value, i.e., $D \leq \max(\xi)$.*

Proof: See Appendix C. ■

Therefore, searching the minimum of $f(x)$ by increasing x from 0 to $\max_{\xi \in \Xi} \xi$, the optimal solution x^* can be obtained.

The optimal solution of Problem 7.5 for a given λ is also the solution of Problems 7.2, 7.3, and 7.4. By solving Problem 7.5 for $g \in G$, the optimal decision D_g^* is attained. The optimal average leftover $w_g(D_g^*)$ and number of refills $y_g(D_g^*)$ for patient $k \in K_g$ can be estimated as a result of optimization as they are represented by the two terms in (7.5a).

$$\begin{aligned}\hat{w}_g(D_g^*) &= \frac{1}{S} \left\{ D_g^* \cdot \left(\sum_s \left\lceil \frac{(\xi_s - D_g^*)_+}{D_g^*} \right\rceil + S \right) - \sum_s \xi_s \right\}, \\ \hat{y}_g(D_g^*) &= \frac{1}{S} \sum_s \left\lceil \frac{(\xi_s - D_g^*)_+}{D_g^*} \right\rceil.\end{aligned}$$

Then, the optimal values $D_g^*, \forall g \in G$, compose the solution of Problem 1, D^* . The objectives, i.e., the expected total leftover $\mathbb{E}_R[\sum_{k \in K} w^{(k)}]$ and number of total refills $\mathbb{E}_R[\sum_{k \in K} y^{(k)}]$, can be obtained by aggregating the objectives of Problem (7.5a) as follows.

$$\begin{aligned}\mathbb{E}_R \left[\sum_{k \in K} w^{(k)} \right] &\approx \sum_{g \in G} \left(|K_g| \cdot \hat{w}_g(D_g^*) \right), \\ \mathbb{E}_R \left[\sum_{k \in K} y^{(k)} \right] &\approx \sum_{g \in G} \left(|K_g| \cdot \hat{y}_g(D_g^*) \right).\end{aligned}\tag{7.6}$$

7.3.3 Demand estimation

Problem 7.1 introduced in Subsection 7.3.1 handles the demand values as uncertain parameters. However, the information of underlying distributions is not given. Since the demand distribution of class g should represent the demand values of new patients who are classified into g , the test data with the predicted classes can be used for density estimation.

The distribution of the random variable $R_g^{(k)}$ representing opioid demands of patients in class g can be estimated in various ways, specifically through one of the following

estimation functions, where x_1, x_2, \dots, x_n are the historical opioid usages of n patients classified into the corresponding class.

- Empirical distribution function: Where $\mathbb{I}\{A\}$ is an indicator of event A ,

$$\hat{f}_{emp}(x) = \frac{\partial}{\partial x} \hat{F}_{emp}(x), \text{ where } \hat{F}_{emp}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{x_i \leq x\}. \quad (\text{EMP})$$

- Kernel density estimation: For a kernel function K and a bandwidth h ,

$$\hat{f}_{kde}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right). \quad (\text{KDE})$$

- One of probability distributions that can fit the data (e.g., Gaussian distribution).

7.4 Results and Solutions

Using a case study at the collaborating community hospital, this section provides a demonstration of the proposed framework.

7.4.1 Resulting expected demand densities

Using the classification results, distributions of opioid demands for the classified patients can be estimated. Figure 7.3 shows the estimated densities by using kernel density estimation (KDE) in Subsection 7.3.3, where a standard normal density function is used for kernel function K and the bandwidth is computed by the rule of thumb of Scott ([155]).

The kernel density plots depicted in Figure 7.3 suggest that the densities of used opioids are slightly skewed to the left, but unimodal, which implies that the demand densities approximately follow a truncated Gaussian. Hence, for analytical simplicity, the truncated Gaussian distribution is used as the family of distributions for opioid demands in this study. The final estimations for distributions of opioid demands in Groups

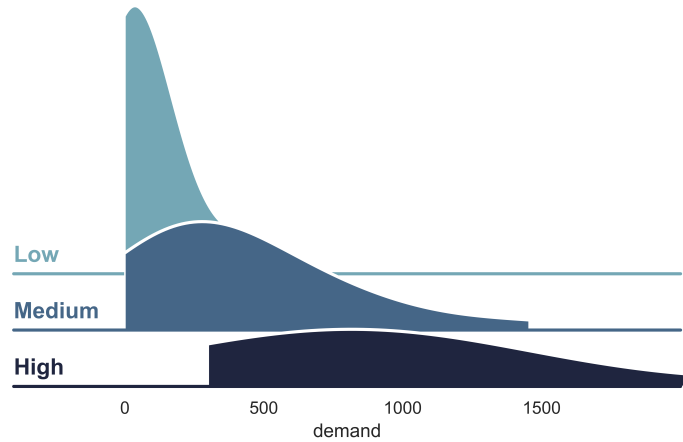


Figure 7.3: Results of opioid demand distribution estimation

L , M , and H are chosen as $\max\{\mathcal{N}(95.75, 158.58^2), 0\}$, $\max\{\mathcal{N}(408.05, 349.28^2), 0\}$, and $\max\{\mathcal{N}(1034.17, 756.62^2), 0\}$, respectively.

7.4.2 Resulting optimal values

By substituting the estimated distributions into the proposed stochastic programming Problem 7.5, we derive optimal solutions, as shown in Figure 7.4 for Group L patients. Similar figures can be obtained for other groups. Since the objective values and optimal solutions are dependent on λ , which is the ratio between two objectives, the total opioid leftovers and the total number of refills, the resulting solutions can differ with respect to λ . Thus, in the figure, we provide different optimal objective values with various λ 's, where the expected opioid leftovers and average numbers of refills are shown in vertical and horizontal axes, respectively.

After obtaining the resulting solutions of Problem 7.5 for all three classes, we can aggregate them into the final solution and objectives of Problem 7.1 by using equation (7.6). For better visualization and comparison purpose, the aggregated results are presented at an individual level by dividing the objectives by number of patients, and are shown as the red star (thick and dark) line for various λ 's in Figure 7.5. The blue dot, orange cross, and green plus lines characterize L , M , and H groups, respectively (note that the

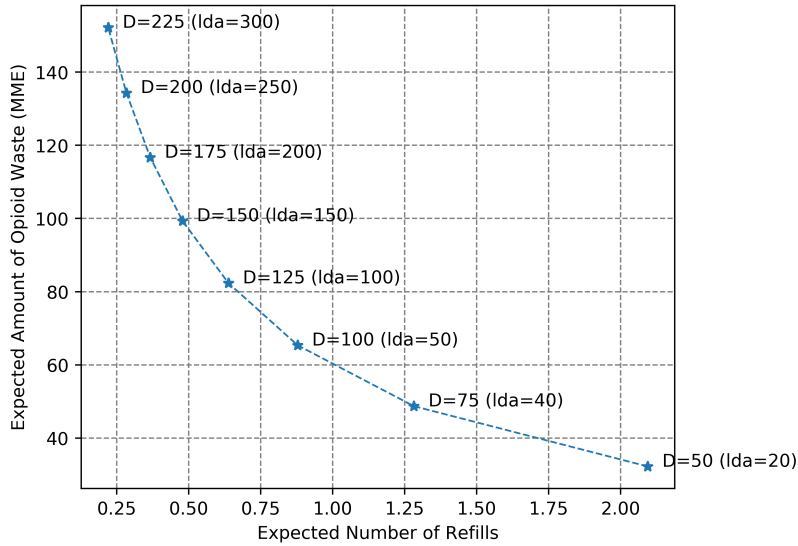


Figure 7.4: Optimal solutions of Problem 7.5 for Group L

orange cross line for Group M is overlapped with the red star line). The horizontal and vertical axes represent $\frac{1}{|K|}\mathbb{E}_R \sum_{k \in K} y^{(k)}$ and $\frac{1}{|K|}\mathbb{E}_R \sum_{k \in K} w^{(k)}$, respectively. Each point in the red star line has a corresponding decision $D^* = (D_L^*, D_M^*, D_H^*)$, and is linked by grey dash lines to the corresponding values of D_L^* , D_M^* , and D_H^* in other lines, which are the solutions of group-wise problem. Using such a figure, practitioners can check if their choices are feasible or not. For instance, it may not be feasible to achieve both less than 150 MME leftover and less than 0.5 refill simultaneously.

Remark 7.2. The weight λ in the multi-objective function (equation (7.1a)) can be viewed as a general gauge of relative importance for each objective (i.e., the expected total opioid leftovers and the expected number of total refills). Even though there is no fundamental guideline for selecting weights for accurate a priori articulation of preference [156], in this context, we can provide a view of weight as an incremental cost of reducing the number of refills in terms of the amount of opioid leftover. Let us consider a group-wise problem. With a fixed λ , suppose that the optimal solution is x^* and the optimal values of the objective function, the expected total leftover, and the expected number of refills are denoted by $F(x^*)$, $F_1(x^*)$, and $F_2(x^*)$, respectively. Then, $F(x^*) = F_1(x^*) + \lambda F_2(x^*)$ by equation (7.2a). If one wants to reduce the expected number of refills by 1.0 by using another

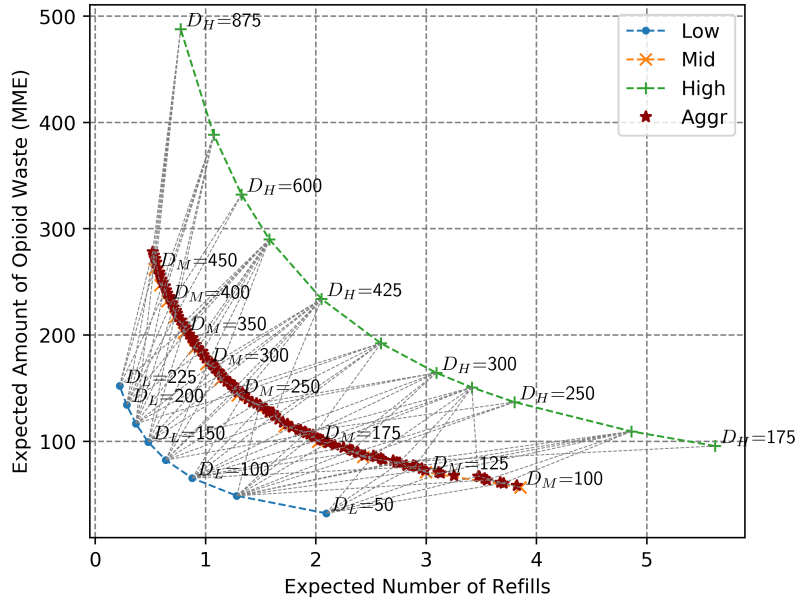


Figure 7.5: Aggregating the optimal solutions into final solution of the original problem

solution x' , then

$$F_1(x^*) + \lambda F_2(x^*) = F(x^*) = \max_x F(x) \leq F_1(x') + \lambda F_2(x') = F_1(x') + \lambda(F_2(x^*) - 1).$$

Thus, $F_1(x') \geq F_1(x^*) + \lambda$ and the expected opioid leftover increases at least λ . Therefore, weight λ can be interpreted as an incremental cost of reducing number of refills in terms of the expected opioid leftovers.

7.4.3 Comparison with one-size-fits-all practice

To demonstrate the advantage of the integrated classification and optimization framework over the prevailing one-size-fits-all approach, the stochastic problem without classifying patients is solved to represent current practice and compared with the proposed solutions. Using the same population in the test data employed in the classification, we estimate the overall opioid-use density as $\max\{\mathcal{N}(452.12, 494.84), 0\}$. The results are shown as cross marks on the purple (thin) line in Figure 7.6. The red dots with lines are the same result in Figure 5, which represent the resulting expected unused opioid amount and number

of refills from the aggregated solutions. As the connections in Figure 7.5 indicate, each red point in Figure 7.6 links to a combination of optimal prescriptions to three classes and the corresponding solution. For example, $D = (D_L, D_M, D_H)$ with parameters $\text{lda} = (\lambda_L, \lambda_M, \lambda_H)$ are annotated to the points in Figure 7.6. Comparing the two approaches, i.e., the proposed optimization (red) and one-size-fits-all (purple) solutions, it is shown that the integrated framework can provide better results.

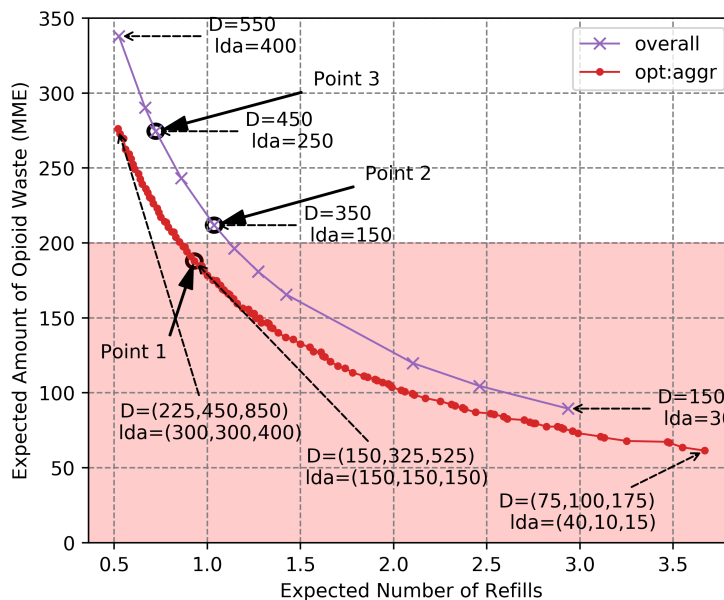


Figure 7.6: Comparison of the proposed framework with one-size-fits-all prescriptions

The selection of a final solution in multi-objective relies on preference of decision maker [157]. Opioid prescribers, e.g., orthopedic surgeons in this case, can decide the final prescription policy when the resulting solutions and their expected outcomes are given as in Figure 6. For example, if a practitioner prefers decisions resulting in less than 200 MME leftovers in average, the shaded area in Figure 7.6 outlines the candidate decisions. Then the practitioner can select a decision in this area, such as Point 1, with the optimal prescriptions for Groups L , M , and H as $D_L = 150$, $D_M = 325$, and $D_H = 525$, with $\lambda = 150$, respectively. Points 2 and 3, which represent the one-size-fits-all decisions, are compared with Point 1 in the figure. As one can see, the decision in Point 2, which uses the same λ ,

results in more average opioid leftovers and more frequent refill requests. Whereas, Point 3 indicating the current standard practice (i.e., 450 MME) will increase the average leftovers by 46%.

Table 7.1 presents the optimal results for the 57 patients in the selected test data set. The resulting values are analyzed based on the classification results on the test data set. Since the future patients who are supposed to be prescribed with D_g are those being classified to class $g \in G$, the resulting values are valid regardless of the classification model's accuracy. Such values include the mean opioid leftovers and average number of refills derived from the chosen optimal decision in the proposed framework. The second and third columns ("Total") are calculated based on the data of all 57 patients, and the results in rest columns are computed using the patient populations predicted into groups L , M , and H . In addition, the actual opioid leftovers and refills realized in the past are provided in the last row. Note that there exists slight variation in opioid prescriptions in actual practice (as shown in Figure 7.1), which also allows minor changes of opioid amount in subsequent refills. Thus, the actual results in the last row are slightly different from those of standard policy ($D = 450$ MME). But the actual amount of opioid leftover is relatively close to that in current standard. Therefore, the decision of $D = 450$ MME, rather than the actual result, will be considered as the current policy in the rest of comparison.

Table 7.1: An example of practicing an optimal solution on the collected data

	Total		Group L		Group M		Group H	
	Left	Refill	Left	Refill	Left	Refill	Left	Refill
Optimal Decision $D_L = 150, D_M = 325, D_H = 525$	145	0.75	114	0.40	155	0.74	137	1.22
One Class Decision $D = 350$	186	0.82	289	0.10	171	0.66	137	2.33
Current Policy $D = 450$	282	0.63	399	0.10	254	0.47	270	1.89
Actual $D \sim 450$ (Adjusted Refill)	249	0.32	354	0.10	252	0.44	122	0.44

The results of average leftover for all patients, represented in the second column, indicate

that the optimal decision (Point 1 in Figure 7.6) derived from the proposed framework results in the least opioid leftover, 145 MME/patient. It is apparent that the opioid leftovers are different when the optimal decision is executed ($p=2.9 \times 10^{-9} < 0.001$ in t-test). Since the number of patients involved in this case study is 285, a total of $(282 - 145) \cdot 285 = 39,045$ MME of opioid leftover can be reduced when compared to the current policy. The optimal decision also leads to less leftovers and less refill requests than those of one-class decision $D = 350$ MME. In addition, patients in all groups will leave less opioids on average when the optimal decision obtained by the proposed framework is practiced. Therefore, the results demonstrate that the proposed framework can contribute to substantial leftover reduction.

A more desirable benefit is revealed when certain group averages of leftovers and refills are compared. Given that only a short period (two weeks) is considered, patients who require more opioids may tend to use them for a longer time. Then high-opioid-use patients (Group H) may continue taking opioids even after the two-week period, so that the opioids kept by a Group H patient, 137 MME in the chosen decision and 270 MME in current policy, might not be the eventual leftovers. But the remaining opioids in Groups L or M are influential since they are less likely to be used in the future. This implies that the remaining opioids of a patient in Group L , which are 114 MME in the chosen optimal decision and 399 MME in current policy, will not be consumed anymore. Thus, the proposed framework can substantially reduce leftovers in Group L compared to one-size-fits-all decisions. Similar results are observed for Group M patients as well. Moreover, the reduction in opioid leftover can be more substantial if a longer period is considered. Hence, the proposed framework can contribute significantly to help curtail the opioid crisis.

7.5 Sensitivity Analysis

One of the main contributions of the proposed opioid prescription optimization model is to integrate patient classification based on their expected opioid-use level. In the classification phase, different models are trained and are mainly evaluated in terms of accuracy or classification error. Since even the best model can hardly achieve an ideal result (i.e., 100% accuracy), classification errors may affect the estimation of opioid density and the subsequent optimal solutions, which can result in extra opioid leftovers or more refills for individual patients in practice. Thus, such errors should be taken into account. To accommodate this, a sensitivity study of optimal decisions with respect to classification accuracy is of importance.

In this study, the classification results and the distributions estimated based on them establish the information of underlying probabilities for uncertain variables in the stochastic model. Instead of directly substituting classification results into a stochastic program, the demand information of classified patient groups is converted into three different opioid demand distributions at the density estimation step. Such a conversion makes the optimization model less sensitive to the classification result, and it can also reflect the distribution of opioid demands for new patients who should be classified into one of the groups. This conversion process breaks down the impact of deviation in classification errors on the quality of optimal solutions and resulting outcomes into two subsequent parts: the resulting opioid demand distributions and the final decisions of stochastic model. The separation of impact facilitates investigation of complicated influence of classification errors in sensitivity analysis by analyzing them in two aspects.

Next, the sensitivity of density estimation is studied first. Then, the sensitivity of optimal decision based on different distributions is analyzed.

7.5.1 Classification accuracy and density estimation

First, we investigate how distributions change with respect to classification accuracy. To test this, different classification models with various accuracy are emulated using the following procedure: When the data with true amounts of used opioids and a target accuracy α are given, we emulate two binary classification models (e.g., classified to Group L and not to Group H) having accuracy $\sqrt{\alpha}$ so that the combined classification result achieves accuracy α . First, labels for all patient records are initialized to 0. Then, randomly sample $(1 - \sqrt{\alpha})$ of patients who use less than or equal to 200 MME opioids, and add 1 to the labels of these samples. Likewise, sample $\sqrt{\alpha}$ of patients who take more than 200 MME opioids, and add 1 to those sample labels. Following these steps, a classification model detecting Group L patients with $\sqrt{\alpha}$ accuracy is emulated. Analogously, sample $\sqrt{\alpha}$ of patients who use more than 700 MME opioids and $(1 - \sqrt{\alpha})$ of whom use less than or equal to 700 MME opioids, and add 1 to the sample labels. Then, patients are classified into three groups by assigning Groups L , M , and H to the records whose labels are equal to 0, 1, and 2, respectively. Finally, we can evaluate the mean μ_g and standard deviation σ_g for each Group $g \in \{L, M, H\}$.

The analysis of resulting means and variations obtained from the emulated classification models with accuracy ranging from nearly 50% to 99.99% indicates a pattern that the resulting three demand distributions become more distinguished as classification accuracy increases. To quantify the distinction, the distances between distributions are measured by the Kullback-Leibler divergence (KL-divergence) ([158]), which is a measure of how a probability distribution is different from another one. For continuous probability distributions P and Q defined in the same probability space, the KL-divergence of Q from P is defined as:

$$D_{\text{KL}}(P||Q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx.$$

As aforementioned, the distribution family for opioid demands in the proposed model

is assumed as Gaussian. Let P and Q be $\mathcal{N}(\mu_p, \sigma_p^2)$ and $\mathcal{N}(\mu_q, \sigma_q^2)$, respectively. Then, the KL-divergence is

$$D_{\text{KL}}(P||Q) = \log \frac{\sigma_q}{\sigma_p} + \frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{2\sigma_q^2} - \frac{1}{2}.$$

Figure 7.7 summarizes the resulting values of KL-divergence according to different classification accuracy. The vertical values of the points in the figure are the averages of $D_{\text{KL}}(\mathcal{N}(\mu_M, \sigma_M^2)||\mathcal{N}(\mu_L, \sigma_L^2))$ and $D_{\text{KL}}(\mathcal{N}(\mu_M, \sigma_M^2)||\mathcal{N}(\mu_H, \sigma_H^2))$, which indicate how the demand distribution for Groups L or H are away from the center Group M . Due to randomness in the emulating procedure, variability exist in both final classification accuracy and resulting average KL-divergence. The points in the figure characterize the mean results, while the horizontal lines represent the standard deviations of the accuracy and horizontal bars are the variances of average KL-Divergence. From Figure 7.7, we observe that as the classification accuracy increases, the resulting distributions of three classes get more distinct from each other. As the demand distributions in the case study are estimated based on 63% accuracy in aggregated classification, the class distinction following this accuracy is not high enough. However, the distinction can be substantially increased if the classification performance becomes better.

