

INFORMS Journal on Optimization

Publication details, including instructions for authors and subscription information:
<http://pubsonline.informs.org>

Data-Driven Percentile Optimization for Multiclass Queueing Systems with Model Ambiguity: Theory and Application

Austin Bren, Soroush Saghaian

To cite this article:

Austin Bren, Soroush Saghaian (2019) Data-Driven Percentile Optimization for Multiclass Queueing Systems with Model Ambiguity: Theory and Application. INFORMS Journal on Optimization 1(4):267-287. <https://doi.org/10.1287/ijoo.2018.0007>

Full terms and conditions of use: <https://pubsonline.informs.org/Publications/Librarians-Portal/PubsOnLine-Terms-and-Conditions>

This article may be used only for the purposes of research, teaching, and/or private study. Commercial use or systematic downloading (by robots or other automatic processes) is prohibited without explicit Publisher approval, unless otherwise noted. For more information, contact permissions@informs.org.

The Publisher does not warrant or guarantee the article's accuracy, completeness, merchantability, fitness for a particular purpose, or non-infringement. Descriptions of, or references to, products or publications, or inclusion of an advertisement in this article, neither constitutes nor implies a guarantee, endorsement, or support of claims made of that product, publication, or service.

Copyright © 2019, INFORMS

Please scroll down for article—it is on subsequent pages



With 12,500 members from nearly 90 countries, INFORMS is the largest international association of operations research (O.R.) and analytics professionals and students. INFORMS provides unique networking and learning opportunities for individual professionals, and organizations of all types and sizes, to better understand and use O.R. and analytics tools and methods to transform strategic visions and achieve better outcomes.

For more information on INFORMS, its publications, membership, or meetings visit <http://www.informs.org>

Data-Driven Percentile Optimization for Multiclass Queueing Systems with Model Ambiguity: Theory and Application

Austin Bren,^a Soroush Saghafian^b

^a School of Computing, Informatics, and Decision Systems Engineering, Arizona State University, Tempe, Arizona 85287; ^b Harvard Kennedy School, Harvard University, Cambridge, Massachusetts 02138

Contact: austin.s.bren@asu.edu (AB); soroush_saghafian@hks.harvard.edu,  <http://orcid.org/0000-0002-9781-6561> (SS)

Received: September 29, 2017

Revised: March 10, 2018; July 13, 2018

Accepted: August 27, 2018

Published Online in Articles in Advance:
April 29, 2019

<https://doi.org/10.1287/ijoo.2018.0007>

Copyright: © 2019 INFORMS

Abstract. Multiclass queueing systems widely used in operations research and management typically experience ambiguity in real-world settings in the form of unknown parameters. For such systems, we incorporate robustness in the control policies by applying a data-driven percentile optimization technique that allows for (a) expressing a controller's optimism level toward ambiguity and (b) utilizing incoming data to learn the true system parameters. We show that the optimal policy under the percentile optimization objective is related to a closed-form, priority-based policy. We also identify connections between the optimal percentile optimization and $c\mu$ -like policies, which, in turn, enables us to establish effective but easy-to-use heuristics for implementation in complex systems. Using real-world data collected from a leading U.S. hospital, we also apply our approach to a hospital emergency department setting and demonstrate the benefits of using our framework for improving current patient flow policies.

Supplemental Material: The online appendices are available at <https://doi.org/10.1287/ijoo.2018.0007>.

Keywords: model ambiguity • data-driven optimization • ED operations

1. Introduction

Multiclass queueing systems require dynamic control in environments in which servers must process multiple types of jobs that vary with respect to holding costs, service rates, and other defining characteristics. These types of queueing systems are widely used to model call centers, hospitals, manufacturing lines, and service operations, in which elements in the queue can be classified based on differing levels of urgency, processing time, or other attributes. For example, in a hospital emergency department (ED), patients are classified through a triage system, which differentiates them based on their severity, medical complexity, or other conditions (see, e.g., Saghafian et al. 2012, 2014, and the references therein). Hence, a natural way to analyze ED patient flow is via a multiclass queueing system that separates patients based on their attributes.¹

In such systems, when all parameters are known, many well-established policies, such as the $c\mu$ rule, have been shown to be optimal for optimizing the system's performance (see, e.g., Van Mieghem 1995 and Buyukkoc et al. 1985). However, the assumption that all the model parameters are perfectly known is often unrealistic, especially in settings with little supporting data, inaugural system launch, or various other sources of ambiguity. A manager with incorrect parameter specifications may enforce policies that perform poorly or may not have confidence in using a policy that is obtained from a model with parameters that the manager does not fully trust. In an effort to combat such mistrust, we consider a form of model ambiguity caused by the ambiguity in parameters termed *parameter ambiguity* and develop strategies that directly take these into account.

Traditionally, robust optimization protects against parameter ambiguity by utilizing a minimax objective on an ambiguity set of parameters that are assumed to contain the true system parameters. However, this type of robustness (a) can result in overly pessimistic policies and (b) ignores the significant potential to learn about the true system parameters from data acquired both before and after system launch. Even when this pessimism is reduced by choosing tighter ambiguity sets, the policies generated are not capable of learning from incoming data. To avoid these deficiencies, we model parameter ambiguity via a partially observable Markov decision process (POMDP), an extension of Markov decision processes (MDPs), which allows for (a) imperfect state knowledge and (b) learning in a Bayesian manner. A POMDP supports the distribution of the underlying system parameters, known as the belief space, and updates this distribution to reflect received observations. This is ideal from a learning perspective; however, in a POMDP, the decision maker is assumed to have an initial prior belief, which is often a subjective value, guided by scarce data, error-prone expert opinion, intuition, or instinct. For these reasons, Bayesian critics distrust such learning mechanisms, citing the unreliability of the prior specification in real-world applications.²

To incorporate robustness to such a prior belief (hence, gaining robustness to parameter ambiguity), we integrate our POMDP model with a *percentile optimization* approach. Percentile optimization is traditionally used to avoid overly conservative policies by offering a certain level of performance over a percentage of the ambiguity set (see, e.g., Delage and Mannor 2010 and Nemirovski and Shapiro 2006). We extend percentile optimization to incorporate robustness to the belief about the model parameters rather than relying on a robustness generated directly from the parameters themselves. In this way, we investigate strategies in which the controller learns the main model parameters (e.g., unknown service rates) while simultaneously controlling the underlying system for superior performance, which contrasts with robust techniques that only focus on parameter ambiguity sets. Thus, our framework allows generating policies that are robust to parameter ambiguities (considering a manager's pessimism level) while simultaneously learning about the true model from data/observation of the system's performance in a Bayesian manner.

Our main contributions stem from extending the robust percentile optimization approach for integration with POMDPs. We find that the percentile optimization objective reduces to the minimax and minimin objectives when the optimism level is set to its lowest and highest values, respectively, and show that the optimal policies under these objectives are myopic $c\mu$ priority policies. Understanding the nonrobust problem (which assumes a specified initial belief) proves to be essential in finding robust policies in which the belief is subject to ambiguity. We find that optimal robust policies can be formed using specific nonrobust policies via a geometric structure known as the *convex floating body*. Therefore, to solve the robust percentile problem, we first solve the nonrobust problem that has a known initial belief. As the rate of observations increases, we find that a priority-based policy that acts as an extension of the well-known $c\mu$ rule becomes asymptotically optimal to the nonrobust problem. This policy, which we term $Ec\mu$, is myopic and prioritizes the class with the largest expected $c\mu$ value. The proposed $Ec\mu$ policy utilizes incoming data for learning (unlike the traditional $c\mu$ rule) and is extremely simple to implement.

Because of its foundation in POMDPs, the robust framework we consider is computationally ambitious and necessitates finding tractable methods for implementation. Using the analytical insights gained from the connection between nonrobust and robust policies, constraints via the convex floating body, and the relation of $Ec\mu$ to the nonrobust objective, we develop a heuristic for the robust problem that (a) is highly scalable to large problem instances and (b) shows strong performance in extensive simulation experiments. We also develop analytical bounds to the nonrobust problem based on queueing systems with fully known parameters. These bounds are (a) tight under a variety of conditions and (b) can be used to more effectively compute optimal robust policies. Furthermore, because the bounds are based on nonlearning policies, they can be computed in an efficient manner.

Finally, we demonstrate the benefits of our approach in a real-world setting by utilizing data that we have collected from a leading U.S. hospital and by establishing the advantages of using our framework in improving the current ED patient-flow policies. Our percentile optimization framework is the first study in the literature to yield data-driven policies for use in EDs that hedge against parameter ambiguity. We find that highly congested EDs are well suited to our percentile optimization framework, especially in geographical areas with uncertain/unstable patient population characteristics. Additionally, our approach explicitly avoids overly conservative policies that focus only on the worst-case scenarios. As a result, we find that percentile optimization performs well over a large spectrum of optimism/pessimism. In particular, our simulations calibrated with hospital data suggest that, by using our approach, an ED manager can significantly improve performance regardless of the manager's disposition.