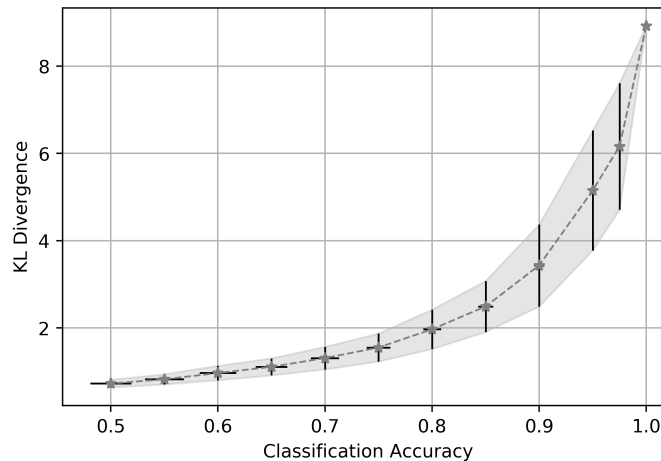


Figure 7.7: Relationship between classification accuracy and resulting density estimation

7.5.2 Classification accuracy and optimization result

From extensive experiments, it has been shown that the resulting optimal decision leads to more reduction of opioid leftovers with the same or a smaller number of refill requests if the distributions of resulting group demands are more distinguished. Hence, it is of interest to know how optimal solutions and resulting outcomes are affected when the classification model changes and thereby different demand distributions are estimated and applied in the stochastic program. The results of sensitivity analysis are provided in Figure 7.8, where the points represent the final solutions with their resulting outcomes (number of refills in x-axis and unused opioid amount in y-axis) when 60%, 80%, 90%, 99.99% classification models are involved in the framework. In the figure, the red points representing the optimal solutions with a 99.99% classification accuracy are the lowest, which signifies that they result in the least leftovers with the same number of refill requests among the solutions using other classification models. As the classification error increases, the aligned line of corresponding solutions rises up in the figure and leads to more leftovers. Therefore, if the classification model gets more accurate, the reduction in opioid leftovers with the same number of refills is improved. The difference in leftover reduction is more significant when the expected number of refills per patient is between 0 and 0.5, which makes easier to be acceptable for practitioners. The error of our final classification model can lead to a substantial leftover reduction, which should be located between blue and orange lines. If the classification model is improved by exploiting a larger data set, more reduction in opioid leftovers with less refills are expected.

7.6 Conclusions

In this study, an analytical framework integrating an opioid-use level classification model and a stochastic programming model is proposed to determine the optimal amounts of opioid prescriptions after receiving total joint replacement surgeries. The framework aims

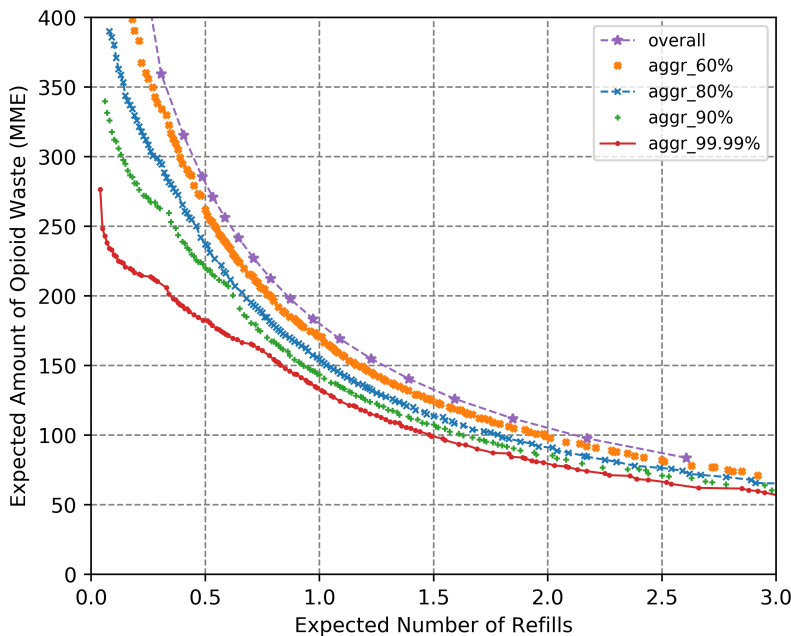


Figure 7.8: Relationship between classification accuracy and resulting optimal objectives

to reduce opioid surplus that can act as a reservoir for potential opioid misuse, while limiting the expected refill requests that lead to a burden for healthcare providers and patients. In the proposed framework, the classification model is trained so that distributions of opioid demands are estimated and a new patient can be classified through the trained model. Then, a stochastic program can be solved based on the estimated distributions, and the practitioners can select optimal solutions of opioid amounts for each group based on their needs. Sensitivity analysis is carried out to understand impact of classification errors on final solutions. As a conclusion, classification models with higher accuracy lead to more distinct opioid demand distributions and result in better solutions with less opioid wastes. To demonstrate advantages of the framework, computational results are provided based on the data collected from the collaborating community hospital. The results prove that the model can help hospitals to curtail opioid over-prescription and reduce opioid leftovers, without increasing the burden of hospitals and patients.

Chapter 8

Improving Care Quality in Primary Care Clinics

8.1 Introduction

To design and evaluate workflow models in healthcare system, this chapter introduces terminating Markov chain models for physician workflow management in primary care clinics. The physician workload is characterized by face-to-face encounters with patients and documentation of electronic health record (EHR) data. Focusing on the two types of tasks, three workflow management policies are considered: preemptive priority (stop ongoing documentation tasks if a new patient arrives); non-preemptive priority (finish ongoing documentation even if a new patient arrives); and batch documentation (start and finish documentation when the desired number of tasks is reached). Analytical formulas are derived to quantify the performance of three management policies, such as physician's daily working time, patient's waiting time, and documentation waiting time. A comparison of the results under three policies is carried out. Finally, a case study in a primary care clinic is carried out to illustrate model applicability. Such a work provides a quantitative tool for primary care physicians to design and manage their workflow to reduce burnout

and improve care quality.

8.2 System Description and Problem Formulation

Physicians in primary care clinics are facing mainly two types of work everyday, clinical meetings with patients and computer documentations, and may have different preferences on workflow management. In many cases, face-to-face encounters with patients have a higher priority than EHR documentation work. Thus, after serving one patient and leaving the associated documentation task in a document queue, a physician will continue to meet with the patients waiting in the patient queue. Only when no patient is waiting, documentation work will be started. Moreover, he/she may stop EHR work to meet with a newly arrived patient. However, such interruptions may lead to loss of information for previous patients after serving multiple ones, which will influence EHR accuracy. To mitigate this impact, some physicians may prefer non-interruption of ongoing documentation. In a crowded clinic, the documentation queue can become relatively long, which again may impact documentation quality. Thus, some physicians may batch process EHR related documentation if the document queue is accumulated to a certain number. In other words, a physician may not start working on documentation until the queue reaches a threshold, and the documentation task will not be interrupted until the batch is finished.

Based on the aforementioned discussions, we formulate three kinds of workflow management strategies or policies. The objective of this study is to evaluate the performances of these policies under different settings and parameters, including patient arrival rate, physician service and documentation rates. To achieve this, we adopt terminating Markov chain models. The following assumptions are introduced.

- (i) A physician has two kinds of tasks: meeting with patients and documentation, where their service times follow general distributions with rates μ_p and μ_d and coefficients of variation (CVs) cv_p and cv_d , respectively.

(ii) There are N patients visiting the physician per day. The inter-arrival process can be described by a general distribution with rate λ and CV cv_a .

(iii) There are three workflow management policies:

- Preemptive priority policy (PEP): After finishing the service to one patient, the physician will continue to meet with the next one if another is waiting. Otherwise, he/she will work on a documentation task. However, if a new patient arrives during documentation, it will be stopped and the patient will be met first.
- Non-preemptive priority policy (NPP): Again the physician will meet with the patient in waiting first. Otherwise, he/she will work on a documentation task and finish it, even if a new patient comes during the process.
- Batch documentation policy (BDC): The physician will not start documentation until the number of documenting tasks in waiting reaches a predefined threshold M , and the physician will finish all M tasks before meeting with new patients. In the case of last batch less than M , all remaining tasks will be processed in a smaller batch.

(iv) All patients must be served and all documentations must be finished everyday. If the total time is longer than the scheduled shift time, the physician will work overtime.

Remark 8.1. *Note that the three policies (PEP, NPP, and BDC) are summarized and abstracted based on observations of physician practice and discussion with physicians and support staff. They are introduced to characterize different working styles. The PEP model is used to describe the scenario that meeting with patients has the highest priority, and the NPP model intends to represent the style of non-interrupted documentation work, while batch processing of documentations is illustrated in the BDC model. In practice, we think no physicians will consistently and exactly follow one policy only. However, we hope that this study can conceptually evaluate and compare the performances of different workflow policies and provide useful insights.*

Under assumptions (i)-(iv), the problem addressed in this chapter is to develop a method to evaluate system performance under three workflow management policies and investigate system properties.

Remark 8.2. *To start the study, we will assume exponential service and inter-arrival times first. Then, using the exponential model results, we will extend to general distribution through approximations.*

8.3 Workflow Modeling

In this section, we formally define the state space and state transitions in workflow management models.

8.3.1 Preemptive Priority (PEP) Model

A physician's daily work can be viewed as a terminating process with a fixed number of patients visiting him/her everyday. To fully characterize the system, we need to know how many patients have been served, how many are currently in the system (being served or waiting), and how many will arrive; as well as how many documentation tasks have been finished, how many are still in the pipeline (being processed or waiting), and how many will come. Therefore, the states of the system can be defined as:

$$S^k = (X_p^k, X_d^k, X_c^k), \quad k = 1, 2, \dots, K,$$

where $K = (N + 1)^3$ characterizes system dimension. In addition, X_p^k is the number of patients currently in the system in state k , which includes the patients waiting in the system and the patient being served now; X_d^k is the number of unfinished documentations (including the in-process one and those waiting to be processed) in the system in state k ; and finally, X_c^k is the number of completed documentations in state k . For example,

state $(3,1,3)$ indicates that there are 3 patients currently in the system, 1 documentation needs to be finished, and the physician has completed 3 documentations already. Note that state $(0, 0, N)$ is the terminating state, which implies all patients have been served and all paperwork has been finished.

However, not all states are feasible. For instance, the total number of patients per day N defines the maximal number of documentation tasks, which includes the finished documentations X_c^k (belong to patients who have been served already), the documentations waiting to be finished X_d^k (also belong to patients being served already), and the number of incoming documentations (which should be larger or equal to the number of patients in the system X_p^k). The sum of the number of finished and unfinished documents and the number of patients in the system should not exceed the total number of patients per day. Therefore, the state space is limited by the following constraints:

- (a) $X_p^k + X_d^k + X_c^k \leq N$: total number of tasks constraint.
- (b) $X_p^k, X_d^k, X_c^k \geq 0$: positive number of tasks constraint.

Only the states satisfying these constraints are feasible. Then, the number of feasible states is:

Lemma 8.1. *The number of feasible states in the PEP model can be calculated by*

$$K^{\text{PEP}} = \frac{1}{6}(N + 1)(N + 2)(N + 3). \quad (8.1)$$

Proof: See Appendix A. ■

For a feasible state S^k , the state transition from state S^k to another state $S^{k'}$ can be triggered by one of the following events:

Transition 1: A patient arrives and the physician is or will be meeting with a patient. In this case, X_p is increased by 1 with transition rate λ .

Transition 2: The physician finishes serving a patient. Then, X_p is decreased by 1 and X_d is increased by 1, while the transition rate is μ_p .

Transition 3: The physician finishes a documentation task. Then, X_d is decreased by 1, while X_c is increased by 1, with transition rate μ_d .

According to the PEP policy, a physician can work on a documentation task only when no patient is waiting or arriving. The conditions based on which these transitions occur are listed in Appendix A.

8.3.2 Non-preemptive Priority (NPP) Model

In this model, a physician may continue to work on an ongoing documentation task even if a new patient comes in. This can happen regardless of whether there is a patient waiting. Thus, different from the state definition in the PEP model, an additional state variable is needed in order to indicate whether a documentation task is processed or not when a patient is waiting. Let $S^k = (X_p^k, X_d^k, X_c^k, D_d^k)$, where $D_d^k \in \{0, 1\}$, and $D_d^k = 1$ indicates that a documentation task is in progress even though there is a waiting patient. For example, state $(1,1,3,1)$ implies that there is one patient in the system, one document to be finished, and three documentation tasks completed, but the physician is working on a documentation task currently (so that the patient is waiting). Again, state $(0, 0, N, 0)$ is a terminating state.

According to the NPP policy, in addition to the total number of tasks constraint (a) and the positive number of tasks constraint (b), the state space is further constrained by:

(c) $D_d^k \leq X_d^k$: constraint of documentation indicator, thus if $X_d^k = 0$, then $D_d^k = 0$.

(d) $D_d^k \leq X_p^k$: constraint of patient indicator, thus if $X_p^k = 0$, then $D_d^k = 0$.

These two constraints suggest that $D_d^k = 1$ only when there exists a patient in the system and the physician is working on a document. Then we obtain

Lemma 8.2. *The number of feasible states in the NPP model can be calculated by*

$$K^{\text{NPP}} = \frac{1}{3}(N + 1)(N^2 + 2N + 3). \quad (8.2)$$

Proof: See Appendix A. ■

In the NPP model, Transitions 1 and 2 are the same as in the PEP model. Transition 3 still holds, and D_d will be decreased to 0 if $D_d = 1$ before. In addition, state S^k can be transited to state $S^{k'}$ by the following event:

Transition 4: During a documentation process, a new patient arrives before the task is finished. In this case, X_p is increased by 1 and D_d will become or remain at 1, with transition rate λ .

The conditions based on which these transitions occur are listed in Appendix A.

8.3.3 Batch Documentation (BDC) Model

The state definition is the same with that in the PEP model, i.e., $S^k = (X_p^k, X_d^k, X_c^k)$. State $(0, 0, N)$ is still a terminating state. However, the state space is limited by more constraints in addition to the total and positive number of tasks constraints (a) and (b):

- (e) $X_d^k \leq M$: documentation queue length constraint, where M is the batch size.
- (f) $((X_d^k + X_c^k) \bmod M) = 0$, if $X_p^k + X_d^k + X_c^k < 0$ and $(X_c^k \bmod M) > 0$: documentation batch work constraint.

The last constraint indicates that a batch of documentation tasks is finished only after the last document is processed. Then it follows

Lemma 8.3. *The number of feasible states in the BDC model can be calculated by*

$$K^{\text{BDC}} = \frac{1}{2}MA(4N - M - 2MA + 5) + \frac{1}{2}(B^2 + 5B + 2), \quad (8.3)$$

where

$$A = \left\lfloor \frac{N}{M} \right\rfloor, \quad B = N - AM.$$

Proof: See Appendix A. ■

The transition from state S^k to state $S^{k'}$ can still be triggered by the similar events in the PEP model. However, the occurrence of transitions is related to the remainder of modular operation “mod”, which characterizes whether the batch processing is finished or not. The feature of batch processing indicates that $(X_c^k \bmod M)$ will be greater than zero and less than M after the first documentation task in a batch is carried out until the M -th one in the same batch is finished. The conditions based on which these transitions occur are listed in Appendix A.

Remark 8.3. *In the BDC model, even when there is no patient in the system, as long as more patients can arrive (i.e., the number of patients served is less than N), the physician will not work on documentation until the batch is fully filled with M tasks. In other words, the physician starts working on documentation only when there are M waiting ones. Thus, the physician may observe more idle time in the BDC model. The only exception is the last batch whose number of documentations could be smaller than M . After the physician has served all the scheduled patients for the day, he/she will work on documentations.*

Remark 8.4. *The three policies are represented by finite-state Markov chains. In each Markov chain, there is only one recurrent state, which corresponds to the terminating state, $(0, 0, N)$ in the PEP*

and BDC models and $(0, 0, N, 0)$ in the NPP model. All the other states are transient.

8.4 Performance Analysis

8.4.1 Solution Approach

In order to compare the efficacy of three workflow management policies, we select the following performance measures under policy $x \in \{\text{PEP}, \text{NPP}, \text{BDC}\}$:

- the physician's average total working time per day to finish all the tasks, T^x ;
- the average patient's waiting time, W^x ; and
- the average waiting time for documentation, Q^x .

The importance of the first two measures is evident, while the third one is related to quality of care [159]. All of the performance measures are functions of arrival and service parameters λ, μ_p, μ_d , as well as number of patients N , or batch size M , under policy x , $x \in \{\text{PEP}, \text{NPP}, \text{BDC}\}$, i.e.,

$$T^x = f_t^x(\lambda, \mu_p, \mu_d, N, M),$$

$$W^x = f_w^x(\lambda, \mu_p, \mu_d, N, M),$$

$$Q^x = f_q^x(\lambda, \mu_p, \mu_d, N, M).$$

All three policy models contain an absorbing state in which a physician finishes all the work and no further transition will occur. Therefore, performance measures cannot be obtained using steady state analysis. However, all three models have an intrinsic property that all states are visited at most once, thus the Markov chain will never return to the states they have visited before. From this property, the performance measures are obtained using the sojourn time and the probability of visiting, which are derived below.

Let $\eta(S^i, S^j)$ be the instantaneous transition rate, i.e., the rate of process transition from state S^i into state S^j . Then, the transition rates corresponding to Transitions 1-3 (or 1-4) in PEP, NPP, and BDC models can be derived (see Appendix A). Using these rates, the sojourn time, which is the amount of time spent in a state before making a transition into another state, can be derived for state S^i . Denote

$$v(S^i) = \sum_{\forall j \neq i} \eta(S^i, S^j).$$

When the sojourn time in state S^i is exponentially distributed with rate $v(S^i)$, it implies that the average sojourn time in state S^i is

$$T(S^i) = \frac{1}{v(S^i)}. \quad (8.4)$$

The visiting probability defines the probability that a state is visited during the process. To derive this probability, let $p(S^i, S^j)$ be the transition probability from state S^i to S^j where $i \neq j$. Then,

$$p(S^i, S^j) = \frac{\eta(S^i, S^j)}{v(S^i)} = \frac{\eta(S^i, S^j)}{\sum_{\forall j, j \neq i} \eta(S^i, S^j)}.$$

With all the transition probabilities and initial state probabilities, we can obtain the visiting probabilities. Denote $P \in \mathbb{R}^{K^x \times K^x}$, $x \in \{\text{PEP}, \text{NPP}, \text{BDC}\}$, as the matrix of transition probabilities with elements

$$P_{ij} = p(S^i, S^j),$$

and introduce $I \in \mathbb{R}^{K^x}$ as the vector of initial probabilities where only the value corresponding to state $S = (0, 0, 0)$ or $S = (0, 0, 0, 0)$ is 1 and all others are 0. Then, define $V \in \mathbb{R}^{K^x}$ as the vector of visiting probabilities $V(S^i)$. Since the i -th elements of $I^\top P^k$ and V are the probability a process enters state S^i after k transitions, and the probability state S^i is visited during the entire process, respectively, and the total number of transitions occurring in one

process is exactly $3N$ (i.e., N patient arrivals, N patient services, and N documentations), we obtain

$$V^{\top} = \sum_{k=0}^{3N} I^{\top} P^k,$$

where superscript “ \top ” indicates transpose.

Using the sojourn time and visiting probability, we can obtain analytical formulas for all three performance measures, T^x , W^x , and Q^x .

8.4.2 Analysis Formulas

Let m , n , l and w represent the numbers of X_p , X_d , X_c and D_d , respectively. Given state $S = (m, n, l)$ in PEP and BDC models or $S = (m, n, l, w)$ in the NPP model, rate $v(S)$ can be obtained by considering the sum of rates of possible transitions. Then sojourn time $T(S)$ can be calculated from (8.4). In addition, visiting probability $V(S)$ of state S can be derived recursively from previous states (see expressions (8.5)-(8.10)).