The rest of the paper is organized as follows. In Section 2, we provide a literature review of the related studies. Section 3 introduces the nonrobust, continuous-time formulation of our problem, which is uniformized into a discrete-time problem in Section 3.1 and lays the foundation for the percentile framework developed in Section 3.2. We provide the majority of our analytical insights in Sections 4 and 5, in which we establish optimal policies for the nonrobust and robust formulations and identify upper/lower bound results. Section 6 introduces a heuristic to the robust problem that is rooted in the analytical insights generated from Section 4. In Section 7, we present various numerical experiments, discuss the application of our work for improving patient flow in EDs, and use real-world data obtained from a leading U.S. hospital to evaluate the potential benefits of our approach. Finally, in Section 8, we present our concluding remarks.

2. Literature Review

The literature surrounding multiclass queueing systems aims to analyze complex structures and discover their optimal control policies, such as the $c\mu$ policy and its variations (see, e.g., Buyukkoc et al. 1985, Van Mieghem 1995, and Saghaian and Veatch 2016 and the references therein). A common tool used to analyze and control such systems is MDPs. However, their use is limited to the unrealistic case in which the decision maker is assumed to completely know all the parameters of the model (e.g., service rates). Most notably, this includes a perfect knowledge assumption of the transition matrices that guide a system's state transitions. This assumption can be

problematic in various practical applications in which service rates (or other parameters) are not perfectly known. Mannor et al. (2007) and Nilim and El Ghaoui (2005) find that small changes in such parameters can result in significant differences in decision-making strategies. However, a synthesis of most studies on dynamic control in queueing systems indicates the use of tools that heavily rely on a full knowledge about the system's parameters. This is despite the fact that, in practice, such parameters are typically unknown and often hard to estimate.

Robust methods applied to queueing models are largely involved with reducing the computational burden of characterizing queueing metrics and policies. Su (2006) studies a fluid approximation of a multiclass queueing model's holding cost under a robust paradigm established by Bertsimas et al. (2004) and Bertsimas and Sim (2004). Bertsimas et al. (2011) focus on finding bounds for performance measures through a method rooted in robust optimization and study the performance of this method on tandem and multiclass single-server queueing networks. Jain et al. (2010) find that a queueing network with control over traffic intensities has a simple threshold-type policy under a robust objective. For more recent studies on robust techniques used in queueing systems, we refer to Pedarsani et al. (2014), Bandi and Bertsimas (2012), and Bandi et al. (2015) and the references therein. This stream of research is mainly aimed at increasing tractability by focusing on worst-case (i.e., fully pessimistic) scenarios and establishing related performance metrics. Unlike this stream, our goal is to provide policies that (a) are more optimistic (i.e., less conservative) and (b) incorporate learning from online, system-run data/observations.

Adding robustness when facing parameter ambiguity is a topic of significant interest to a variety of fields, including economics, operations research/management, computer science, and decision theory among others. Typically, robustness in MDPs is added using a minimax objective because this often results in tractable analyses as shown in Nilim and El Ghaoui (2005) and Iyengar (2005) and the references therein. Other studies, such as Chen and Farias (2013), deal with ambiguities by considering policies that offer guarantees on expected performance. Still other methods of incorporating robustness include regret minimization (Lim et al. 2012), relative entropy (Bagnell et al. 2001), and martingale-based approaches (Hansen and Sargent 2007) that provide less conservative and, hence, potentially more realistic alternatives to minimax techniques. In particular, Delage and Mannor (2010) identify a robust approach applied to MDPs called percentile optimization that effectively avoids overconservatism (see also Nemirovski and Shapiro 2006 and Wiesemann et al. 2013 for related studies). Instead of finding policies that are tailored to work well in worst-case scenarios, the percentile optimization method finds policies that maximize performance with respect to a level of belief about the true parameters for a given level of optimism.

Chow et al. (2018) also utilize this type of robustness to develop risk-constrained policies for MDPs. However, a significant deficit in current percentile optimization approaches is the lack of ability to *learn* about the true parameters over time. Delage and Mannor (2007) work to fill this gap via a similar formulation to our approach, and find second-order approximations to MDPs that experience transition parameter uncertainty. However, the Dirichlet-type uncertainty assumed in transition parameters does not fit our queueing problem, and in our work, we extend the percentile optimization approach with respect to ambiguity in the initial belief. Thus, system data/observations can be used for learning the true operational model, and as we show, this ability to learn itself adds a strong layer of robustness for controlling queueing systems (e.g., hospital patient flows) that face parameter ambiguity. Learning to overcome ambiguities is also discussed in Bassamboo and Zeevi (2009), which models a call center application using a data-driven technique. However, their work (a) does not include any notion of robustness and (b) focuses on near-optimal policies with performance bounds. Our work differs in modeling approach by our joint focus on learning and robustness and in methodology by our contributions in characterizing the exact optimal policies.

Data-driven parameter learning has been incorporated in POMDPs: Ross et al. (2011) explores a finite-horizon POMDP model that updates a posterior of its parameter belief in a Bayesian manner, and Thrun (1999) investigates a POMDP in continuous action and state spaces that relies on particle-filtering techniques to determine the belief state. Unlike learning mechanisms, robust methods are almost nonexistent in POMDP frameworks. Osogami (2015) shows that traditional minimax approaches with convex ambiguity sets can be extended to POMDPs while still retaining their structural features (such as convexity). In a new approach, Saghaian (2018) extends POMDPs to a new class termed ambiguous POMDPs (APOMDPs), which incorporates ambiguity in transition and observation probabilities in a robust fashion. The robustness in Saghaian (2018) is achieved by considering α -maximin (α -MEU) preferences and by incorporating the decision maker's temperament toward model ambiguity. Different from the APOMDP approach of Saghaian (2018), we utilize a percentile optimization objective to hedge against ambiguities.

3. The Multiclass Queueing Control Problem with Parameter Ambiguity

We begin by considering a continuous-time, multiclass, queueing-control problem with preemption, in which a single server is responsible for serving n classes of customers over an infinite time horizon. Unlike the traditional version of this model, we assume the controller does not know the main parameters of the system and, hence, is

faced with parameter ambiguity. We focus on the case in which the ambiguity is on service rates. To this end, we start by excluding dynamic arrivals to the system and instead consider a *clearing system*³ version of the problem. We relax this assumption in Sections 7 and A.4 by allowing for dynamic arrivals and find that many of our major results are transferable from the clearing system. Our general approach can also be used for systems in which arrival rates or other parameters are ambiguous by modifying the underlying dynamic program to include these components along with their learning mechanisms. However, this appears to increase the problem's complexity without providing additional insights.

With $\mathcal{N} = \{1, \dots, n\}$ denoting the set of customer classes, we assume each customer of class $i \in \mathcal{N}$ accrues a cost $\hat{c}_i > 0$ for each unit of time spent in the system. Let $\hat{\mathbf{c}} = (\hat{c}_1, \hat{c}_2, \dots, \hat{c}_n)$ be the cost vector, $\alpha \in (0, \infty)$ the discount rate, and $\mathbf{X}(t) = (X_1(t), X_2(t), \dots, X_n(t))$ the vector of the number of customers in the system, where $X_i(t)$ is number of class i customers in the system at time t . In line with many robust approaches, we begin by outlining an ambiguity set (i.e., a “cloud” of models) that is assumed to include the true model. To this end, and for tractability, we assume service times for each class are independent and identically distributed exponential⁴ random variables with unknown rates for each class. The true service rate for each class $i \in \mathcal{N}$ is chosen by nature at time $t = 0$ and lies within ambiguity set $\mathcal{M}_i = \{\hat{\mu}_{i,1}, \dots, \hat{\mu}_{i,m_i}\}$. We further assume that service times for different classes are independent. For future notational convenience, we let $\mathcal{J}_i = \{1, \dots, m_i\}$. Throughout the paper, we assume $m_i \in \mathbb{N}$, and $\hat{\mu}_{i,j} \neq \hat{\mu}_{i,k}$ for each $i \in \mathcal{N}$ and distinct $j, k \in \mathcal{J}_i$. Although the ambiguity sets \mathcal{M}_i are discrete, the continuous case can be approximated arbitrarily closely by increasing the number of potential service rates m_i to make the mesh size of \mathcal{M}_i close to zero.