Let $S = (m, n, l)$ in PEP and BDC models and $S = (m, n, l, w)$ in NPP model.

$$v(S) = \begin{cases} \lambda\delta(m+n+l < N) + \mu_p\delta(m > 0) + \mu_d\delta(m=0, n > 0), & \text{PEPmodel,} \\ \lambda\delta(m+n+l < N) + \mu_p(1-w)\delta(m > 0) + w\mu_d \\ \quad + \mu_d(1-w)\delta(m=0, n > 0), & \text{NPPmodel,} \\ \lambda\delta(m+n+l < N) + \mu_p\delta(m > 0, n < M, (l \bmod M) = 0) \\ \quad + \mu_d\delta((l \bmod M) \neq 0 \text{ or } l = M \text{ or } (n+l = N, n > 0)), & \text{BDCmodel,} \end{cases} \quad (8.5)$$

For $m, n, l \in \mathbb{Z}^+$ and $0 < m+n+l \leq N$,

$$V(S) = \begin{cases} \frac{\lambda V((m-1, n, l))}{\lambda + \alpha} + \frac{\mu_p V((m+1, n-1, l))}{\mu_p + \beta} + \frac{\mu_d \delta(m=0) V((m, n+1, l-1))}{\mu_d + \gamma}, & \text{PEP model,} \\ \frac{\lambda [1 - \delta(m=1, n > 0)] V((m-1, n, l, w))}{\lambda + \alpha} + \frac{\mu_p (1-w) V((m+1, n-1, l, w))}{\mu_p + \beta} \\ + \frac{\mu_d (1-w) \delta(m=0) V((m, n+1, l-1, w))}{\mu_d + \gamma} + \frac{w \mu_d V((m, n+1, l-1, w))}{\mu_d + \gamma} \\ + \frac{\mu_d (1-w) V((m, n+1, l-1, w+1))}{\mu_d + \gamma} + \frac{\lambda w \delta(m=1) V((m-1, n, l, w-1))}{\lambda + \mu_d}, & \text{NPP model,} \\ \frac{\lambda V((m-1, n, l))}{\lambda + \alpha} + \frac{\mu_p \delta((l \bmod M)=0, n \leq M) V((m+1, n-1, l))}{\mu_p + \beta} \\ + \frac{\mu_d \delta((n+l \bmod M)=0 \text{ or } n+l=N) V((m, n+1, l-1))}{\mu_d + \gamma}, & \text{BDC model,} \end{cases} \quad (8.6)$$

where

$$\delta(A) = \begin{cases} 1, & \text{if } A \text{ is true,} \\ 0, & \text{otherwise,} \end{cases} \quad \beta = \begin{cases} 0, & \text{if } m + n + l = N, \\ \lambda, & \text{otherwise,} \end{cases} \quad \gamma = \begin{cases} 0, & \text{if } m + n + l = N, \\ \lambda, & \text{otherwise,} \end{cases} \quad (8.7)$$

$$\alpha = \begin{cases} 0, & \text{if } \begin{cases} m = 1, n = 0, & \text{PEP and NPP models,} \\ (l \bmod M) = 0, m = 1, n < M, & \text{BDC model,} \end{cases} \\ \mu_d, & \text{if } \begin{cases} m = 1, n \neq 0, & \text{PEP model,} \\ m = 1, n \neq 0 \text{ or } w = 1, & \text{NPP model,} \\ (n + l \bmod M) = 0, n \neq 0, & \text{BDC model,} \end{cases} \\ \mu_p, & \text{otherwise,} \end{cases} \quad (8.8)$$

with initial conditions

$$V((0, 0, 0)) = 1 \quad (\text{PEP and BDC models}) \quad \text{or} \quad V((0, 0, 0, 0)) = 1 \quad (\text{NPP model}), \quad (8.9)$$

and boundary condition

$$V(S) = 0, \quad \text{if } \min(m, n, l, w) < 0 \text{ or } m + n + l > N, \quad (8.10)$$

Equations (8.5)-(8.10) imply that the visiting probability of state S is a weighted sum of visiting probabilities of other states. For instance, in the PEP model, the visiting probability of state (m, n, l) is a combination of visiting probabilities of states $(m-1, n, l)$, $(m+1, n-1, l)$ and $(m, n+1, l-1)$, weighted by the occurring probabilities of Transitions 1, 2 and 3, respectively. Note that Transition 3 in the PEP model occurs only when there is no patient, so an indicator function is multiplied in the last term.

Using these results, the performance measures of PEP, NPP, and BDC policies can be obtained (see formulas (8.11)-(8.13)).

$$T^x = \sum_{\text{State } S} T(S)V(S), \quad x = \text{PEP, NPP, BDC}, \quad (8.11)$$

$$W^x = \begin{cases} \sum_{\substack{S=(m,n,l) \\ m>1}} \frac{(m-1)T(S)V(S)}{N}, & x = \text{PEP}, \\ \sum_{\substack{S=(m,n,l,w) \\ m>1, w=0}} \frac{(m-1)T(S)V(S)}{N} + \sum_{\substack{S=(m,n,l,w) \\ w=1}} \frac{mT(S)V(S)}{N}, & x = \text{NPP}, \\ \sum_{S=(m,n,l)} \frac{mT(S)V(S)}{N} - \sum_{(l \bmod M)=0, m>0, n<M} \frac{T(S)V(S)}{N}, & x = \text{BDC}, \end{cases} \quad (8.12)$$

$$Q^x = \begin{cases} \sum_{\substack{S=(m,n,l) \\ m>0}} \frac{nT(S)V(S)}{N} + \sum_{\substack{S=(m,n,l) \\ m=0,n>1}} \frac{(n-1)T(S)V(S)}{N}, & x = \text{PEP}, \\ \sum_{\substack{S=(m,n,l,w) \\ m>0,w=0}} \frac{nT(S)V(S)}{N} + \sum_{\substack{S=(m,n,l,w) \\ w=1 \text{ or} \\ m=0,n>0,w=0}} \frac{(n-1)T(S)V(S)}{N}, & x = \text{NPP}, \\ \sum_{S=(m,n,l)} \frac{nT(S)V(S)}{N} - \sum_{\substack{S=(m,n,l) \\ (l \bmod M) \neq 0}} \frac{T(S)V(S)}{N} \\ - \sum_{\substack{S=(m,n,l) \\ (l \bmod M)=0,n=M}} \frac{T(S)V(S)}{N} - \sum_{\substack{S=(m,n,l) \\ (l \bmod M)=0,m=0 \\ n<M,m+n+l=N}} \frac{T(S)V(S)}{N}, & x = \text{BDC}. \end{cases} \quad (8.13)$$

- In all models, the average total working time per day T in equation (8.11) can be obtained by a weighted sum of average sojourn time of all states, where the weights are defined by visiting probabilities.
- In the PEP model, the average patient waiting time W in equation (8.12) can be calculated by summing up the multiplication of average sojourn time, visiting probability, and number of waiting patient for all states that have waiting patients (i.e. $X_p > 1$), and then normalizing by the maximum number of patients N . Meanwhile, in the NPP model, a waiting patient will be served after the previous patient's visit is finished or a documentation task is done (when $X_p > 0$ and $D_d = 0$). However, if a new patient comes when the physician is still working on a documentation task (i.e., $D_d = 1$), he/she needs to wait until the documentation is finished. In the BDC model, a patient is served when $(X_c \bmod M) = 0$, $X_p > 0$, and $X_d < M$. Thus, the average patient waiting time excludes this scenario.
- To obtain the average documentation waiting time Q in equation (8.13), in the PEP model, we need to take into account all waiting documentation tasks in a state when a patient is being served, since a documentation task can only be carried out when there is no patient in the system. In the NPP model, when a patient is in service, all documentation tasks need to wait. If there is no patient in the system (i.e., $X_p = 0$), or

if the physician is working on a documentation task (i.e., $D_d = 1$), the documentation waiting time does not include the task that is in process. Finally, in the BDC model, the average documentation waiting time can be obtained by taking into account all documentation tasks and subtracting the ongoing ones.

Remark 8.5. *In the analysis, the three proposed performance measures T^x , W^x , and Q^x are calculated by the weighed sums of sojourn times. As mentioned earlier, the sojourn time in each state follows exponential distribution. Thus, the performance measures follow mixture of exponential distributions.*

Remark 8.6. *The average patient's waiting time W^x aggregates the waiting times of all patient visited in a day and does not consider individual values. However, in some cases, it is important to prevent a patient from waiting longer than a certain limit. Then, one should analyze the maximum of average waiting time among the visiting patients, i.e., $\max_{i \in \{1, \dots, N\}} W_i^x$ where W_i^x is the i -th patient's expected waiting time. As patients are served by 'first come, first served' basis and each state can tell which patients are waiting for services, each individual patient's waiting time is obtained by the following equations.*

$$W_i^x = \begin{cases} \sum_{\substack{S=(m,n,l) \\ m+n+l \geq i, i > n+l+1}} T(S)V(S), & x = \text{PEP}, \\ \sum_{\substack{S=(m,n,l,w) \\ m+n+l \geq i, i > n+l+1, \\ w=0}} T(S)V(S) + \sum_{\substack{S=(m,n,l,w) \\ m+n+l \geq i, i \geq n+l+1, \\ w=1}} T(S)V(S), & x = \text{NPP}, \\ \sum_{\substack{S=(m,n,l,w) \\ m+n+l \geq i, i \geq n+l+1}} T(S)V(S) - \sum_{\substack{S=(m,n,l) \\ m+n+l \geq i, i = n+l+1, \\ (l \bmod M) = 0, n < M}} T(S)V(S), & x = \text{BDC}, \end{cases} \quad (8.14)$$

In addition, note that the first patient's waiting time is always 0 in all workflow policies.

8.4.3 System Properties and Comparisons

To understand the efficacy and insights of workflow models, we investigate the system properties and compare their performances under different parameter settings.

8.4.3.1 Average daily working time

Comparing the results in all three models, as shown in Figure 8.1(a), the average daily working times are the same in PEP and NPP models (solid/dot lines and star lines, respectively), but longer in the BDC model (circle lines). Since the difference between PEP and NPP models is whether the physician will stop documentation immediately or not if a new patient arrives, the sojourn times are the same in both models. In the BDC model, the physician will not start documentation work if the number of tasks is not accumulated to the batch size, even if there is no patient in the clinic. This makes the working time longer. In all three models, the daily working time T^x , $x=PEP, NPP, BDC$, is monotonically decreasing with respect to patient arrival rate λ and the decreasing speed is higher when patient service rate μ_p or documentation service rate μ_d is larger. Since the total number of patients N is a constant, when patients arrive more intensively, the physician will have less idle time, and the total working hours will be reduced. Similarly, T^x decreases when patient service rate μ_p increases. Again, higher documentation service rate μ_d will result in shorter T^x , i.e., faster documentation leads to shorter working time.

Through extensive numerical experiments to compare T^{PEP} , T^{NPP} and T^{BDC} , we obtain

Numerical Fact 1. *Under assumptions (i)-(iv),*

$$T^{PEP} = T^{NPP} \leq T^{BDC}, \quad (8.15)$$

where the inequality holds as equality when $M = 1$.

Remark 8.7. *Note when $M = 1$, the BDC model implies the physician will start documentation immediately after meeting with a patient no matter whether there are other patients waiting or not.*

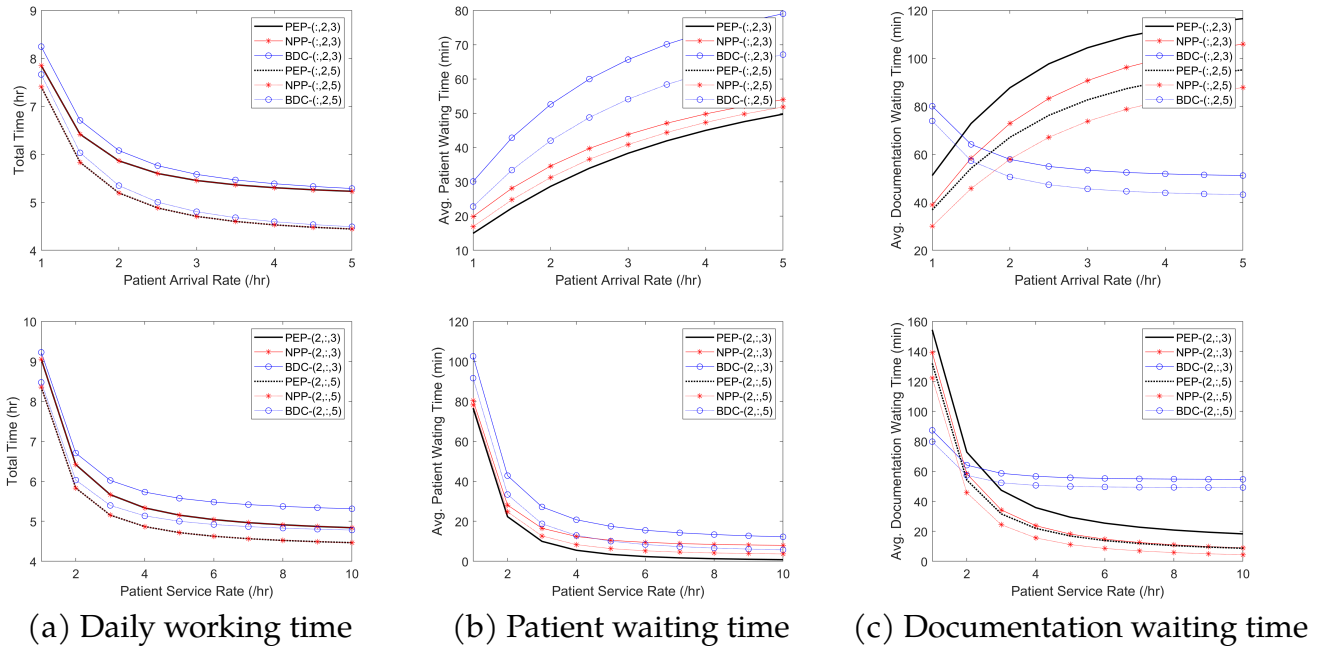


Figure 8.1: Performance comparison in all models

In this case, the only idle time is when no patient or documentation exists in the system. Thus, the working time should be the same in all three models.

8.4.3.2 Average patient waiting time

As shown in Figure 8.1(b), the average patient waiting time is the shortest in the PEP model (solid/dot lines) and the longest in the BDC model (circle lines), while the NPP model (star lines) is in the middle. This is because, in the PEP model, serving a patient has the highest priority so that documentation service rate does not have any impact, while in NPP and BDC models, a patient could wait for one or M documentation times, respectively.

In addition, it is observed that the average patient waiting times, W^{PEP} and W^{NPP} , both become longer when arrival rate λ increases. However, they become shorter if patient service rate μ_p increases. This implies that more intensive arrivals lead to longer patient waiting time. Also, the faster services, the less waiting time. Although W^{PEP} does not change with respect to μ_d . W^{NPP} and W^{BDC} are dependent on μ_d . In other words, faster documentation work does not affect patient waiting time in the PEP model, but leads to

less waiting in NPP and BDC models. This is due to the fact that a patient has a higher priority than documentation in the PEP model so that no patient will ever wait because of documentation task, while waiting during a documentation service is possible in the other two models.

Comparing W^x with the same parameters $N, M, \lambda, \mu_p, \mu_d$, we have

Numerical Fact 2. *Under assumptions (i)-(iv),*

$$W^{\text{PEP}} \leq W^{\text{NPP}} \leq W^{\text{BDC}}. \quad (8.16)$$

Note that the inequality holds as equality in a few extreme cases. For instance, if every patient enters the system after the previous patient's service and documentation and $M = 1$, then $W^{\text{PEP}} = W^{\text{NPP}} = W^{\text{BDC}} = 0$.

8.4.3.3 Average documentation waiting time

Unlike daily working time and patient waiting time that have similar trends among all models, we can find different patterns in documentation waiting time among the three models. As one can see from Figure 8.1(c), which compares the average documentation times under three policies, Q^{PEP} (solid/dot lines) and Q^{NPP} (star lines) have similar trends but Q^{NPP} is shorter, which is because the physician will finish an ongoing documentation first. Thus,

Numerical Fact 3. *Under assumptions (i)-(iv),*

$$Q^{\text{PEP}} \geq Q^{\text{NPP}}. \quad (8.17)$$

Again, only in extreme cases where no documentation is interrupted, the inequality holds as equality.

In addition, it can be shown that Q^{PEP} and Q^{NPP} increase with respect to λ and decrease in μ_p . Thus, when patients arrive more intensively, more documentation tasks pile up and longer waiting time is expected. However, Q^{BDC} has an opposite pattern with respect to patient arrival rate. When patients arrive more frequently, the waiting time is reduced since the number of accumulated documentation tasks will increase faster and the physician will start to work on documentations earlier. Thus, we can find a decreasing pattern of Q^{BDC} with respect to arrival rate. If the service time is reduced, less waiting will be observed in all three models. Moreover, if μ_d becomes larger, i.e., faster processing of documents, then waiting time is reduced in all models.

Note that the BDC model is introduced to characterize the batch processing style used by some physicians. Using such a model, we can evaluate the differences in performance measures comparing with other models, and seek appropriate thresholds to achieve a trade off between performance shortage and documentation interruptions.

Figure 8.1 is drawn based on $N = 6$ (representing a half-day or less crowded flow). In addition, we also consider the case of $N = 14$ that characterizes a full-day or more crowded scenario. The behaviors are much similar to those when $N = 6$.

8.5 Extensions

The performance measures introduced above are obtained under the exponential assumption of patient inter-arrival time, service time and documentation time. However, in practice, they may not follow exponential distributions. Thus, developing a method to study the non-exponential scenarios is of importance.

To consider non-exponential cases, two questions need to be answered. First, are the performance measures dependent on distribution types or not? If the answer is yes, then we may need to develop a method for each distribution type. If the answer is no, then is it possible to develop an empirical formula independent of distributions type? To answer

them, we use simulations to evaluate the outcomes of non-exponential distributions. Three widely used distribution types, gamma, Weibull, and lognormal, are considered since they all have two parameters, which provide more freedom for selecting mean and CV values.

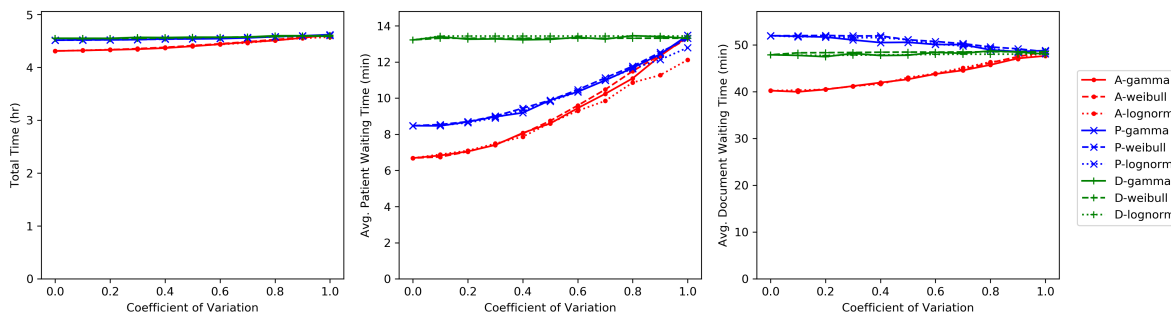
In simulation experiments, we assume $CV < 1$ for patient inter-arrival, physician service and documentation times. As shown in [160], if the service rate is an increasing function of time, then the CV of service time will be less than 1. In other words, the longer a service has been carried out, the higher probability the service will finish. In primary care clinics, many services (such as physical exams, well child check, typical disease diagnosis, and follow-up visits) are standard procedures, thus the probability to finish a service is increasing with the time of service that has been carried out, i.e., the service rate increases with time. The documentation process has a similar property. Since most of the patients come based on schedules (although with variations), the arrival rate is also increasing with time, i.e., the longer time from last arrival, the higher probability the next patient will come. This again leads to the CV of inter-arrival time less than 1. The data analyses in a number of primary care and specialty clinics (for instance, [113,116,161]) all verify such a property. Thus, the CVs are, most likely, smaller than 1, and the simulation experiments and follow-up empirical formulas are designed for this scenario.

Specifically, we consider CVs increasing from 0 to 1 with 0.1 increments. Using Python and package *SimPy*, the simulation models of PEP, NPP, and BDC policies are developed. The simulations are executed 10,000 times of physician's daily work to obtain the sample means of physician working time, patient waiting time and documentation waiting time.

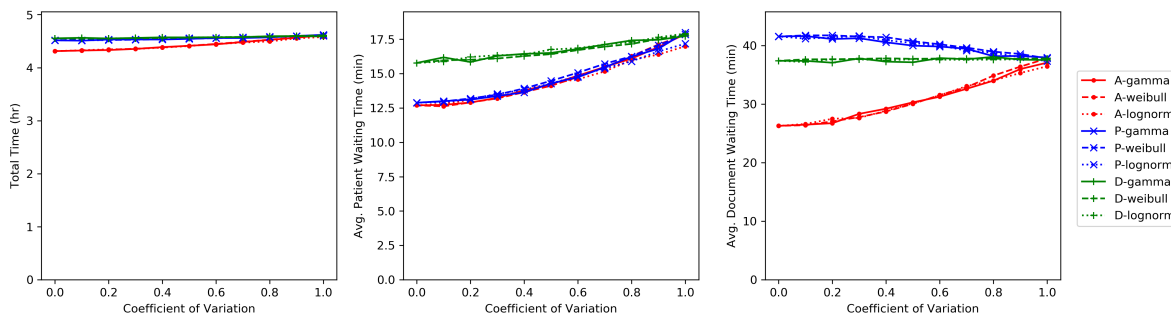
8.5.1 Dependence on Distribution Type

To investigate whether the performance measures are dependent on distribution types, the patient inter-arrival and service times, and documentation time with different distributions and CVs are generated. Figure 8.2 illustrates the simulation results when $\lambda = 2$, $\mu_p = 3$, $\mu_d = 4$ and $N = 6$, in PEP, NPP, and BDC models. The lines represent performance

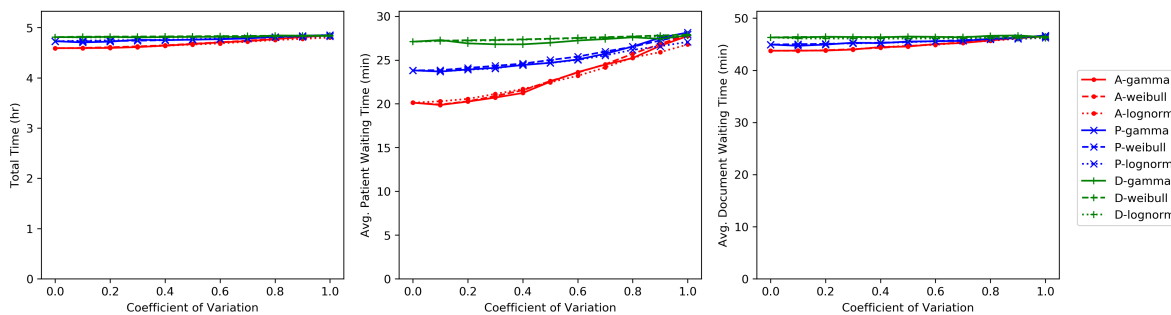
measures with respect to CVs (from 0 to 1) of patient inter-arrival time (denoted as “A” in legends), service time (denoted as “P”), and documentation time (denoted as “D”) when they follow gamma, Weibull, and lognormal distributions. Examining these results, for each rate (A, P, and D), we observe that the performance measures are practically the same for different distributions when CVs are identical.



(a) PEP model



(b) NPP model



(c) BDC model

Figure 8.2: Performance measures with non-exponential distributions

To further analyze the differences in performance measures, let ζ_z be the maximum relative difference in performance z , $z \in \{T, W, Q\}$, under distribution g , $g \in \{\text{gamma, Weibull, lognormal}\}$,

and CVs between 0 and 1, i.e.,

$$\zeta_z = \max_{cv} \left(\frac{\max_g z_{g,cv} - \min_g z_{g,cv}}{z_{exp}} \right) \cdot 100\%,$$

$$g \in \{\text{gamma, Weibull, lognormal}\}, \quad cv \in \{0, 0.1, 0.2, \dots, 1\},$$

where z_{exp} represents the corresponding measure in exponential case.

The results in PEP, NPP, and BDC models with different rates of λ , μ_p , and μ_d are examined. It is shown that the maximal differences are very small with $\zeta_z \leq 1.22\%$. This implies that the system performance is primarily dependent on the first and second moments (i.e., mean and CV) rather than the complete distribution. Such a result is consistent with observations in other healthcare systems studies [114,116,161] as well as in manufacturing systems [162].

8.5.2 Empirical Formulas

As shown in Figure 8.2, the results suggest that the performance measures exhibit close to quadratic functions of CV in most cases. However, the impacts of variations are different. Thus, empirical formulas are proposed for each performance measure under a given policy. Since $CV = 1$ stands for exponential distribution, performance measures when all CVs equal to 1 are obtained using the analytical formulas derived in Section 8.4.1 under exponential assumptions. When all CVs are 0, it implies no variation exists in any processes so the performance measures can be derived directly. Using these points, the empirical solutions with respect to CVs between 0 and 1 are developed.

8.5.2.1 Daily working time

From Figure 8.2, it is inferred that a physician's daily working time slightly increases with respect to CVs, while the increasing rates are proportional to their reciprocals, and the growth rate is higher for CV of inter-arrival time. This is likely because the variation of

inter-arrival time can directly affect the physician's work ending time. For example, if a patient arrives late, then the physician will start service and documentation late, which delays work ending time. Meanwhile, it is possible that the service and documentation can be finished before the next patient's arrival. In this case, an increase of patient service or documentation time may not affect work ending time. Thus, the daily working time should be a function of the CVs.

In addition, empirical formulas are found in other studies of similar primary and specialty care systems (e.g., [113, 114, 116]) to approximate the performance measures in non-exponential cases by a straight line between the exponential performance and the non-variation one (i.e., $CV = 0$) with the coefficient of squared CVs. Similarly, in the approximation formulas for a G/G/1 model, the sum of the squared CVs of inter-arrival time and service time are used to adjust the performance of an M/M/1 queue. Therefore, we develop the empirical formula to consider a linear function between exponential and non-variation cases with a coefficient of a summation of squared CVs of inter-arrival, patient service, and documentation service times, weighted by the corresponding rates to account them equally. A similar rationale is applied to waiting time measures.

Specifically, the approximation formula is introduced as follows:

$$\tilde{T}^x = T_0^x + \frac{(T^x - T_0^x) \left(\frac{2cv_a^2}{\lambda} + \frac{cv_p^2}{\mu_p} + \frac{cv_d^2}{\mu_d} \right)}{\frac{2}{\lambda} + \frac{1}{\mu_p} + \frac{1}{\mu_d}}, \quad x \in \{\text{PEP}, \text{NPP}, \text{BDC}\}, \quad (8.18)$$

where T^x and T_0^x denote the total working times when $CV = 1$ (exponential) and $CV = 0$ (no variation), respectively, and \tilde{T}^x indicates non-exponential ($0 < CV < 1$) case. The algorithm to calculate T_0^x is presented in Appendix B.

To evaluate the accuracy of the empirical formula, define ϵ^x as the difference between the simulated working time, $\tilde{T}^{\text{sim},x}$, and the calculated working time using (8.18), \tilde{T}^x , for policy x , $x \in \{\text{PEP}, \text{NPP}, \text{BDC}\}$, and let $\tilde{\epsilon}^x$ be the normalized difference by dividing the

total working time under exponential assumption, T^x .

$$\begin{aligned}\epsilon^x &= | \tilde{T}^{\text{sim},x} - \tilde{T}^x |, \\ \tilde{\epsilon}^x &= \frac{\epsilon^x}{T^x} \cdot 100\%.\end{aligned}\tag{8.19}$$

Then the expected difference $E[\epsilon^x]$ and the expected normalized difference $E[\tilde{\epsilon}^x]$ for each parameter setting of N , λ , μ_p , and μ_d are estimated by the averages of ϵ^x and $\tilde{\epsilon}^x$ resulting from 1,000 experiments with varying values of cv_a , cv_p and cv_d from 0 to 1 with 0.1 increments. The mean and maximal values of ϵ^x and $\tilde{\epsilon}^x$ of the experimental results using different parameter setting for each model are presented in Table 8.1. In these experiments, different combinations of λ , μ_p , and μ_d are used and N is fixed as 6. The average differences are around 0.06 to 0.08 (i.e., only 4-5 minutes) and the largest one is about 10 minutes, while the average relative error is less than 2% and the largest is within 4%. These indicate that empirical formula (8.18) can provide accurate estimates of physician daily working time in non-exponential environment.

Table 8.1: Accuracy of total working time formula

x	$E[\epsilon^x]$		$E[\tilde{\epsilon}^x]$	
	Mean (hr)	Max (hr)	Mean (%)	Max (%)
PEP	0.08	0.17	1.86	3.91
NPP	0.08	0.17	1.85	3.91
BDC	0.06	0.09	1.13	1.92

8.5.2.2 Patient waiting time

In the PEP model, it is shown that patient waiting time is not dependent on documentation rate. In fact, a patient in the system will be served first no matter there exists a documentation work or not. Thus, variations in documentation time will not affect patient waiting time. Thus, cv_d should not be included in patient waiting time calculation in the PEP model.

The patient waiting time in the NPP model increases quadratically as any CV increases. Thus all CVs need to be included in deriving an empirical formula. However, there is a critical issue when CVs are close to zero. For instance, if a patient is scheduled to arrive at the exact instant when a physician service finishes, an infinitesimal increase on patient inter-arrival time may make the physician starting documentation rather than meeting the patient so that the patient needs to wait for the whole documentation time, and such delays could propagate to next one. To accommodate this effect, we modify the waiting time when CVs are zero as follows:

$$\widetilde{W}_0^{\text{NPP}} = \frac{W_0^{\text{NPP}} + \frac{N-1}{N\mu_d}}{2}, \quad (8.20)$$

where $\frac{N-1}{N\mu_d}$ represents a heuristic approximation of possible waiting time due to CV perturbation, and W_0^{NPP} is the waiting time when all CVs are zero, which can be obtained using the procedure in Appendix B.

When the BCD model is considered, no matter what values of cv_a , cv_p , cv_d , N and M are, the impact of variation in documentation time is much less than in others, as implied by Figure 8.2. Thus, we discount cv_d . Finally, the empirical formulas are given as:

$$\begin{aligned} \widetilde{W}^{\text{PEP}} &= W_0^{\text{PEP}} + \frac{(W^{\text{PEP}} - W_0^{\text{PEP}})\left(\frac{cv_a^2}{\lambda} + \frac{cv_p^2}{\mu_p}\right)}{\frac{1}{\lambda} + \frac{1}{\mu_p}}, \\ \widetilde{W}^{\text{NPP}} &= \widetilde{W}_0^{\text{NPP}} + \frac{(W^{\text{NPP}} - \widetilde{W}_0^{\text{NPP}})\left(\frac{cv_a^2}{\lambda} + \frac{cv_p^2}{\mu_p} + \frac{cv_d^2}{\mu_d}\right)}{\frac{1}{\lambda} + \frac{1}{\mu_p} + \frac{1}{\mu_d}}, \\ \widetilde{W}^{\text{BDC}} &= W_0^{\text{BDC}} + \frac{(W^{\text{BDC}} - W_0^{\text{BDC}})\left(\frac{cv_a^2}{\lambda} + \frac{cv_p^2}{\mu_p} + \frac{cv_d^2}{10\mu_d}\right)}{\frac{1}{\lambda} + \frac{1}{\mu_p} + \frac{1}{10\mu_d}}, \end{aligned} \quad (8.21)$$

where all W_0^x can be obtained through the algorithm in Appendix B.

Using patient waiting time W to define ϵ^x and $\tilde{\epsilon}^x$ similar to (8.19), under the same parameter setting, we obtain the mean and maximal values in different parameter combinations for each model, as shown in Table 8.2. The average differences are 1-2 minutes

only, and the largest one is 6.52 minutes. For the relative error, the average is 5%-7% and the largest is 22.76% (which occurs in an extreme parameter setting, $\lambda = 2, \mu_p = 2, \mu_d = 2$). Under other more reasonable parameters, the absolute difference is almost negligible, which implies that the empirical formulas for patient waiting time capture the general patterns well.

Table 8.2: Accuracy of patient waiting time formula

x	$E[\epsilon^x]$		$E[\tilde{\epsilon}^x]$	
	Mean (min)	Max (min)	Mean (%)	Max (%)
PEP	2.23	6.52	7.25	22.76
NPP	1.23	4.88	6.48	12.38
BDC	1.52	4.47	5.05	6.76

8.5.2.3 Documentation waiting time

While the total working time and patient waiting time show relatively consistent movements in terms of CV, the documentation waiting time exhibits more complicated trends. Both increasing or decreasing patterns with respect to different CVs can be observed. For example, in the PEP model, when there is no variation, all documentations may need to wait until all patients are served. However, when CV is not zero, it could be possible that a documentation can be finished before a new patient arrives. Similar scenarios can be observed in the NPP model. In both models, the documentation waiting time could decrease or increase concavely with respect to CVs in most cases. In the BDC model, the documentation waiting time follows a trend similar to that of patient waiting time.

Therefore, the empirical formulas have the following formats:

$$\begin{aligned}
\tilde{Q}^{\text{PEP}} &= Q^{\text{PEP}} - (Q^{\text{PEP}} - Q_0^{\text{PEP}}) \cdot \frac{\frac{(1-cv_a)^2}{\lambda} + \frac{(1-cv_p)^2}{\mu_p} + \frac{(1-cv_d)^2}{\mu_d}}{\frac{1}{\lambda} + \frac{1}{\mu_p} + \frac{1}{\mu_d}}, \\
\tilde{Q}^{\text{NPP}} &= Q^{\text{NPP}} - (Q^{\text{NPP}} - \tilde{Q}_0^{\text{NPP}}) \cdot \frac{\frac{(1-cv_a)^2}{\lambda} + \frac{(1-cv_p)^2}{\mu_p} + \frac{(1-cv_d)^2}{\mu_d}}{\frac{1}{\lambda} + \frac{1}{\mu_p} + \frac{1}{\mu_d}}, \\
\tilde{Q}^{\text{BDC}} &= Q_0^{\text{BDC}} + (Q^{\text{BDC}} - Q_0^{\text{BDC}}) \cdot \frac{\frac{cv_a^2}{\lambda} + \frac{cv_p^2}{\mu_p} + \frac{cv_d^2}{10\mu_d}}{\frac{1}{\lambda} + \frac{1}{\mu_p} + \frac{1}{10\mu_d}},
\end{aligned} \tag{8.22}$$

where

$$\tilde{Q}_0^{\text{NPP}} = \frac{Q_0^{\text{NPP}} + \frac{N-1}{N\mu_d}}{2}, \tag{8.23}$$

and Q_0^x can be evaluated using the algorithm in Appendix B.

Table 8.3 gives the mean and maximal values of ϵ^x and $\tilde{\epsilon}^x$ (defined similar to (8.19) but replaced with \tilde{Q}) using various parameters for each model. The average differences are within 4 minutes and the largest one is less than 5 minutes, while the average relative error is around 7% and the largest is about 12%. Thus, acceptable estimation of documentation waiting time is obtained.

Table 8.3: Accuracy of documentation waiting time formula

x	$E[\epsilon^x]$		$E[\tilde{\epsilon}^x]$	
	Mean (min)	Max (min)	Mean (%)	Max (%)
PEP	2.18	4.55	7.41	12.68
NPP	3.83	7.41	7.41	11.87
BDC	0.88	1.58	1.69	3.60

Remark 8.8. *More experiments have been carried for $N = 14$, and similar results to Tables 8.1-8.3 are obtained, which indicate that the empirical formulas can achieve acceptable accuracy.*

Remark 8.9. *In this study, the workflow models have single distributions for patient interarrival time, service time, and documentation time, respectively. However, in reality, depending on patient's*

underlying conditions and services in needs, the duration of their services and the corresponding documentations can be longer or shorter. If patients can be classified based on their expected service times or documentation times (e.g. by machine learning models or by manual coordination), the distributions of service and documentation times can be determined according to the patient's class. Each distribution can be first assumed to be exponential, then, by separating the states depending on patient classes or by using mixture of exponential distributions, general distributions of arrival and service times can be obtained. Then, applying the analysis in this section, workflow policies can be evaluated.

8.6 Case Study

8.6.1 Background

To illustrate the applicability of the models developed above, a case study is carried out at a primary care clinic. The clinic belongs to one of the largest integrated health care delivery systems in the country. It provides medical and health services from diagnosis and treatment for minor and chronic illnesses to immunizations and x-rays, for convenient care in the local community. By shadowing the physicians, nurses, and medical assistants (MA) in the clinic, we observe the types and duration of their activities. Specifically, over the first four hours of the observation, all detailed activities such as physical exam, talking with patients, EHR writing, communicating with other staffs, and transitions were observed. Then the activity grains were grouped into patient service and documentation processes, and the main activities to be tracked and analyzed were identified. For the remaining observation periods, each physician and the associated nurse or MA were followed by two different researchers. Their activities and duration on the workflow of each role were recorded for more than 20 hours. Using the observation data, EHR data, and interviews with healthcare providers and supporting staff in the clinic, we obtain the estimates of workflow model parameters.

Although the patient visits are scheduled every 30 minutes, they often come earlier or later than their scheduled time so that the patient arrivals can still be viewed as a random process. Based on the observation study, the mean value and standard deviation of inter-arrival time are 27.1 minutes and 7.49 minutes, respectively. The average face-to-face encounter time a physician spends with a patient is observed as 14.6 minutes with standard deviation as 5.04 minutes. Moreover, it takes 12.3 minutes to work on a documentation task in average and the standard deviation is 9.93 minutes.

In addition, the physician typically works from 8am to 6:30pm, but the slots between 12pm to 1pm and 4pm and 6:30pm are not reserved for patients. Thus, we can assume that between 8am to 12pm and 1pm to 4pm, the physician is scheduled to serve 8 and 6 patients (i.e., $N = 8$ and $N = 6$), respectively. As the waiting queues could be smoothed out during the breaks, two models for 8 and 6 patients can be developed. Then the final performance measures can be obtained by combining the results.

8.6.2 Model Results

Using these parameters and the empirical formulas derived above, we analyze the performance for different workflow models and compare the results. Tables 8.4 and 8.5 show the resulting performance measures. For the BDC model, two different values of batch size are used, $M = 3$ and $M = 5$.

Table 8.4: Performance measures in a morning shift ($N = 8$)

	\tilde{T} (hour)	\tilde{W} (min)	\tilde{Q} (min)
PEP	4.21	0.87	26.24
NPP	4.21	3.11	20.31
BDC ($M = 3$)	4.58	11.00	28.99
BDC ($M = 5$)	4.95	14.87	59.92

As one can see, in both morning and afternoon shifts, the physician has to work overtime. The waiting times (both patient and documentation) are shorter in the afternoon since fewer patients are scheduled so that less accumulations exist. In addition, the results

Table 8.5: Performance measures in an afternoon shift ($N = 6$)

	\tilde{T} (hour)	\tilde{W} (min)	\tilde{Q} (min)
PEP	3.28	0.74	19.73
NPP	3.28	2.63	14.75
BDC ($M = 3$)	3.65	7.39	34.06
BDC ($M = 5$)	4.02	9.03	66.15

indicate that either PEP or NPP policy is acceptable since these two models lead to short waiting times. Particularly, the patient has almost no waiting in the PEP model, because the model allows interruption of documentation in order to meet with patients first. However, it leads to longer documentation time. While in the NPP model, the patient waiting time is only slightly longer but tolerable, and the documentation queue is reduced substantially. For the BDC model, it results in longer physician working hours, and much longer waiting times for both patients and documentations. However, it should be noted that if the documentation service time becomes longer, it is possible that the BDC model can achieve better performance (as discussed in Subsection 8.4.3).

By summing up the working hours, and calculating the weighted average of waiting times, we obtain the daily measurements, as shown in Table 8.6. The resulting average daily working times from all models are longer than seven hours. Adding the two non-reserved time slots (3.5 hours), the total business hours are over 11 hours, which reflects the reality that a physician typically spends several hours after his/her regular office works to complete the unfinished work. The PEP and NPP models achieve similar performance. If the PEP policy is adopted, the patient waiting time can be ignored, while the documentation waiting time is longer than that when the NPP policy is followed. In case that reducing the possibility of information loss and improving the quality of EHR data become more critical, the NPP policy can achieve a good balance by increasing a minimal patient waiting time.

Table 8.6: Performance measures in a full day ($N = 14$)

	\bar{T} (hour)	\bar{W} (min)	\bar{Q} (min)
PEP	7.49	0.81	23.45
NPP	7.49	2.90	17.93
BDC ($M = 3$)	8.23	9.45	31.16
BDC ($M = 5$)	8.97	12.37	62.59

8.6.3 Sensitivity Analysis

To investigate the sensitivity of the results, we first increase the number of patient visits per day to understand the impact of higher patient volume on system performance. For illustration purpose, only the results of the morning shift (from 8am to 12pm, i.e., $N = 8$) are presented here. Specifically, given fixed parameters on service and documentation rates and CVs, the number of patient visits in the morning is increased from 8 to 9 and 10. Meanwhile, the inter-arrival times are reduced accordingly (i.e., with 11% and 20%, respectively). The results are summarized in Table 8.7, where the values in parentheses indicate the changes compared with those in Table 8.4. Apparently, longer working times and patient waiting times are observed in all models. In the PEP model, the patient waiting time only shows a slight increment but the documentation waiting time is extended almost two or three times when $N = 9$ and 10, respectively. Whereas the NPP model has the same working time as the PEP model, it results in an acceptable range of patient and documentation waiting times. The BDC model in both scenarios leads to longer patient waiting times, while the changes in documentation waiting times are not monotonic. A possible reason is that for $M = 3$, when $N = 9$, all documents need to wait until a complete batch before processing. While when $N = 10$, the last document does not need to wait for a full batch, which reduces the average waiting time. Moreover, the interacting impact of variations in both N and λ may also result in non-monotonicity. In summary, the NPP model could be preferred if a higher patient volume is expected.

In order to accommodate the possible increase in patient volumes (from $N = 8$ to $N = 9$ or 10), the service and documentation times need to be reduced. To investigate

Table 8.7: Sensitivity Analysis with increased number of patient visits

N	Model	\bar{T} (hour)	\bar{W} (min)	\bar{Q} (min)
9	PEP	4.56 (+0.35)	1.05 (+0.18)	48.51 (+22.27)
	NPP	4.56 (+0.35)	6.03 (+2.92)	29.35 (+9.04)
	BDC ($M = 3$)	4.86 (+0.28)	16.72 (+5.72)	30.60 (+1.61)
	BDC ($M = 5$)	5.15 (+0.20)	18.24 (+3.37)	58.98 (-0.94)
10	PEP	4.95 (+0.74)	1.25 (+0.28)	69.12 (+42.88)
	NPP	4.95 (+0.74)	6.89 (+3.78)	41.23 (+20.92)
	BDC ($M = 3$)	5.17 (+0.59)	26.28 (+15.28)	26.75 (-2.24)
	BDC ($M = 5$)	5.40 (+0.45)	21.74 (+6.87)	61.61 (+1.69)

the effects of such reductions, Table 8.8 presents the results when the patient service and documentation times are decreased by 10%, while the other parameters are kept the same. For completeness, the scenarios of 10% increase in service and documentation times, which are analogous to reducing inter-arrival times, are also presented. It is observed that the effect of increased service and documentation times is almost identical to that of increasing number of patients. In other words, increasing N by 1 and reducing service and documentation times by 10% are indifferent options in terms of system performance. Again, the NPP model is more favorable in both cases.

The above results indicate that reducing service and documentation times can be a viable way to balance the impact of higher patient volume. To investigate the effect of such practices, we increase N to 9, and reduce mean service and documentation times by 10%. The results are shown in Table 8.9. As one can see, the differences between the results in Table 9 and Table 4 are inconsiderable. This indicates that, in order to accept one more patient in the morning session, the mean service and documentation times could be reduced by 10% to mitigate the impact of patient volume increment.

Table 8.8: Sensitivity Analysis with varying patient service and documentation times

Change	Model	\tilde{T} (hour)	\tilde{W} (min)	\tilde{Q} (min)
+10%	PEP	4.53 (+0.32)	1.05 (+0.18)	43.84 (+17.60)
	NPP	4.53 (+0.32)	5.84 (+2.73)	26.32 (+6.01)
	BDC ($M = 3$)	4.87 (+0.29)	16.76 (+5.76)	30.58 (+1.59)
	BDC ($M = 5$)	5.20 (+0.25)	18.47 (+3.60)	62.44 (+2.52)
-10%	PEP	4.12 (-0.09)	0.67 (-0.20)	20.69 (-5.55)
	NPP	4.12 (-0.09)	2.42 (-0.69)	15.65 (-4.66)
	BDC ($M = 3$)	4.38 (-0.20)	7.8 (-3.20)	28.97 (-0.02)
	BDC ($M = 5$)	4.68 (-0.27)	11.35 (-3.52)	57.26 (-2.66)

Table 8.9: Sensitivity Analysis with increased patient visits and reduced service and documentation times

N	Model	\tilde{T} (hour)	\tilde{W} (min)	\tilde{Q} (min)
9	PEP	4.17 (-0.04)	0.82 (-0.05)	31.22 (+4.98)
	NPP	4.17 (-0.04)	3.18 (+0.07)	21.00 (+0.69)
	BDC ($M = 3$)	4.51 (-0.07)	9.67 (-1.33)	28.44 (-0.55)
	BDC ($M = 5$)	4.84 (-0.11)	13.71 (-1.16)	55.82 (-4.10)

The above analyses are carried out with respect to mean values, while all CVs, i.e., cv_a , cv_p , and cv_d , are fixed. To investigate the impact of process variability, we increase and decrease values of cv_a , cv_p and cv_d by 10% without changing the mean values. In other words, the standard deviations are changed by $\pm 10\%$. As shown in Table 8.10, the resulting performance measures change in the same direction of CV variations, which is a result from the empirical formulas of \tilde{T} , \tilde{W} , and \tilde{Q} developed in Section 8.5. As the resulting performances in all models vary with similar rates, the model preference does not change

in response to slight variations of CVs.

Table 8.10: Sensitivity Analysis with varying CVs

Change	Model	\tilde{T} (hour)	\tilde{W} (min)	\tilde{Q} (min)
+10%	PEP	4.23 (+0.02)	1.03 (+0.16)	27.65 (+1.41)
	NPP	4.23 (+0.02)	3.7 (+0.59)	21.4 (+1.09)
	BDC ($M = 3$)	4.59 (+0.01)	11.41 (+0.41)	29.05 (+0.06)
	BDC ($M = 5$)	4.95 (0.00)	15.14 (+0.27)	59.95 (+0.03)
-10%	PEP	4.18 (-0.03)	0.69 (-0.18)	24.3 (-1.94)
	NPP	4.18 (-0.03)	2.51 (-0.60)	18.81 (-1.50)
	BDC ($M = 3$)	4.55 (-0.03)	10.64 (-0.36)	28.83 (-0.16)
	BDC ($M = 5$)	4.93 (-0.02)	14.66 (-0.21)	59.75 (-0.17)

8.7 Conclusions

This chapter introduces terminating Markov chain models to study physician workflow management policies in primary care clinics. Two types of physician's work are considered: face-to-face encounters with patients, and documentation (EHR data inputs, summary, communication, etc.). Three workflow models, preemptive priority (stop documentation if a new patient arrives), non-preemptive priority (finishing ongoing documentation even if a new patient arrives), and batch documentation (starting documentation only until a given threshold of accumulated documents is achieved), are formulated and analyzed. The physician daily working hours, patient waiting times, and documentation waiting times are evaluated and compared. Then the study is extended to consider non-Markovian scenarios. Empirical formulas are derived to evaluate the performance based on mean and

coefficients of variations of service and interarrival times. Finally, a case study at a primary care clinic is carried out to illustrate and validate model applicability.

Chapter 9

Conclusions and Extensions

9.1 Summary of Contributions

This dissertation has introduced various analytical models. Such models and frameworks are utilized to achieve a better delivery of care in the following phases: (1) At discharge from hospital, a patient's status or risk (risk of readmission or disease exacerbation) is predicted by means of machine learning models. (2) Interventions or treatments to the patient is determined based on the patient identification, and optimization methods are utilized to support the decision processes. (3) Further, stochastic modeling methods are involved in care practice to implement operation policies. Through integrating the classification, optimization and stochastic process models, data analytics based patient specific decision support can be expected to help healthcare providers to reduce COPD readmission rates, avoid opioid over-prescription to curtail opioid epidemic, and improve physician workflow. To conclude this dissertation, we discuss how the models or tools introduced in each chapter contributes to the improvement of the current healthcare management.

Frameworks for COPD readmission reduction: In the dissertation, two different frameworks are developed for health management of COPD patients and following readmission reduc-

tion. The first framework classifies COPD patients into high or low risk of readmissions, then based on the risk group, an intervention resource allocation is determined. In the second framework, a causal model describing relationships between factors and readmission is integrated to a Markov decision process to propose a dynamic intervention planning.

A causal Bayesian network developed in Chapter 3 is utilized in both frameworks by classifying COPD patients based on their risk of readmission for the subsequent analysis in Chapter 5 and by analyzing the effect of variable manipulation in order to provide optimal interventions in the subsequent analysis in Chapter 6. The developed causal model can analyze how the patients' readmission risk is affected as they go through post-discharge intervention, which is one of key contributions of the work.

Chapter 5 assesses the cost-effectiveness of intervention strategies based on risk predictions. Specifically, an optimization model introduced in 5 investigates the post-discharge intervention process and determines the optimal resource allocation for interventions by taking into account the patient's predicted risk. The key contribution of the work is that this work presents a novel framework where the readmission risk prediction model provides clinically relevant stratification of risk and presents information to trigger a transitional and patient-tailored care intervention.

In order to identify the readmission risk factors and design a personalized intervention for each patient, a dynamic intervention decision framework for reducing COPD readmission risks is also developed in Chapter 6. The proposed model can determine the optimal decision on intervention planning minimizing the risk of readmissions and update the decision dynamically based on the patient's altered status. Another critical key contribution of this entire work is to develop such a novel framework and provide an algorithm to solve the MDP having the high-dimensional structure and its structural properties with proofs.

There are various potential applications of the proposed framework as it provides a means of establishing dynamic intervention plans for different health conditions and can be extended to other diseases to evaluate various treatments and intervention procedures.

In such different applications, the framework can be potentially used to support decision making for personalized care.

A framework for opioid prescription optimization: For opioid prescription optimization, a semi-supervised classification model developed in Chapter 4 predicts patients' expected opioid consumption levels, and a stochastic program provided in Chapter 7 determines how many opioids should be prescribed to each class of patients to avoid opioid over-prescription.

In order to effectively manage the opioid prescription, the opioid doses are differentiated to each patient by introducing an opioid usage prediction model and classifying the patients based on the expected opioid usage levels. By collaborating with SSM Health, a survey study is first conducted to collect exact amount of consumed opioids. Based on the survey results and EHR data, a novel approach for semi-supervised learning is proposed in Chapter 4 that results in improved classification performance, which is the first contribution of this work. The model classifies patients into three groups based on their expected opioid consumption levels.

Then, a stochastic programming model is formulated in Chapter 7 to determine the optimal amount of opioids prescribed at discharge to patients in three groups, respectively. The results substantially reduced excess opioids without increasing number of refills. Thus, the main contributions of this work is to provide a framework that enables patient-tailored decision making for opioid prescription. As the existing literature does not provide a mathematical framework for opioid prescriptions nor consider individual patient information in the prescription models, the proposed model can contribute to analyzing consequences of prescription plan, reducing over-prescriptions, and curtailing opioid crisis.

Physician's workflow management: In order to facilitate primary care physician's work efficiency and patient's satisfaction, physician's work process needs to be carefully managed.

Thus, in Chapter 8, various physician's workflow management policies are proposed through stochastic modeling where terminating Markov chain models are developed.

The proposed Markov model considers an absorbing or terminating state which indicates the state where all tasks of the day have been completed by the physician. Closed formulas to evaluate the efficacy of such workflows in terms of physician working time, patient and documentation waiting times, are derived and the corresponding performance measures are compared. Both Markovian and non-Markovian scenarios are studied. Finally, a case study in a primary care clinic is carried out to illustrate model applicability. As a contribution, this work provides a quantitative tool for primary care physicians to design and manage their workflow to improve care quality.

9.2 Extensions

This section describes two potential extensions to the work introduced in this dissertation. The first problem considers the extension of the CNMDP model proposed in Chapter 6 to generalize the patient's condition variation to cope with various circumstances and longer monitoring periods. The second problem relates to integrating data analysis in workflow models provided in Chapter 8 in order to consider different patient and service types which yield varied distributions. The details of each problem are described in the following subsections.

9.2.1 Extension 1: Compliance in the CNMDP Model

In Chapter 6, a causal Bayesian network is extended to be a dynamic network through the proposed assumptions. In the analysis, we assume that patients will be compliant with physician's instruction and the patients' conditions will not be deteriorate (Assumption 6.3). Based on the assumptions, conditional probabilities corresponding to inter-timeslice

edge are set up as in Lemmas 6.4 and 6.5 and the following mathematical properties are derived.

If patients become not compliant with the instruction for their health management (e.g. resuming smoking) or their conditions deteriorate over time (e.g. getting anxious about their health or losing home equipment), the assumption does not hold. In such cases, since the corresponding variables x 's can transit to worse states, the conditional probabilities given in Lemmas 6.4 and 6.5 should be modified so that the CNMDP model can cope with the generalized conditions. When conditional probabilities are set up reflecting other circumstances, different structural properties, such as the existence of a monotone optimal policy, will be derived.

9.2.2 Extension 2: Data-Driven Workflow Models

The current workflow models in Chapter 8 do not distinguish patients, services, or documentation types. However, as it is referred in Remark 8.9, depending on patient's underlying conditions and services in needs, the duration of their services and the corresponding documentations can be longer or shorter. In the future work, as in other framework introduced in this dissertation, machine learning can be utilized in the workflow analysis in order to improve the current analysis. Using machine learning models, patients can be classified based on their expected service times or documentation times. Then, based on the classes and the corresponding expected times, patients can be scheduled so that physician's idle time and overwork time can be minimized. In addition, different from the proposed analysis where documentations are worked as 'first come, first served' basis. However, depending on the expected time, it would also be possible to schedule the documentation task in order to minimize patient's or documentation task's waiting time. These scheduling policies can be a part of the physician's workflow models. This integration of data analytics into the workflow model will enable more patient-centered analysis and improve the current workflow system of primary care clinics.

Appendices

Appendix A: Proofs of Chapter 6

Proof of Lemma 6.1: The defined relation \preceq is a partial order as it satisfies the conditions of a partial order as follows: $s \preceq s$ (reflectivity); if $s \preceq s'$ and $s' \preceq s$, then $s = s'$ (antisymmetry); if $s \preceq s'$ and $s' \preceq s''$, then $s \preceq s''$ (transitivity). ■

Proof of Lemma 6.2: Let $s^t = ((s_{x_1}^t, \dots, s_{x_{|X|}}^t), s_y^t)$. By Assumption 6.3, $s_x^1 \leq s_x^2$ for all $x \in \mathcal{X}$ and $s_y^1 \leq s_y^2$. Then $s^1 \preceq s^2$. Likewise, $s_x^t \leq s_x^{t+1}$ for all $x \in \mathcal{X}$ and $s_y^t \leq s_y^{t+1}$ so that $s^t \preceq s^{t+1}$ for all t . Therefore, $s_x^1 \leq s_x^t$ for all $x \in \mathcal{X}$ and $s_y^1 \leq s_y^t$ for all t . Then, $S^T \subseteq S^{T-1} \subseteq \dots \subseteq S^2 = \{s \mid s \succeq s^1, s \in S\}$. ■

Proof of Lemma 6.4: By Assumption 6.2 and the given condition in Lemma 6.4, the state does not change at the next epoch $t + 1$. Thus, (6.8) holds. Meanwhile, the condition $\text{Parent}_G(x)^{t+1} \neq \text{Parent}_G(x)^t$ in (6.9) implies that an intervention affecting x 's ancestor variable has been taken and improved a state of at least one of parent nodes of x . In this case, the transition probability regarding x depends on the underlying Bayesian network. By Assumptions 6.1 and 6.3, the transition only occurs to better states, thus the conditional probability is normalized by $\sum_{s \geq v} p_G(x^{t+1} = s \mid \text{Parent}_G(x)^{t+1})$. Thus, (6.9) holds. ■

Proof of Lemma 6.5: Due to the inference rule of a causal Bayesian network summa-

rized in 6.4, when x is manipulated by an intervention a_x , x 's original dependency on its parents and the corresponding conditional probability is not considered. Therefore, $p(x^{t+1} | x^t, \text{Parent}_G(x)^{t+1}, \text{Parent}_G(x)^t, a_x) = p(x^{t+1} | x^t, a_x)$, and $p(x^{t+1} | x^t, a_x)$ is defined by the efficacy of intervention a_x , $p_{a_x, (v', v)}$. Because of Assumption 6.3, $p_{a_x, (v', v)} = 0$ if $v' < v$. Thus, (6.10) holds. ■

Proof of Lemma 6.6: For all x' such that $x' \in \mathcal{X} \setminus L(x)$, x' does not have x as its ancestor. Then by Assumption 6.2, x' does not vary at the next time epoch. Thus, $p(s'_{x'} | s_{x'}, s'_{\text{Parent}_G(x')}, s_{\text{Parent}_G(x')}, a_x)$ equals to 1 if $s'_{x'} = s_{x'}$ and 0 otherwise. ■

Proof of Lemma 6.7:

$$\begin{aligned} \sum_{s'_{\mathcal{X}} \in S_{\mathcal{X}}} p_G(y = 1 | s'_{\mathcal{X}}) \cdot p(s'_{\mathcal{X}} | s_{\mathcal{X}}, a) &\leq \max_{s''_{\mathcal{X}} \in S_{\mathcal{X}}} \{p_G(y = 1 | s''_{\mathcal{X}})\} \cdot \sum_{s'_{\mathcal{X}} \in S_{\mathcal{X}}} p(s'_{\mathcal{X}} | s_{\mathcal{X}}, a) \\ &= \max_{s''_{\mathcal{X}} \in S_{\mathcal{X}}} \{p_G(y = 1 | s''_{\mathcal{X}})\} \end{aligned}$$

By Assumption 6.1 and 6.3, $\max_{s''_{\mathcal{X}} \in S_{\mathcal{X}}} \{p_G(y = 1 | s''_{\mathcal{X}})\} \leq p_G(y = 1 | s_{\mathcal{X}})$. Thus, $\sum_{s'_{\mathcal{X}} \in S_{\mathcal{X}}} p_G(y = 1 | s'_{\mathcal{X}}) \cdot p((s'_{\mathcal{X}}, 1) | (s_{\mathcal{X}}, 0), a) \leq p_G(y = 1 | s_{\mathcal{X}})$. ■

Proof of Proposition 6.1: By Lemma 6.6, variables outside the path, $x' \in \mathcal{X} \setminus P(x)$ are fixed with the same value given any intervention related to $P(x)$. It enables us to analyze a reduced form of CNMDP regarding path $P(x) \subset G$ as the associated network G' . Further, as the intervention costs are the same, the negative total expected readmission rate can replace the value function. Denote $z \in P(x)$ as $z_1 = x, \dots, z_n = y$ where z_{i+1} is a child of z_i in $P(x)$ for all i and $n = |P(x)|$. Define $u_{t,a}(s)$ as the total expected readmission probability given state $s \in S_{P(x)}$ and action a at t , and $u_t(s)$ so that $u_t(s) = \min_a u_{t,a}(s)$

We prove the result by the backward induction. Let $S(k)$ denote a group of states where $z_{n-k} = 1$ and its descendants are 0 ($z_{n-k+1} = z_{n-k+2} = \dots = z_{n-1} = 0$). For example, states

in $S(3)$ have $z_{n-3} = 1$ and $z_{n-1} = z_{n-2} = 0$.

At $t = T$, when $s \in S(1)$ (i.e., $z_{n-1} = 1$), no intervention among the candidates is applicable as the state of the closest variable to y is optimal. Therefore, a_0 is the optimal decision and $u_T(s) = P_{G'}(y = 1 | z_{n-1} = 1)$ for $s \in S(1)$. If $s \in S(2)$, intervention $a_{z_{n-1}}$ is only the applicable intervention. Thus, $a_{z_{n-1}}$ is the optimal decision and $u_T(s) = p_a \cdot P_{G'}(y = 1 | z_{n-1} = 1) + (1 - p_a) \cdot P_{G'}(y = 1 | z_{n-1} = 0)$ for $s \in S(2)$. If $s \in S(3)$, both interventions $a_{z_{n-1}}$ and $a_{z_{n-2}}$ are applicable and the total expected readmission probability given each intervention is

$$\begin{aligned} u_{T,a_{z_{n-1}}}(s) &= p_a \cdot P_{G'}(y = 1 | z_{n-1} = 1) + (1 - p_a) \cdot P_{G'}(y = 1 | z_{n-1} = 0), \\ u_{T,a_{z_{n-2}}}(s) &= p_a [P_{G'}(y = 1 | z_{n-1} = 1)P_{G'}(z_{n-1} = 1 | z_{n-2} = 1) \\ &\quad + P_{G'}(y = 1 | z_{n-1} = 0)P_{G'}(z_{n-1} = 0 | z_{n-2} = 1)] \\ &\quad + (1 - p_a) \cdot P_{G'}(y = 1 | z_{n-1} = 0). \end{aligned}$$

Since $P_{G'}(y = 1 | z_{n-1} = 1) < P_{G'}(y = 1 | z_{n-1} = 0)$, $u_{T,a_{z_{n-1}}}(s) < u_{T,a_{z_{n-2}}}(s)$ then the optimal decision is $a_{z_{n-1}}$ and $u_T(s) = u_T(s')$ for any $s \in S(3)$, $s' \in S(2)$. As such, for all $s \in s(k)$ where $k > 1$, the optimal intervention is $a_{z_{n-1}}$ and $u_T(s) = u_T(s')$ for any $s \in S(k)$, $s' \in S(2)$.

Note that for $t < T$, the probability that readmission occurs at t given s and a is the same with $u_{T,a}(s)$. At $t = T - 1$, When $s^1 \in S(1)$, a_0 is the optimal decision and $u_{T-1}(s^1) = P_{G'}(y = 1 | z_{n-1} = 1) + P_{G'}(y = 0 | z_{n-1} = 1) \cdot u_T(s^1) = u_T(s^1) + P_{G'}(y = 0 | z_{n-1} = 1) \cdot u_T(s)$. If $s^2 \in S(2)$, $a_{z_{n-1}}$ is the optimal decision and $u_{T-1}(s^2) = u_T(s^2) + p_a \cdot P_{G'}(y = 0 | z_{n-1} = 1) \cdot u_T(s^1) + (1 - p_a) \cdot P_{G'}(y = 0 | z_{n-1} = 0) \cdot u_T(s^2)$. If $s^3 \in S(3)$, both interventions $a_{z_{n-1}}$ and $a_{z_{n-2}}$ are applicable and the total expected readmission probability given each intervention

is

$$\begin{aligned}
u_{T-1,az_{n-1}}(s^3) &= u_{T,az_{n-1}}(s^3) + p_a P_{G'}(y = 0 \mid z_{n-1} = 1) u_T(s^1) \\
&\quad + (1 - p_a) P_{G'}(y = 1 \mid z_{n-1} = 0) u_t(s^3), \\
u_{T-1,az_{n-2}}(s^3) &= u_{T,az_{n-2}}(s^3) + p_a [P_{G'}(y = 0 \mid z_{n-1} = 1) P_{G'}(z_{n-1} = 1 \mid z_{n-2} = 1) u_T(s^1) \\
&\quad + P_{G'}(y = 0 \mid z_{n-1} = 0) P_{G'}(z_{n-1} = 0 \mid z_{n-2} = 1) u_T(s^2)] \\
&\quad + (1 - p_a) P_{G'}(y = 1 \mid z_{n-1} = 0) u_T(s^3).
\end{aligned}$$

Since $u_{T,az_{n-1}}(s) < u_{T,az_{n-2}}(s)$, $u_T(s(1)) < u_T(s(2)) = u_T(s(3))$, $u_{T-1,az_{n-1}}(s) < u_{T-1,az_{n-2}}(s)$. Thus, the optimal decision is $a_{z_{n-1}}$ and $u_{T-1}(s^3) = u_{T-1}(s^2)$. Likewise, for all $s \in s(k)$ where $k > 1$, the optimal intervention is $a_{z_{n-1}}$ and $u_t(s^k) = u_t(s^2) > u_t(s^1)$ for all $s^k \in S(k)$ and time t . ■

Proof of Proposition 6.2: Define $q_t(s, a)$ be the total expected reward when action a is chosen at state s and time t , so that $v_t(s) = \max_a q_t(s, a)$. Then,

$$q_t(s, a) = -c_a - c_r \sum_{s'} p(y = 1 \mid s') p(s' \mid s, a) + \sum_{s'} p(y = 0 \mid s') p(s' \mid s, a) v_{t+1}(s').$$

If taking no intervention is preferred to taking a for every state over the entire period, even when there are other interventions added to the action space, a won't be chosen over the period. Therefore, in order to compare the actions of taking intervention (a) and taking no intervention (a_0), we can ignore the existence of other interventions and only need to show that the proposition holds when $A = \{a_0, a\}$.

If a_0 is the optimal at state s and time $t + 1, t + 2, \dots, T$, the state will remain as s over the time, and

$$v_{t+1}(s) = -c_r p(y = 1 \mid s) + p(y = 0 \mid s) v_{t+2}(s), \dots, v_T(s) = -c_r p(y = 1 \mid s).$$

Then, $v_{t+1}(s)$ is represented by the following equation.

$$\begin{aligned} v_{t+1}(s) &= -c_r \cdot p(y = 1 | s) \cdot (1 + p(y = 0 | s) + \cdots + p(y = 0 | s)^{T-t-1}) \\ &= -c_r \cdot p(y = 1 | s) \cdot \frac{1 - p(y = 0 | s)^{T-t}}{1 - p(y = 0 | s)} = -c_r [1 - p(y = 0 | s)^{T-t}]. \end{aligned}$$

Then,

$$\begin{aligned} q_t(s, a_0) - q_t(s, a) &= c_a - c_r \{p(y = 1 | s) + p(y = 0 | s)(1 - p(y = 0 | s)^{T-t}) \\ &\quad - \sum_{s'} p(s' | s, a)(p(y = 1 | s') + p(y = 0 | s')(1 - p(y = 0 | s')^{T-t}))\} \\ &= c_a - c_r \{-p(y = 0 | s)^{T-t+1} + \sum_{s'} p(s' | s, a)p(y = 0 | s')^{T-t+1}\}. \end{aligned}$$

In order that $q_t(s, a_0) - q_t(s, a) \geq 0$,

$$\frac{c_a}{c_r} \geq \sum_{s'} p(s' | s, a)p(y = 0 | s')^{T-t+1} - p(y = 0 | s)^{T-t+1}.$$

Thus, when the initial state at time $t = 1$ is s , if the following holds, intervention a is ignored and a_0 is the optimal policy for the entire decision period.

$$\frac{c_a}{c_r} \geq \max_{t \in \{1, \dots, T\}} \sum_{s'} p(s' | s, a)p(y = 0 | s')^t - p(y = 0 | s)^t. \quad (1)$$

Further, if the ratio between intervention cost and readmission cost, $\frac{c_a}{c_r}$, is larger than $\max_{s,t} [\sum_{s'} p(s' | s, a)p(y = 0 | s')^t - p(y = 0 | s)^t]$, the intervention a won't be taken over the entire decision period for every state s . Particularly, when Assumption 6.4 holds, the threshold in (1) can be represented by a function of intervention efficacy p_a as follows.

$$\max_t \sum_{s'} p(s' | s, a)p(y = 0 | s')^t - p(y = 0 | s)^t = p_a \cdot \max_t [\sum_{s' \neq s} p(s' | s, a)p(y = 0 | s')^t - p(y = 0 | s)^t].$$

■

Proof of Proposition 6.3: Suppose there are two interventions a_1 and a_2 . As in the proof of Proposition 6.2, define $q_t(s, a)$ as

$$q_t(s, a) = -c_a - c_r \sum_{s'} p(y = 1 | s') p(s' | s, a) + \sum_{s'} p(y = 0 | s') p(s' | s, a) v_{t+1}(s').$$

If $q_t(s, a_1) \geq q_t(s, a_2)$, intervention a_1 dominates a_2 at state s and t . If Assumption 6.4 holds so that a affects x_a and p_a is efficacy of intervention a ,

$$\begin{aligned} \sum_{s'} p(y = 1 | s') p(s' | s, a) &= (1 - p_a) p(y = 1 | s) + p_a \sum_{s'} p(y = 1 | s') p(s' | s, s'_{x_a}), \\ \sum_{s'} p(y = 0 | s') p(s' | s, a) v_{t+1}(s') &= (1 - p_a) p(y = 0 | s) v_{t+1}(s) + p_a \sum_{s'} p(y = 0 | s') p(s' | s, s'_{x_a}) v_{t+1}(s'). \end{aligned}$$

In order that $q_t(s, a_1) \geq q_t(s, a_2)$,

$$\begin{aligned} & p_{a_1} [-c_r \{-p(y = 1 | s) + \sum_{s'} p(y = 1 | s') p(s' | s, s'_{x_1})\} \\ & \quad - p(y = 0 | s) v_{t+1}(s) + \sum_{s'} p(y = 0 | s') p(s' | s, s'_{x_1}) v_{t+1}(s')] \\ & \geq p_{a_2} [-c_r \{-p(y = 1 | s) + \sum_{s'} p(y = 1 | s') p(s' | s, s'_{x_2})\} \\ & \quad - p(y = 0 | s) v_{t+1}(s) + \sum_{s'} p(y = 0 | s') p(s' | s, s'_{x_2}) v_{t+1}(s')] \end{aligned}$$

Thus, if $\frac{p_{a_1}}{p_{a_2}}$ is larger than the threshold $R_{s,t}$ given below, p_{a_1} is preferred to p_{a_2} . Otherwise, p_{a_2} dominates p_{a_1} . ■

Appendix B: Probability Tables in Case Study of Chapter 6

Table .1: Conditional Probability Table $p_G(y | x_p)$

	$x_p = 0$	$x_p = 1$		
$y = 0$	0.0545	0.8440	$x_p = 0$: readmitted in the past	$y = 0$: not admit
$y = 1$	0.9455	0.1560	$x_p = 1$: not readmitted in the past	$y = 1$: admit

Table .2: Conditional Probability Table $p_G(x_p | x_a, x_e)$

	$x_a = 0$	$x_a = 1$	$x_a = 0$	$x_a = 1$		
	$x_e = 0$	$x_e = 0$	$x_e = 1$	$x_e = 1$		
$x_p = 0$	0.0763	0.0625	0.2348	0.0310	$x_a = 0$: have anxiety	$x_e = 0$: not have an equipment
$x_p = 1$	0.9237	0.9375	0.7652	0.9690	$x_a = 1$: not feel anxiety	$x_e = 1$: have an equipment

Table .3: Conditional Probability Table $p_G(x_a | x_w)$

	$x_w = 0$	$x_w = 1$	$x_w = 2$	$x_w = 3$
$x_a = 0$	0.6034	0.3950	0.2816	0.3366
$x_a = 1$	0.3966	0.6050	0.7184	0.6634

Table .4: Conditional Probability Table $p_G(x_e | x_s)$

	$x_s = 0$	$x_s = 1$		
$x_e = 0$	0.7108	0.3824	$x_s = 0$: not have support	
$x_e = 1$	0.2892	0.6176	$x_s = 1$: have support	

Appendix C: Proofs of Chapter 7

Proof of Proposition 7.1: Rewrite the objective of Problem 7.2, we have

$$\begin{aligned}
\min \quad & \mathbb{E}_R \left[\sum_{k \in K} w^{(k)} + \lambda \sum_{k \in K} y^{(k)} \right] \\
= \quad & \min \sum_{g \in G} \mathbb{E}_{R_g} \left[\sum_{k \in K_g} w^{(k)} + \lambda \sum_{k \in K_g} y^{(k)} \right] = \sum_{g \in G} \min \mathbb{E}_{R_g} \left[\sum_{k \in K_g} w^{(k)} + \lambda \sum_{k \in K_g} y^{(k)} \right] \\
= \quad & \min \mathbb{E}_{R_L} \left[\sum_{k \in K_L} w^{(k)} + \lambda \sum_{k \in K_L} y^{(k)} \right] + \min \mathbb{E}_{R_M} \left[\sum_{k \in K_M} w^{(k)} + \lambda \sum_{k \in K_M} y^{(k)} \right] \\
& + \min \mathbb{E}_{R_H} \left[\sum_{k \in K_H} w^{(k)} + \lambda \sum_{k \in K_H} y^{(k)} \right],
\end{aligned}$$

where the first equality is due to linearity of expectation. Therefore, the objective can be decomposed to group-wise objectives. Constraints can also be decoupled because $w^{(k)}$ and $y^{(k)}$ are not intertwined for different k 's. Thus, the problem can be decomposed into three sub-problems. ■

Proof of Proposition 7.2: Problem 7.2 in Proposition 7.1 can be simplified as follows:

$$\begin{aligned}
\min_{D_g} \quad & \mathbb{E}_{R_g} \left[\sum_{k \in K_g} (D_g(y^{(k)} + 1) - R_g^{(k)}) + \lambda \sum_{k \in K_g} y^{(k)} \right] \\
\text{s.t.} \quad & D_g(y^{(k)} + 1) \geq R_g^{(k)}, \quad \forall k \in K_g, \\
& D_g \in \mathbb{R}_+, \quad y^{(k)} \in \mathbb{Z}_+, \quad \forall k \in K_g.
\end{aligned}$$

Here, random variables $R_g^{(1)}, R_g^{(2)}, \dots, R_g^{(N_g)}$ are IID. Since $y^{(k)}$ depends on $R_g^{(k)}$, then $y^{(1)}, \dots, y^{(N_g)}$ are also IID. Thus,

$$\min_{D_g} \mathbb{E}_{R_g} \left[\sum_{k \in K_g} (D_g(y^{(k)} + 1) - R_g^{(k)}) + \lambda \sum_{k \in K_g} y^{(k)} \right] = \min_{D_g} |K_g| \cdot \mathbb{E}_{R_g} [(D_g(y + 1) - R_g) + \lambda y].$$

By removing the summations, and considering a problem of one individual patient rather than a group, the problem can be further simplified. Since the mean of opioid demands

in class g , $\mathbb{E}_{R_g} R_g$, is given as a constant, the optimal value of the variables are obtained by solving Problem 7.3. ■

Proof of Theorem 7.1: Denote

$$f(x) = x \cdot \left(\sum_{s=1}^S \left\lceil \frac{(\xi_s - x)_+}{x} \right\rceil + S \right) + \lambda \sum_{s=1}^S \left\lceil \frac{(\xi_s - x)_+}{x} \right\rceil.$$

Define $\Psi := \{\xi \mid \xi \geq x, \xi \in \Xi\}$. Through expansion, we have

$$\begin{aligned} f(x) &= x \cdot \left(\sum_{\xi \in \Xi} \left\lceil \frac{(\xi - x)_+}{x} \right\rceil + S \right) + \lambda \sum_{\xi \in \Xi} \left\lceil \frac{(\xi - x)_+}{x} \right\rceil = (x + \lambda) \cdot \sum_{\xi \in \Xi} \left(I\{\xi \geq x\} \cdot \left\lceil \frac{\xi - x}{x} \right\rceil \right) + Sx \\ &= (x + \lambda) \cdot \sum_{\xi \in \Psi} \left(\left\lceil \frac{\xi}{x} \right\rceil - 1 \right) + Sx = (S - |\Psi|)x - \lambda|\Psi| + (x + \lambda) \cdot \sum_{\xi \in \Psi} \left\lceil \frac{\xi}{x} \right\rceil. \end{aligned}$$

Now, define sets of samples according to the relationship with x as follows: $\Psi_n := \{\xi \mid \xi \geq nx, \xi \in \Xi\}$ for $n = 2, 3, \dots$. Then,

$$\sum_{\xi \in \Psi} \left\lceil \frac{\xi}{x} \right\rceil = |\Psi| + \sum_{n=1}^{\infty} |\Psi_n|.$$

Therefore,

$$f(x) = (S - |\Psi|)x - \lambda|\Psi| + (x + \lambda) \cdot (|\Psi| + \sum_{n=1}^{\infty} |\Psi_n|) = Sx + (x + \lambda) \sum_{n=1}^{\infty} |\Psi_n|$$

The sets Ψ_n , $n = 1, 2, 3, \dots$ can differ depending on the value of x . Therefore, $f(x)$ is not a continuous function, and discontinuities occur at the point where x becomes reciprocals of integers multiplied by any sample. Other than the discontinuity points, a derivative of $f(x)$ exists and

$$f'(x) = S + \sum_{n=1}^{\infty} |\Psi_n| > 0.$$

Thus, except where the discontinuity points locate, function $f(x)$ increases in x , which implies that, $f(x)$ will keep increasing till the rightmost discontinuity point of $x = \max_{\xi \in \Xi} \xi$.

Therefore, the solution of Problem 7.5, $x^* = \operatorname{argmin}_x f(x)$, is bounded by $\max_{\xi \in \Xi} \xi$. ■

Appendix D: Transitions in Chapter 8

In PEP model, the transitions occur based on the following conditions:

- (P1) If $X_p^k = X_d^k = 0$ and $X_c^k < N$, the physician is idle (no patient and documentation are in process or waiting); thus only Transition 1 can occur.
- (P2) If $X_p^k > 0$ and $X_p^k + X_d^k + X_c^k < N$, there are patients in the system so that the physician is meeting with them; thus both Transitions 1 and 2 can occur. But Transition 3 cannot occur because the physician is working on a patient.
- (P3) If $X_p^k > 0$ and $X_p^k + X_d^k + X_c^k = N$, there are patients in the system and no more patient will come; thus only Transition 2 can occur.
- (P4) If $X_p^k = 0$, $X_d^k > 0$ and $X_d^k + X_c^k < N$, no patient is in the system, but documentation is not finished and still more patients will come; thus both Transitions 1 and 3 can occur.
- (P5) If $X_p^k = 0$, $X_d^k > 0$ and $X_d^k + X_c^k = N$, no patient is in the system, documentation is not done yet but no more patient will come; thus only Transition 3 can occur.

In NPP model, the transitions occur based on the following rules:

- (N1) If $X_p^k = X_d^k = 0$ and $X_c^k < N$, the physician is idle; thus only Transition 1 can occur.
- (N2) If $D_d^k = 0$, $X_p^k > 0$ and $X_p^k + X_d^k + X_c^k < N$, there are patients in the system, and more will come, but no documentation is in process; thus both Transitions 1 and 2 can occur.
- (N3) If $D_d^k = 0$, $X_p^k > 0$ and $X_p^k + X_d^k + X_c^k = N$, there are patients in the system, but no patient will come and no documentation is in process; thus only Transition 2 can occur.

- (N4) If $D_d^k = 0$, $X_p^k = 0$, $X_d^k > 0$ and $X_d^k + X_c^k < N$, no patient is in the system but more will come, and documentation is in process; thus both Transitions 3 and 4 can occur.
- (N5) If $D_d^k = 0$, $X_p^k = 0$, $X_d^k > 0$ and $X_d^k + X_c^k = N$, all patients have been served and documentation is in process; thus only Transition 3 can occur.
- (N6) If $D_d^k = 1$ and $X_p^k + X_d^k + X_c^k < N$, a new patient comes while documentation is in process, and more patients will come; thus both Transitions 3 and 4 can occur.
- (N7) If $D_d^k = 1$ and $X_p^k + X_d^k + X_c^k = N$, the last patient arrives while documentation is in process; thus only Transition 3 can occur.

In BDC model, the transitions occur in the following scenarios:

- (B1) If $(X_c^k \bmod M) = 0$, $X_p^k = 0$, $X_d^k < M$ and $X_p^k + X_d^k + X_c^k < N$, no patient is in the system, and previous documentation batch is finished but current one is not started since the batch is not filled, so the physician is idle; thus only Transition 1 can occur.
- (B2) If $(X_c^k \bmod M) = 0$, $X_p^k > 0$, $X_d^k < M$ and $X_p^k + X_d^k + X_c^k < N$, there are patients in the system, and previous documentation batch is finished but current one is not started due to less number of tasks in the batch, so the physician is meeting with a patient; thus both Transitions 1 and 2 can occur.
- (B3) If $(X_c^k \bmod M) = 0$, $X_d^k = M$ and $X_p^k + X_c^k < N - M$, the documentation batch size is reached so the physician will start working on the batch, and more patients will come; thus both Transitions 1 and 3 can occur.
- (B4) If $(X_c^k \bmod M) = 0$, $X_d^k = M$ and $X_p^k + X_c^k = N - M$, the documentation batch size is reached so the physician will start working on the batch, but no patient will come; thus only Transition 3 can occur.
- (B5) If $(X_c^k \bmod M) = 0$, $X_p^k > 0$, $X_d^k < M$ and $X_p^k + X_d^k + X_c^k = N$, there are patients in the system, and previous documentation batch is finished and current one is not

started due to less number of tasks in the batch, so the physician is meeting with a patient, but no new patient will come; thus only Transition 2 can occur.

(B6) If $(X_c^k \bmod M) \neq 0$ and $X_p^k + X_d^k + X_c^k < N$, the current documentation batch is in process and more patients will come; thus both Transitions 1 and 3 can occur.

(B7) If $(X_c^k \bmod M) = 0$, $X_p^k = 0$ and $X_d^k + X_c^k = N$, all patients have been served, so the physician will work on documentation of the last batch no matter the batch size is not reached; thus only Transition 3 can occur.

(B8) If $(X_c^k \bmod M) \neq 0$ and $X_p^k + X_d^k + X_c^k = N$, the current documentation batch is in process and no new patient will come; thus only Transition 3 can occur.

For the transition rate, $\forall l = 0, 1, \dots, N-1; n = 0, 1, \dots, N-l-1; m = 0, 1, \dots, N-l-n-1$, in the PEP model, we have the following transition rates corresponding to Transitions 1 to 3 described in (P1)-(P5):

$$\text{From P1,2,4: } \eta((m, n, l), (m + 1, n, l)) = \lambda,$$

$$\text{From P2,3: } \eta((m + 1, n, l), (m, n + 1, l)) = \mu_p,$$

$$\text{From P4,5: } \eta((0, n + 1, l), (0, n, l + 1)) = \mu_d.$$

In the NPP model, it follows from Transitions 1 to 3 explained in (N1)-(N7) that

$$\text{From N1: } \eta((0, n, l, 0), (1, n, l, 0)) = \lambda, \quad \text{if } n = 0,$$

$$\text{From N4: } \eta((0, n, l, 0), (1, n, l, 1)) = \lambda, \quad \text{if } n > 0,$$

$$\text{From N2: } \eta((m, n, l, 0), (m + 1, n, l, 0)) = \lambda, \\ \text{if } m > 0,$$

$$\text{From N6: } \eta((m, n, l, 1), (m + 1, n, l, 1)) = \lambda, \\ \text{if } m > 0, n > 0,$$

$$\text{From N2,3: } \eta((m + 1, n, l, 0), (m, n + 1, l, 0)) = \mu_p,$$

$$\text{From N4,5: } \eta((0, n + 1, l, 0), (0, n, l + 1, 0)) = \mu_d,$$

$$\text{From N6,7: } \eta((m, n + 1, l, 1), (m, n, l + 1, 0)) = \mu_d, \\ \text{if } m > 0.$$

In the BDC model, from Transitions 1 to 3 outlined in (B1)-(B8), we obtain,

$$\text{From B1-3,6: } \eta((m, n, l), (m + 1, n, l)) = \lambda,$$

$$\text{From B2,5: } \eta((m + 1, n, l), (m, n + 1, l)) = \mu_p, \\ \text{if } (l \bmod M) = 0 \text{ and } n < M,$$

$$\text{From B3,6: } \eta((m, n, l), (m, n - 1, l + 1)) = \mu_d, \\ \text{if } ((l + n) \bmod M) = 0 \text{ and } 1 \leq n \leq M,$$

$$\text{From B4,7,8: } \eta((0, n, l), (0, n - 1, l + 1)) = \mu_d, \\ \text{if } l + n = N \text{ and } 1 \leq n \leq M.$$

Appendix E: Proofs of Chapter 8

Proof of Lemma 8.1: For a given number of completed documentation task l , the number of patients in the system, m , ranges from 0 to $N - l$. With given l and m , the number of unfinished documents is from 0 to $N - l - m$ (i.e., $N - l - m + 1$ cases). Thus,

$$\sum_{m=0}^{N-l} (N - l - m + 1) = \frac{(N + 2 - l)(N + 1 - l)}{2}.$$

For $l \in [0, N]$, we obtain

$$\begin{aligned} K^{PEP} &= \sum_{l=0}^N \frac{(N + 2 - l)(N + 1 - l)}{2} \\ &= \frac{1}{6}(N + 1)(N + 2)(N + 3). \end{aligned}$$

■

Proof of Lemma 8.2: The state space constraints for the number of patients in the system m , the number of unfinished documentation tasks n and the number of documentations tasks l in the NPP model are identical to those in the PEP model. Therefore, when $D_d^k = 0$, the number of feasible states equals to K^{PEP} . When $D_d = 1$, m and n cannot be 0 due to the constraints of D_d . The number of states for $m = 0$ or $n = 0$ is

$$\frac{(N + 1)(N + 2)}{2} + \frac{(N + 1)(N + 2)}{2} - (N + 1) = (N + 1)^2.$$

Thus, we have

$$\begin{aligned} K^{NPP} &= 2K^{PEP} - (N + 1)^2 \\ &= \frac{1}{3}(N + 1)(N^2 + 2N + 3) \end{aligned}$$

■

Proof of Lemma 8.3: For a given number of completed documentation tasks l , when l is a multiple of the batch size M , the number of unfinished documentation tasks n varies from 0 to $\min(M, N - l)$. Furthermore, with a fixed n , the number of patients in the system m varies from 0 to $N - n - l$ (i.e., $N - n - l + 1$ cases).

Let $A = \lfloor \frac{N}{M} \rfloor$ and $B = N - AM$. If $l = Mk$, where k is a non-negative integer, then

$$\begin{aligned} & \sum_{n=0}^M (N - n - Mk + 1) \\ &= \frac{(M + 1)[2N + 2 - M(1 + 2k)]}{2}, \quad \text{if } k < A, \\ & \sum_{n=0}^{N-MA} (N - n - MA + 1) = (B + 1) + B + \cdots + 1 \\ &= \frac{(B + 2)(B + 1)}{2}, \quad \text{if } k = A. \end{aligned}$$

If $l = Mk + j$ where j is a positive integer and $j < M$, then, n takes $M - j$ so that $n < M$ and $((n + l) \bmod M) = 0$ due to state space constraints, and m varies from 0 to $N - M(1 + k)$. Therefore, when $Mk < l < \min(M(1 + k), N)$, the number of feasible states is $(M - 1)(N + 1 - M(1 + k))$ if $k < A$, and B if $k = A$. Thus, we can obtain

$$\begin{aligned} K^{BDC} &= \sum_{k=0}^{A-1} \left[\frac{(M + 1)(2N + 2 - M(1 + 2k))}{2} \right. \\ &\quad \left. + (M - 1)(N + 1 - M(1 + k)) \right] \\ &\quad + \frac{(B + 2)(B + 1)}{2} + B \\ &= \frac{MA(4N - M - 2MA + 5) + (B^2 + 5B + 2)}{2}. \end{aligned}$$

■

Appendix F: Algorithms for $CV = 0$ in Chapter 8

To evaluate the performance measures when $CV = 0$, calculation algorithms can be introduced for each model. In these algorithms, ta , ts , and td denote inter-arrival, service, documentation times, and N and M represent total number of patients and documentation batch size, respectively. Note that closed formulas for T_0^{PEP} , T_0^{NPP} , T_0^{BDC} and W_0^{PEP} can be obtained through the algorithm.

Algorithm 4: PEP model

```

current  $\leftarrow$  0, pwait(0 : N)  $\leftarrow$  0, dwait(0 : N)  $\leftarrow$  0;
arrival  $\leftarrow$  0, dstartwait(0 : N)  $\leftarrow$  0, doctime(0 : N)  $\leftarrow$  0;
for i = 1 to N do
  arrival  $\leftarrow$  i · ta ;
  if current  $\leq$  arrival then
    pwait(i)  $\leftarrow$  0;
    current  $\leftarrow$  arrival + ts;
  else
    pwait(i)  $\leftarrow$  current - arrival;
    current  $\leftarrow$  current + ts;
  end
  dstartwait(i)  $\leftarrow$  current;
  doctime(i)  $\leftarrow$  td;
  j  $\leftarrow$  0;
  if i=N then
    for j=1 to N do
      if doctime(j) > 0 then
        dwait(j)  $\leftarrow$  dwait(j) + (current - dstartwait(j));
        current  $\leftarrow$  current + doctime(j);
        doctime(j)  $\leftarrow$  0;
      end
    end
  else
    while (current < arrival + ta) and (j < i) do
      j  $\leftarrow$  j + 1;
      if doctime(j) > 0 then
        dwait(j)  $\leftarrow$  dwait(j) + (current - dstartwait(j));
        if current + doc(j)  $\leq$  arrival + ta then
          current  $\leftarrow$  current + doctime(j);
          doctime(j)  $\leftarrow$  0;
        else
          doctime(j)  $\leftarrow$  doctime(j) - (arrival + ta - current);
          current  $\leftarrow$  arrival + ta;
          dstartwait(j)  $\leftarrow$  current;
        end
      end
    end
  end
  Ttotal  $\leftarrow$  current;
  Wpatient  $\leftarrow$  average(pwait(0 : N));
  Qdoc  $\leftarrow$  average(dwait(0 : N));
end

```

Algorithm 5: NPP model

```

current  $\leftarrow$  0, pwait(0 : N)  $\leftarrow$  0, dwait(0 : N)  $\leftarrow$  0;
arrival  $\leftarrow$  0, service(0 : N)  $\leftarrow$  0;
for i = 1 to N do
  arrival  $\leftarrow$  i · ta ;
  if current  $\leq$  arrival then
    service(i)  $\leftarrow$  arrival;
    pwait(i)  $\leftarrow$  0;
    current  $\leftarrow$  arrival + ts;
  else
    service(i)  $\leftarrow$  current;
    pwait(i)  $\leftarrow$  current - arrival;
    current  $\leftarrow$  current + ts;
  end
  for j = 1 to i do
    if doc(j) = 0 then
      if (current < arrival + ta) or (i = N) then
        dwait(j)  $\leftarrow$  current - (service(j) + ts);
        current  $\leftarrow$  current + td;
      end
    end
  end
  Ttotal  $\leftarrow$  current;
  Wpatient  $\leftarrow$  average(pwait(0 : N));
  Qdoc  $\leftarrow$  average(dwait(0 : N));
end

```

Algorithm 6: BDC model

```

current  $\leftarrow$  0, pwait(0 : N)  $\leftarrow$  0, dwait(0 : N)  $\leftarrow$  0;
arrival  $\leftarrow$  0, service(0 : N)  $\leftarrow$  0, batchnum  $\leftarrow$  0;
for i = 1 to N do
    arrival  $\leftarrow$  i · ta ;
    if current  $\leq$  arrival then
        | service(i)  $\leftarrow$  arrival;
        | pwait(i)  $\leftarrow$  0;
        | current  $\leftarrow$  arrival + ts;
    else
        | service(i)  $\leftarrow$  current;
        | pwait(i)  $\leftarrow$  current - arrival;
        | current  $\leftarrow$  current + ts;
    end
    if (i mod M = 0) or (i = N) then
        | batchnum  $\leftarrow$   $\lfloor \frac{i}{M} \rfloor$ ;
        | for j = (batchnum+1) to min(N,(batchnum+M)) do
            | dwait(j)  $\leftarrow$  current - (service(j) + ts);
            | current  $\leftarrow$  current + td;
        | end
    end
    Ttotal  $\leftarrow$  current;
    Wpatient  $\leftarrow$  average(pwait(0 : N));
    Qdoc  $\leftarrow$  average(dwait(0 : N));
end

```

Bibliography

- [1] COPD Foundation, "What is COPD? [online]. Available from: <https://www.copdfoundation.org/What-is-COPD/Understanding-COPD/What-is-COPD.aspx>. [Accessed 21 August 2020]," 2020.
- [2] GBD 2015 Disease and Injury Incidence and Prevalence Collaborators, "Global, regional, and national incidence, prevalence, and years lived with disability for 310 diseases and injuries, 1990-2015: a systematic analysis for the global burden of disease study 2015," *Lancet*, vol. 388, no. 10053, pp. 1545–1602, 2016.
- [3] GBD 2015 Mortality and Causes of Death Collaborators, "Global, regional, and national life expectancy, all-cause mortality, and cause-specific mortality for 249 causes of death, 1980-2015: a systematic analysis for the global burden of disease study 2015," *Lancet*, vol. 388, no. 10053, pp. 1459–1544, 2016.
- [4] C. Mathers and D. Loncar, "Projections of global mortality and burden of disease from 2002 to 2030," *PLoS Medicine*, vol. 3, no. 11, pp. 2011–2029, 2006.
- [5] B. Lomborg, *Global problems, local solutions: Costs and benefits*. Cambridge University Press, 2013.
- [6] J. Vestbo, S. Hurd, A. Agusti, P. Jones, C. Vogelmeier, A. Anzueto, P. Barnes, L. Fabbri, F. Martinez, M. Nishimura, R. Stockley, D. Sin, and R. Rodriguez-Roisin, "Global strategy for the diagnosis, management, and prevention of chronic obstructive pulmonary disease: Gold executive summary," *American Journal of Respiratory and Critical Care Medicine*, vol. 187, no. 4, pp. 347–365, 2013.
- [7] S. Braman, "Hospital readmissions for copd: we can meet the challenge," *Chronic Obstructive Pulmonary Diseases: Journal of the COPD Foundation*, vol. 2, no. 1, pp. 4–7, 2015.
- [8] K. Bahadori and J. M. FitzGerald, "Risk factors of hospitalization and readmission of patients with COPD exacerbation systematic review," *International Journal of Chronic Obstructive Pulmonary Disease*, vol. 3, no. 3, pp. 241–251, 2007.

- [9] V. Prieto-Centurion, M. A. Markos, N. I. Ramey, H. A. Gussin, S. M. Nyenhuis, M. J. Joo, B. Prasad, N. Bracken, R. DiDomenico, P. O. Godwin, and H. A. Jaffe, "Interventions to reduce rehospitalizations after chronic obstructive pulmonary disease exacerbations: a systematic review," *International Journal of Chronic Obstructive Pulmonary Disease*, vol. 11, no. 3, pp. 417–424, 2014.
- [10] T. Shah, M. M. Churpek, M. C. Perrillon, and R. T. KonetzkaShah, "Understanding why patients with COPD get readmitted: a large national study to delineate the Medicare population for the readmissions penalty expansion," *Chest*, vol. 147, no. 5, pp. 1219–1226, 2015.
- [11] O. Hasan, D. O. Meltzer, S. A. Shaykevich, C. M. Bell, P. J. Kaboli, A. D. Auerbach, T. B. Wetterneck, V. M. Arora, J. Zhang, and J. L. Schnipper, "Hospital readmission in general medicine patients: a prediction model," *Journal of General Internal Medicine*, vol. 25, no. 3, pp. 211–219, 2010.
- [12] D. Kansagara, H. Englander, A. Salanitro, D. Kagen, C. Theobald, M. Freeman, and S. Kripalani, "Risk prediction models for hospital readmission - a systematic review," *JAMA*, vol. 306, no. 15, pp. 1688–1698, 2011.
- [13] L. Zeng, S. Neogi, J. Rogers, S. Seidensticker, C. Clark, L. Sonstein, R. Trevino, and G. Sharma, "Statistical models for hospital readmission prediction with application to chronic obstructive pulmonary disease (copd) patients," *Proceedings of the 2014 International Conference on Industrial Engineering and Operations Management*, 2014.
- [14] B. J. Make, G. Eriksson, P. M. Calverley, C. R. Jenkins, D. S. Postma, S. Peterson, O. Ostlund, and A. Anzueto, "A score to predict short-term risk of COPD exacerbations (SCOPEX)," *International Journal of Chronic Obstructive Pulmonary Disease*, vol. 10, pp. 201–209, 2015.
- [15] A. Sanduzzi, P. Balbo, P. Candoli, G. Catapano, P. Contini, A. Mattei, G. Puglisi, L. Santoiemma, and A. Stanziola, "COPD: adherence to therapy," *Multidisciplinary Respiratory Medicine*, vol. 9, no. 1, pp. 1–9, 2014.
- [16] J. Bourbeau and S. Bartlett, "Patient adherence in COPD," *Thorax*, vol. 63, no. 9, pp. 831–838, 2008.
- [17] A. Kolodny, D. T. Courtwright, C. S. Hwang, P. Kreiner, J. L. Eadie, T. W. Clark, and G. C. Alexander, "The prescription opioid and heroin crisis: A public health approach to an epidemic of addiction," *Annual Review of Public Health*, vol. 36, pp. 559–574, 2015.
- [18] C. A. Kahlenberg, J. G. Stepan, A. Premkumar, F. D. Lovecchio, and M. B. Cross, "Institutional Guidelines Can Decrease the Amount of Opioids Prescribed After Total Joint Replacement," *HSS Journal*, vol. 15, no. 1, pp. 27–30, 2019.
- [19] D. Dowell, E. Arias, K. Kochanek, R. Anderson, G. P. Guy, J. L. Losby, and G. Baldwin, "Contribution of opioid-involved poisoning to the change in life expectancy in the United States, 2000–2015," *JAMA*, vol. 318, no. 11, pp. 1065–1067, 2017.

- [20] H. Hedegaard, B. A. Bastian, J. P. Trinidad, M. Spencer, and M. Warner, "Drugs most frequently involved in drug overdose deaths: United States, 2011-2016," *National Vital Statistics Reports*, vol. 67, no. 9, 2018.
- [21] Centers for Disease Control and Prevention, "2018 annual surveillance report of drug-related risks and outcomes—United States," 2018.
- [22] L. Scholl, P. Seth, M. Kariisa, N. Wilson, and G. Baldwin, "Drug and opioid-involved overdose deaths — united states, 2013–2017," *Morbidity and Mortality Weekly Report*, vol. 67, no. 5152, pp. 1419–1427, 2019.
- [23] P. S. Huang and S. N. Copp, "Oral Opioids Are Overprescribed in the Opiate-Naive Patient Undergoing Total Joint Arthroplasty," *Journal of the American Academy of Orthopaedic Surgeons*, vol. 27, no. 15, pp. e702–e708, 2019.
- [24] Centers for Disease Control and Prevention, "Opioid Basics [online]. Available from: <https://www.cdc.gov/drugoverdose/opioids/index.html>. [Accessed 1 April 2020]," 2020.
- [25] National Institute on Drug Abuse, "Opioids-Brief Description [online]. Available from: <https://www.drugabuse.gov/drugs-abuse/opioids>. [Accessed 1 April 2020]," 2020.
- [26] P. Voon, M. Karamouzian, and T. Kerr, "Chronic pain and opioid misuse: A review of reviews," *Substance Abuse Treatment, Prevention, and Policy*, vol. 12, no. 1, p. 36, 2017.
- [27] Substance Abuse and Mental Health Services Administration, "Results from the 2008 national survey on drug use and health: national findings," *Office of Applied Studies, NSDUH Series H-36, HHS Publications No. SMA 09-4434*, 2009.
- [28] C. M. Carey, A. B. Jena, and M. L. Barnett, "Patterns of potential opioid misuse and subsequent adverse outcomes in medicare, 2008 to 2012," *Annals of Internal Medicine*, vol. 168, no. 12, pp. 837–845, 2018.
- [29] H. N. Overton, M. N. Hanna, W. E. Bruhn, S. Hutfless, M. C. Bicket, M. A. Makary, B. Matlaga, C. Johnson, J. Sheffield, R. Shechter, *et al.*, "Opioid-prescribing guidelines for common surgical procedures: an expert panel consensus," *Journal of the American College of Surgeons*, vol. 227, no. 4, pp. 411–418, 2018.
- [30] N. D. Volkow, T. A. McLellan, J. H. Cotto, M. Karithanom, and S. R. B. Weiss, "Characteristics of opioid prescriptions in 2009," *JAMA*, vol. 305, no. 13, p. 1299–1301, 2011.
- [31] C. M. Brummett, J. F. Waljee, J. Goesling, S. Moser, P. Lin, M. J. Englesbe, A. S. Bohnert, S. Kheterpal, and B. K. Nallamotheu, "New persistent opioid use after minor and major surgical procedures in us adults," *JAMA Surgery*, vol. 152, no. 6, pp. e170504–e170504, 2017.

- [32] E. Y. Chen, A. Marcantonio, and P. Tornetta, "Correlation between 24-hour pre-discharge opioid use and amount of opioids prescribed at hospital discharge," *JAMA Surgery*, vol. 153, no. 2, pp. e174859–e174859, 2018.
- [33] A. E. Feinberg, T. R. Chesney, S. Srikandarajah, S. A. Acuna, R. S. McLeod, and B. P. i. S. Group, "Opioid Use After Discharge in Postoperative Patients," *Annals of Surgery*, vol. 267, no. 6, pp. 1056–1062, 2018.
- [34] C. P. Hannon, T. E. Calkins, J. Li, C. Culvern, B. Darrith, D. Nam, T. Gerlinger, A. Buvanendran, and C. J. D. Valle, "Large Opioid Prescriptions are Unnecessary after Total Joint Arthroplasty: A Randomized Controlled Trial," *The Journal of Arthroplasty*, vol. 34, no. 7, pp. S4–S10, 2019.
- [35] M. C. Bicket, E. White, P. J. Pronovost, C. L. Wu, M. Yaster, and G. C. Alexander, "Opioid oversupply after joint and spine surgery," *Anesthesia Analgesia*, vol. 128, no. 2, p. 358–364, 2019.
- [36] D. Dowell, T. M. Haegerich, and R. Chou, "Cdc guideline for prescribing opioids for chronic pain—united states, 2016," *JAMA*, vol. 315, no. 15, pp. 1624–1645, 2016.
- [37] C. Steiner, R. Andrews, M. Barrett, and W. A, *HCUP Projections Report 2012-03*. AHRQ, 2012.
- [38] S. M. Kurtz, E. Lau, K. Ong, K. Zhao, K. M, and K. J. Bozic, "Future young patient demand for primary and revision joint replacement: National projections from 2010 to 2030," *Clinical Orthopedic Related Research*, vol. 467, no. 10, pp. 2606–2612, 2009.
- [39] P. D. Franklin, J. A. Karbassi, W. Li, W. Yang, and D. C. Ayers, "Reduction in narcotic use after primary total knee arthroplasty and association with patient pain relief and satisfaction," *The Journal of Arthroplasty*, vol. 25, no. 6, pp. 12–16, 2010.
- [40] N. A. Bedard, A. J. Pugely, R. W. Westermann, K. R. Duchman, N. A. Glass, and J. J. Callaghan, "Opioid use after total knee arthroplasty: trends and risk factors for prolonged use," *The Journal of Arthroplasty*, vol. 32, no. 8, pp. 2390–2394, 2017.
- [41] M. C. Inacio, C. Hansen, N. L. Pratt, S. E. Graves, and E. E. Roughead, "Risk factors for persistent and new chronic opioid use in patients undergoing total hip arthroplasty: a retrospective cohort study," *BMJ Open*, vol. 6, no. 4, p. e010664, 2016.
- [42] J. Goesling, S. E. Moser, B. Zaidi, A. L. Hassett, P. Hilliard, B. Hallstrom, D. J. Clauw, and C. M. Brummett, "Trends and predictors of opioid use following total knee and total hip arthroplasty," *Pain*, vol. 157, no. 6, p. 1259, 2016.
- [43] S. Kim, N. Choudhry, J. Franklin, K. Bykov, M. Eikermann, J. Lii, M. Fischer, and B. Bateman, "Patterns and predictors of persistent opioid use following hip or knee arthroplasty," *Osteoarthritis and Cartilage*, vol. 25, no. 9, pp. 1399–1406, 2017.
- [44] J. C. Rozell, P. M. Courtney, J. R. Dattilo, C. H. Wu, and G.-C. Lee, "Preoperative opiate use independently predicts narcotic consumption and complications after total joint arthroplasty," *The Journal of Arthroplasty*, vol. 32, no. 9, pp. 2658–2662, 2017.

- [45] B. J. Zarling, S. S. Yokhana, D. T. Herzog, and D. C. Markel, "Preoperative and postoperative opiate use by the arthroplasty patient," *The Journal of Arthroplasty*, vol. 31, no. 10, pp. 2081–2084, 2016.
- [46] V. Wylde, S. Hewlett, I. D. Learmonth, and P. Dieppe, "Persistent pain after joint replacement: Prevalence, sensory qualities, and postoperative determinants," *Pain*, vol. 152, no. 3, pp. 566–572, 2011.
- [47] J. A. Singh and D. G. Lewallen, "Predictors of use of pain medications for persistent knee pain after primary total knee arthroplasty: a cohort study using an institutional joint registry," *Arthritis Research & Therapy*, vol. 14, no. 6, p. R248, 2012.
- [48] A. M. Valdes, S. C. Warner, H. L. Harvey, G. S. Fernandes, S. Doherty, W. Jenkins, M. Wheeler, and M. Doherty, "Use of prescription analgesic medication and pain catastrophizing after total joint replacement surgery," *Seminars in Arthritis and Rheumatism*, vol. 45, no. 2, pp. 150 – 155, 2015.
- [49] A. C. of Physicians, "The impending collapse of primary care medicine and its implications for the state of the nation's health care," 2006.
- [50] T. Bodenheimer, "Primary Care — Will It Survive?," *The New England Journal of Medicine*, vol. 355, no. 9, pp. 861–864, 2006.
- [51] C. A. Sinsky, R. Willard-Grace, A. M. Schutzbank, T. A. Sinsky, D. Margolius, and T. Bodenheimer, "In Search of Joy in Practice: A Report of 23 High-Functioning Primary Care Practices," *The Annals of Family Medicine*, vol. 11, no. 3, pp. 272–278, 2013.
- [52] L. N. Dyrbye and T. D. Shanafelt, "Physician Burnout: A Potential Threat to Successful Health Care Reform," *JAMA*, vol. 305, no. 19, pp. 2009–2010, 2011.
- [53] C. Sinsky, L. Colligan, L. Li, M. Prgomet, S. Reynolds, L. Goeders, J. Westbrook, M. Tutty, and G. Blike, "Allocation of Physician Time in Ambulatory Practice: A Time and Motion Study in 4 Specialties," *Annals of Internal Medicine*, vol. 165, no. 11, p. 753, 2016.
- [54] B. G. Arndt, J. W. Beasley, M. D. Watkinson, J. L. Temte, W.-J. Tuan, C. A. Sinsky, and V. J. Gilchrist, "Tethered to the EHR: Primary Care Physician Workload Assessment Using EHR Event Log Data and Time-Motion Observations," *The Annals of Family Medicine*, vol. 15, no. 5, pp. 419–426, 2017.
- [55] L. Hansen, R. Young, K. Hinami, A. Leung, and M. Williams, "Interventions to reduce 30-day rehospitalization: a systematic review," *Annals of Internal Medicine*, vol. 155, no. 8, pp. 520–528, 2011.
- [56] A. Leppin, M. Gionfriddo, M. Kessler, J. Brito, F. Mair, K. Gallacher, Z. Wang, P. Erwin, T. Sylvester, K. Boehmer, H. Ting, M. Murad, N. Shippee, and V. Montori, "Preventing 30-day hospital readmissions: a systematic review and meta-analysis of randomized trials," *JAMA Internal Medicine*, vol. 174, no. 7, pp. 1095–1107, 2014.

- [57] S. Kripalani, C. Theobald, B. Anctil, and E. Vasilevskis, "Reducing hospital readmission rates: current strategies and future directions," *Annual Review of Medicine*, vol. 65, pp. 471–485, 2014.
- [58] E. Bradley, L. Curry, L. Horwitz, H. Sipsma, Y. Wang, M. Walsh, D. Goldmann, N. White, I. Pina, and H. Krumholz, "Hospital strategies associated with 30-day readmission rates for patients with heart failure," *Circulation: Cardiovascular Quality and Outcomes*, vol. 6, no. 4, pp. 444–450, 2013.
- [59] E. Wick, A. Shore, K. Hirose, A. Ibrahim, S. Gearhart, J. Efron, J. Weiner, and M. Makary, "Readmission rates and cost following colorectal surgery," *Diseases of the Colon & Rectum*, vol. 54, no. 12, pp. 1475–1479, 2011.
- [60] R. C. Martin, R. Brown, L. Puffer, S. Block, G. Callender, A. Quillo, C. R. Scoggins, and K. M. McMasters, "Readmission rates after abdominal surgery: the role of surgeon, primary caregiver, home health, and subacute rehab," *Annals of surgery*, vol. 254, no. 4, pp. 591–597, 2011.
- [61] B. Zmistowski, C. Restrepo, J. Hess, D. Adibi, S. Cangoz, and J. Parvizi, "Unplanned readmission after total joint arthroplasty: rates, reasons, and risk factors," *JBJS*, vol. 95, no. 20, pp. 1869–1876, 2013.
- [62] I. De Alba and A. Amin, "Pneumonia readmissions: risk factors and implications," *Ochsner Journal*, vol. 14, no. 4, pp. 649–654, 2014.
- [63] I. Koulouridis, L. Price, N. Madias, and B. Jaber, "Hospital-acquired acute kidney injury and hospital readmission: a cohort study," *American Journal of Kidney Diseases*, vol. 65, no. 2, pp. 275–282, 2015.
- [64] D. Rubin, "Hospital readmission of patients with diabetes," *Current Diabetes Reports*, vol. 15, no. 17, pp. 1–9, 2015.
- [65] J. Donnelly, S. Hohmann, and H. Wang, "Unplanned readmissions after hospitalization for severe sepsis at academic medical center-affiliated hospitals," *Critical Care Medicine*, vol. 43, no. 9, pp. 1916–1927, 2015.
- [66] P. Almagro, B. Barreiro, A. O. de Echagüen, S. Quintana, M. R. Carballeira, J. L. Heredia, and J. Garau, "Risk factors for hospital readmission in patients with chronic obstructive pulmonary disease," *Respiration*, vol. 73, no. 3, pp. 311–317, 2006.
- [67] Z. Cao, K. Ong, P. Eng, W. Tan, and T. Ng, "Frequent hospital readmissions for acute exacerbation of copd and their associated factors," *Respirology*, vol. 11, no. 2, pp. 188–195, 2006.
- [68] P. Coventry, I. Gemmell, and C. Todd, "Psychosocial risk factors for hospital readmission in copd patients on early discharge services: a cohort study," *BMC Pulmonary Medicine*, vol. 11, no. 1, pp. 1–10, 2011.

- [69] M. Tsui, F. Lun, L. Cheng, A. Cheung, V. Chan, W. Leung, and C. Chu, "Risk factors for hospital readmission for copd after implementation of the gold guidelines," *The International Journal of Tuberculosis and Lung Disease*, vol. 20, no. 3, pp. 396–401, 2016.
- [70] J. Simmering, L. Polgreen, A. Comellas, J. Cavanaugh, and P. Polgreen, "Identifying patients with copd at high risk of readmission," *Chronic Obstructive Pulmonary Diseases: Journal of the COPD Foundation*, vol. 3, no. 4, pp. 729–738, 2016.
- [71] P. Spirtes, "Introduction to causal inference," *Journal of Machine Learning Research*, vol. 11, pp. 1643–1662, 2010.
- [72] A. Gelman and J. Hill, *Data Analysis Using Regression and Multilevel/Hierarchical Models*. Cambridge University Press, 2007.
- [73] M. Ahangaran, M. Jahed-Motlagh, and B. Minaei-Bidgoli, "A Novel Method for Predicting the Progression Rate of ALS Disease Based on Automatic Generation of Probabilistic Causal Chains," *Artificial Intelligence in Medicine*, p. 101879, 2020.
- [74] S. Marini, E. Trifoglio, N. Barbarini, F. Sambo, B. D. Camillo, A. Malovini, M. Manfrini, C. Cobelli, and R. Bellazzi, "A Dynamic Bayesian Network model for long-term simulation of clinical complications in type 1 diabetes," *Journal of Biomedical Informatics*, vol. 57, 2015.
- [75] D. Koch, R. S. Eisinger, and A. Gebharter, "A causal Bayesian network model of disease progression mechanisms in chronic myeloid leukemia," *Journal of Theoretical Biology*, vol. 433, pp. 94–105, 2017.
- [76] M. B. Sesen, A. E. Nicholson, R. Banares-Alcantara, T. Kadir, and M. Brady, "Bayesian Networks for Clinical Decision Support in Lung Cancer Care," *PLoS ONE*, vol. 8, no. 12, p. e82349, 2013.
- [77] B. Yet, K. Bastani, H. Raharjo, S. Lifvergren, W. Marsh, and B. Bergman, "Decision support system for Warfarin therapy management using Bayesian networks," *Decision Support Systems*, vol. 55, no. 2, pp. 488–498, 2013.
- [78] O. Alagoz, H. Hsu, A. J. Schaefer, and M. S. Roberts, "Markov Decision Processes: A Tool for Sequential Decision Making under Uncertainty," *Medical Decision Making*, vol. 30, no. 4, pp. 474–483, 2010.
- [79] A. J. Schaefer, M. D. Bailey, S. M. Shechter, and M. S. Roberts, "Modeling medical treatment using markov decision processes," in *Operations research and health care*, pp. 593–612, Springer, 2005.
- [80] J.-H. Ahn and J. C. Hornberger, "Involving patients in the cadaveric kidney transplant allocation process: A decision-theoretic perspective," *Management Science*, vol. 42, no. 5, p. 629–641, 1996.
- [81] O. Alagoz, L. M. Maillart, A. J. Schaefer, and M. S. Roberts, "The optimal timing of living-donor liver transplantation," *Management Science*, vol. 50, no. 10, pp. 1420–1430, 2004.

- [82] P. Magni, S. Quaglini, M. Marchetti, and G. Barosi, "Deciding when to intervene: a markov decision process approach," *International Journal of Medical Informatics*, vol. 60, no. 3, pp. 237–253, 2000.
- [83] S. M. Shechter, M. D. Bailey, A. J. Schaefer, and M. S. Roberts, "The optimal time to initiate hiv therapy under ordered health states," *Operations Research*, vol. 56, no. 1, p. 20–33, 2008.
- [84] B. T. Denton, M. Kurt, N. D. Shah, S. C. Bryant, and S. A. Smith, "Optimizing the start time of statin therapy for patients with diabetes," *Medical Decision Making*, vol. 29, no. 3, p. 351–367, 2009.
- [85] J. Chhatwal, O. Alagoz, and E. S. Burnside, "Optimal breast biopsy decision-making based on mammographic features and demographic factors," *Operations research*, vol. 58, no. 6, pp. 1577–1591, 2010.
- [86] T. R. Hylan, M. Von Korff, K. Saunders, E. Masters, R. E. Palmer, D. Carrell, D. Cronkite, J. Mardekian, and D. Gross, "Automated prediction of risk for problem opioid use in a primary care setting," *The Journal of Pain*, vol. 16, no. 4, pp. 380–387, 2015.
- [87] R. Vunikili, B. S. Glicksberg, K. W. Johnson, J. Dudley, L. Subramanian, and S. Khader, "Predictive modeling of susceptibility to substance abuse, mortality and drug-drug interactions in opioid patients," *bioRxiv*, p. 506451, 2018.
- [88] Z. Che, J. S. Sauver, H. Liu, and Y. Liu, "Deep learning solutions for classifying patients on opioid use," in *AMIA Annual Symposium Proceedings*, vol. 2017, p. 525, American Medical Informatics Association, 2017.
- [89] R. J. Ellis, Z. Wang, N. Genes, and A. Ma'ayan, "Predicting opioid dependence from electronic health records with machine learning," *BioData Mining*, vol. 12, p. 3, Jan 2019.
- [90] M. M. Ghassemi, S. E. Richter, I. M. Eche, T. W. Chen, J. Danziger, and L. A. Celi, "A data-driven approach to optimized medication dosing: a focus on heparin," *Intensive Care Medicine*, vol. 40, no. 9, pp. 1332–1339, 2014.
- [91] S. Nemati, M. M. Ghassemi, and G. D. Clifford, "Optimal medication dosing from suboptimal clinical examples: A deep reinforcement learning approach," in *2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pp. 2978–2981, IEEE, 2016.
- [92] H. Bastani and M. Bayati, "Online decision making with high-dimensional covariates," *Operations Research*, vol. 68, no. 1, 2019.
- [93] M. V. Hill, M. L. McMahon, R. S. Stucke, and R. J. Barth, "Wide variation and excessive dosage of opioid prescriptions for common general surgical procedures," *Annals of Surgery*, vol. 265, no. 4, pp. 709–714, 2017.

- [94] M. C. Bicket, G. A. Brat, S. Hutfless, C. L. Wu, S. A. Nesbit, and G. C. Alexander, "Optimizing opioid prescribing and pain treatment for surgery: Review and conceptual framework," *American Journal of Health-System Pharmacy*, vol. 76, no. 18, pp. 1403–1412, 2019.
- [95] R. A. Young, S. K. Burge, K. A. Kumar, J. M. Wilson, and D. F. Ortiz, "A Time-Motion Study of Primary Care Physicians' Work in the Electronic Health Record Era," *Family Medicine*, vol. 50, no. 2, pp. 91–99, 2018.
- [96] L. Pizziferri, A. F. Kittler, L. A. Volk, M. M. Honour, S. Gupta, S. Wang, T. Wang, M. Lippincott, Q. Li, and D. W. Bates, "Primary care physician time utilization before and after implementation of an electronic health record: A time-motion study," *Journal of Biomedical Informatics*, vol. 38, no. 3, pp. 176–188, 2005.
- [97] J. W. Beasley, T. B. Wetterneck, J. Temte, J. A. Lapin, P. Smith, A. J. Rivera-Rodriguez, and B.-T. Karsh, "Information Chaos in Primary Care: Implications for Physician Performance and Patient Safety," *The Journal of the American Board of Family Medicine*, vol. 24, no. 6, pp. 745–751, 2011.
- [98] J. Bae and W. E. Encinosa, "National estimates of the impact of electronic health records on the workload of primary care physicians," *BMC Health Services Research*, vol. 16, no. 1, p. 172, 2016.
- [99] T. B. Wetterneck, J. A. Lapin, D. J. Krueger, G. T. Holman, J. W. Beasley, and B.-T. Karsh, "Development of a primary care physician task list to evaluate clinic visit workflow," *BMJ Quality & Safety*, vol. 21, no. 1, p. 47, 2012.
- [100] G. T. Holman, J. W. Beasley, B.-T. Karsh, J. A. Stone, P. D. Smith, and T. B. Wetterneck, "The myth of standardized workflow in primary care," *Journal of the American Medical Informatics Association*, vol. 23, no. 1, pp. 29–37, 2016.
- [101] P. Mishra, J. C. Kiang, and R. W. Grant, "Association of Medical Scribes in Primary Care With Physician Workflow and Patient Experience," *JAMA Internal Medicine*, vol. 178, no. 11, p. 1467, 2018.
- [102] M. L. Brandeau, F. Sainfort, and W. P. Pierskalla, *Operations research and health care: a handbook of methods and applications*, vol. 70. Springer Science & Business Media, 2004.
- [103] R. Hall, *Patient flow: Reducing delay in healthcare delivery*. Springer, 2006.
- [104] H. Yang and E. K. Lee, *Healthcare Analytics: from data to knowledge to healthcare improvement*. Wiley, 2018.
- [105] Y. Yih, *Handbook of Healthcare Delivery Systems*. CRC Press, 2016.
- [106] J. Li, N. Kong, X. Xie, Y. Teng, N. Kong, and T. Reimer, "Stochastic Modeling and Analytics in Healthcare Delivery Systems," pp. 253–280, 2017.

- [107] S. Fomundam and J. Herrmann, "A survey of queuing theory applications in health-care," in *ISR Technical Report 24*, 2007.
- [108] D. Gupta and B. Denton, "Appointment scheduling in health care: Challenges and opportunities," *IIE Transactions*, vol. 40, no. 9, pp. 800–819, 2008.
- [109] B. T. Denton, O. Alagoz, A. Holder, and E. K. Lee, "Medical decision making: open research challenges," *IIE Transactions on Healthcare Systems Engineering*, vol. 1, no. 3, pp. 161–167, 2011.
- [110] H. Xie, T. J. Chausalet, and P. H. Millard, "A continuous time Markov model for the length of stay of elderly people in institutional long-term care," *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, vol. 168, no. 1, pp. 51–61, 2005.
- [111] M. J. Côté and W. E. Stein, "A stochastic model for a visit to the doctor's office," *Mathematical and Computer Modelling*, vol. 45, no. 3-4, pp. 309–323, 2007.
- [112] J. Wang, S. Quan, J. Li, and A. M. Hollis, "Modeling and analysis of work flow and staffing level in a computed tomography division of University of Wisconsin Medical Foundation," *Health Care Management Science*, vol. 15, no. 2, pp. 108–120, 2012.
- [113] X. Zhong, J. Song, J. Li, S. M. Ertl, and L. Fiedler, "Design and analysis of gastroenterology (GI) clinic in Digestive Health Center of University of Wisconsin Health," *Flexible Services and Manufacturing Journal*, vol. 28, no. 1-2, pp. 90–119, 2016.
- [114] X. Zhong, J. Li, S. M. Ertl, C. Hassemer, and L. Fiedler, "A System-Theoretic Approach to Modeling and Analysis of Mammography Testing Process," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 46, no. 1, pp. 126–138, 2016.
- [115] X. Zhong, H. K. Lee, M. Williams, S. Kraft, J. Sleeth, R. Welnick, L. Hauschild, and J. Li, "Workload balancing: staffing ratio analysis for primary care redesign," *Flexible Services and Manufacturing Journal*, vol. 30, no. 1-2, pp. 6–29, 2018.
- [116] H. K. Lee, X. Zhong, J. Li, A. J. Musa, and P. A. Bain, "Joint visit in primary care clinics: Modeling, analysis, and an application study," *IIE Transactions on Healthcare Systems Engineering*, vol. 8, no. 2, pp. 93–109, 2017.
- [117] L. V. Green and S. Savin, "Reducing Delays for Medical Appointments: A Queueing Approach," *Operations Research*, vol. 56, no. 6, pp. 1526–1538, 2008.
- [118] L. Jiang and R. E. Giachetti, "A queueing network model to analyze the impact of parallelization of care on patient cycle time," *Health Care Management Science*, vol. 11, no. 3, pp. 248–261, 2008.
- [119] N. Liu and T. D'Aunno, "The Productivity and Cost-Efficiency of Models for Involving Nurse Practitioners in Primary Care: A Perspective from Queueing Analysis," *Health Services Research*, vol. 47, no. 2, pp. 594–613, 2012.

- [120] N. Liu, S. R. Finkelstein, and L. Poghosyan, "A new model for nurse practitioner utilization in primary care," *Health Care Management Review*, vol. 39, no. 1, pp. 10–20, 2014.
- [121] B. Zeng, H. Zhao, and M. Lawley, "The impact of overbooking on primary care patient no-show," *IIE Transactions on Healthcare Systems Engineering*, vol. 3, no. 3, pp. 147–170, 2013.
- [122] X. Zhong, J. Li, P. A. Bain, and A. J. Musa, "Electronic Visits in Primary Care: Modeling, Analysis, and Scheduling Policies," *IEEE Transactions on Automation Science and Engineering*, vol. 14, no. 3, pp. 1451–1466, 2017.
- [123] X. Zhong, P. Hoonakker, P. A. Bain, A. J. Musa, and J. Li, "The impact of e-visits on patient access to primary care," *Health Care Management Science*, vol. 21, no. 4, pp. 475–491, 2018.
- [124] N. Chen, X. Xie, P. A. Bain, M. P. Mundt, L. Zheng, and J. Li, "An Analytical Framework for Modeling, Analysis, and Improvement of Team Communication and Collaboration Process in Primary Care Clinics," *IEEE Transactions on Automation Science and Engineering*, vol. 16, no. 3, pp. 1148–1162, 2019.
- [125] B. Yet, Z. B. Perkins, T. E. Rasmussen, N. R. Tai, and D. W. R. Marsh, "Combining data and meta-analysis to build Bayesian networks for clinical decision support," *Journal of Biomedical Informatics*, vol. 52, pp. 373–385, 2014.
- [126] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. Morgan Kaufmann Publishers, 1997.
- [127] G. C. Sakellaropoulos and G. C. Nikiforidis, "Development of a bayesian network for the prognosis of head injuries using graphical model selection techniques," *Methods of Information in Medicine*, vol. 38, no. 1, pp. 37–42, 1999.
- [128] D. Koller and N. Friedman, *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- [129] J. Pearl, *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2000.
- [130] C. F. Vogelmeier, G. J. Criner, F. J. Martinez, A. Anzueto, P. J. Barnes, J. Bourbeau, B. R. Celli, R. Chen, M. Decramer, L. M. Fabbri, P. Frith, D. M. G. Halpin, M. V. López Varela, M. Nishimura, N. Roche, R. Rodriguez-Roisin, D. D. Sin, D. Singh, R. Stockley, J. Vestbo, J. A. Wedzicha, and A. Agustí, "Global strategy for the diagnosis, management, and prevention of chronic obstructive lung disease 2017 report. GOLD executive summary," *American Journal of Respiratory and Critical Care Medicine*, vol. 195, no. 5, pp. 557–582, 2017.
- [131] D. Bell, W. Liu, J. Cheng, R. Greiner, and J. Kelly, "Learning bayesian networks from data: An information-theory based approach," *Artificial Intelligence*, vol. 137, no. 1-2, pp. 43–90, 2002.

- [132] C. Chow and C. Liu, "Approximating discrete probability distributions with dependence trees," *IEEE Transactions on Information Theory*, vol. 14, no. 3, pp. 462–467, 1968.
- [133] G. F. Cooper and E. Herskovits, "A bayesian method for the induction of probabilistic networks from data," *Machine Learning*, vol. 9, no. 4, pp. 309–347, 1992.
- [134] D. Chickering, D. Geiger, and D. Heckerman, "Learning bayesian networks: search methods and experimental results," in *Proceedings of the 5th Conference on Artificial Intelligence and Statistics*, pp. 112–128, 1995.
- [135] F. Glover and M. Laguna, "Modern heuristic techniques for combinatorial problems," ch. Tabu Search, pp. 70–150, New York, NY, USA: John Wiley & Sons, Inc., 1993.
- [136] D. Heckerman, "A bayesian approach to learning causal networks," in *Proceedings of the 11th Conference on Uncertainty in Artificial Intelligence*, pp. 285–295, Morgan Kaufmann Publishers Inc., 1995.
- [137] M. Scutari, "Learning Bayesian networks with the bnlearn R package," *Journal of Statistical Software*, vol. 35, no. 3, pp. 1–22, 2010.
- [138] D. Barber, *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.
- [139] G. F. Cooper and E. Herskovits, "A Bayesian method for the induction of probabilistic networks from data," *Machine Learning*, vol. 9, no. 4, pp. 309–347, 1992.
- [140] D. Heckerman, "A tutorial on learning with Bayesian networks," in *Innovations in Bayesian Networks*, pp. 33–82, 2008.
- [141] S. J. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*. Pearson Education, 2 ed., 2003.
- [142] F. Glover, "Future paths for integer programming and links to artificial intelligence," *Computers & Operations Research*, vol. 13, pp. 533–549, May 1986.
- [143] G. Schwarz, "Estimating the dimension of a model," *The Annals of Statistics*, vol. 6, no. 2, pp. 461–464, 1978.
- [144] M. B. Sesen, A. E. Nicholson, R. Banares-Alcantara, T. Kadir, and M. Brady, "Bayesian networks for clinical decision support in lung cancer care," *PLoS ONE*, vol. 8, no. 12, 2013.
- [145] M. Boule, "Khiops: A statistical discretization method of continuous attributes," *Machine learning*, vol. 55, no. 1, pp. 53–69, 2004.
- [146] C. M. Bishop, *Pattern Recognition and Machine Learning, 5th Edition*. Information science and statistics, Springer, 2007.

- [147] O. Chapelle, B. Scholkopf, and A. Zien, *Semi-Supervised Learning*. The MIT Press, 1st ed., 2010.
- [148] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [149] G. Sharma, Y. Kuo, J. Freeman, D. Zhang, and J. Goodwin, "Outpatient follow-up visit and 30-day emergency department visit and readmission in patients hospitalized for chronic obstructive pulmonary disease," *Archives of Internal Medicine*, vol. 170, no. 18, pp. 1664–1670, 2010.
- [150] M. Puhan, M. Scharplatz, T. Troosters, and J. Steurer, "Respiratory rehabilitation after acute exacerbation of copd may reduce risk for readmission and mortality – a systematic review," *Respiratory Research*, vol. 6, no. 54, pp. 1–12, 2005.
- [151] S. Lee, S. Wang, P. A. Bain, C. Baker, T. Kundinger, C. Sommers, and J. Li, "Reducing copd readmissions: A causal bayesian network model," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 4046–4053, 2018.
- [152] C. Meek, "Causal inference and causal explanation with background knowledge," in *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence, UAI'95*, (San Francisco, CA, USA), pp. 403–410, Morgan Kaufmann Publishers Inc., 1995.
- [153] J. Pearl and T. Verma, "A theory of inferred causation," in *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning, KR'91*, (San Francisco, CA, USA), pp. 441–452, Morgan Kaufmann Publishers Inc., 1991.
- [154] M. K. Dwyer, C. M. Tumpowsky, N. L. Hiltz, J. Lee, W. L. Healy, and H. S. Bedair, "Characterization of Post-Operative Opioid Use Following Total Joint Arthroplasty," *The Journal of Arthroplasty*, vol. 33, no. 3, pp. 668–672, 2018.
- [155] D. W. Scott, *Multivariate Density Estimation: Theory, Practice, and Visualization*. John Wiley & Sons, 2015.
- [156] R. T. Marler and J. S. Arora, "The weighted sum method for multi-objective optimization: new insights," *Structural and Multidisciplinary Optimization*, vol. 41, no. 6, p. 853–862, 2010.
- [157] J. L. Cohon, *Multiobjective programming and planning*. Courier Corporation, 2004.
- [158] S. Kullback and R. A. Leibler, "On information and sufficiency," *The Annals of Mathematical Statistics*, vol. 22, no. 1, pp. 79–86, 1951.
- [159] J. S. Hahn, J. A. Bernstein, R. B. McKenzie, B. J. King, and C. A. Longhurst, "Rapid Implementation of Inpatient Electronic Physician Documentation at an Academic Hospital," *Applied Clinical Informatics*, vol. 03, no. 02, pp. 175–185, 2012.
- [160] J. Li and S. M. Meerkov, "On the coefficients of variation of uptime and downtime in manufacturing equipment," *Mathematical Problems in Engineering*, vol. 2005, 2005.

- [161] X. Zhong, M. Williams, J. Li, S. A. Kraft, and J. S. Sleeth, "Discrete-event simulation for primary care redesign: Review and a case study," in *Healthcare Analytics: From Data to Knowledge to Healthcare Improvement*, pp. 361–388, 2018.
- [162] J. Li and S. M. Meerkov, *Production systems engineering*. Springer Science & Business Media, 2008.