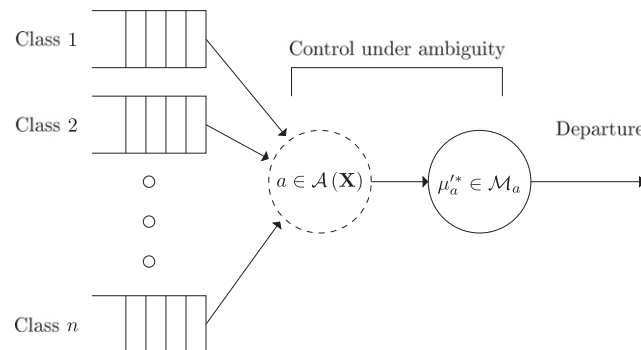
Over time, the controller can learn the true service rates by observing the process history, which includes all previous service durations, control actions, and observations of service completions. For Markovian systems with incomplete information, it has been shown in Bertsekas (1995) that the Bayesian belief on the unknown parameters with respect to the observed process history is a *sufficient statistic*. We let \mathcal{B} be the set of all such sufficient statistics, that is, the set of possible belief distributions on the system's service parameters. Letting $m = \sum_{i \in \mathcal{N}} m_i$, each $\mathbf{b} \in \mathcal{B}$ is an m -dimensional vector of the form $\mathbf{b} = (b_{1,1}, b_{1,2}, \dots, b_{1,m_1}, b_{2,1}, \dots, b_{n,m_n})$ with the condition that each $b_{i,j} \geq 0$ and that $\sum_{j=1}^{m_i} b_{i,j} = 1$ for each $i \in \mathcal{N}$. In this setting, if $\hat{\mu}_i^* \in \mathcal{M}_i$ is the true (unknown) service rate for class $i \in \mathcal{N}$, $\Pr(\hat{\mu}_{i,j} = \hat{\mu}_i^* | \mathbf{b}) = b_{i,j}$. We further assume that the observation made after serving one class does not affect the belief about another. This is aligned with the assumption that service time of one class is independent of that of another class.

To find policies that optimally prescribe which customer class the server should serve at any time given (a) the available information summarized in the current belief about the service rates and (b) the number of customers in each queue, it is known that one can restrict attention to policies that are deterministic, stationary, and Markovian (see, e.g., Sondik 1971, Smallwood and Sondik 1973, and Bertsekas 1995). Consequently, an admissible non-anticipative policy π maps the current belief and queue-length information (information state) to the set of actions: $\pi : \mathbb{Z}_+^n \times \mathcal{B} \rightarrow \mathcal{N} \cup \{0\}$ with the additional condition that π can serve only customer classes that have nonempty queues and serves the fictitious class “0” when the server is idled (e.g., when all the queues are empty). Our model described is schematically illustrated in Figure 1.

We let Π denote the set of all admissible policies and $\mathbf{X}^\pi(t) = (X_1^\pi(t), X_2^\pi(t), \dots, X_n^\pi(t)) \in \mathbb{Z}_+^n$ represent the number of customers in the system under policy $\pi \in \Pi$ at time t . In Online Appendix B, Lemma 17 shows that idling the server when at least one customer class queue is nonempty is always suboptimal; hence, we consider only nonidling policies in our analysis. For a given policy π , the expected discounted *true cost* the system experiences is

$$\mathbb{E}_\pi \left[\int_{t=0}^{\infty} e^{-\alpha t} \hat{\mathbf{c}} \mathbf{X}^\pi(t)^\top dt | \mathbf{X}(0) \right],$$

Figure 1. The Server Serves a Class a Customer with an Unknown Rate Belonging to Ambiguity Set \mathcal{M}_a



given the true transition parameters chosen by nature at time $t = 0$, where the notation “ T ” represents transpose and E_π is expectation with respect to the probability measure induced by π . However, because the controller does not know the true transition matrix (as service rates are unknown), we are interested in the expected cost with respect to the controller’s belief:

$$J^\pi(\mathbf{X}(0), \mathbf{b}(0)) = E_{\pi, \mathbf{b}(0)} \left[\int_{t=0}^{\infty} e^{-\alpha t} \hat{\mathbf{c}} \mathbf{X}^\pi(t)^T dt | \mathbf{X}(0) \right], \quad (1)$$

where $E_{\pi, \mathbf{b}(0)}$ denotes expectation with respect to both the initial belief $\mathbf{b}(0)$ and π . We refer to $J^\pi(\mathbf{X}(0), \mathbf{b}(0))$ as the *nonrobust cost* under policy π because it assumes a perfectly assigned $\mathbf{b}(0)$ (which is inevitably hard to quantify for any decision maker who is faced with model ambiguity). The optimal nonrobust cost is then given by $J(\mathbf{X}(0), \mathbf{b}(0)) = \inf_{\pi \in \Pi} J^\pi(\mathbf{X}(0), \mathbf{b}(0))$. In what follows, we first use uniformization to work with the discrete-time model of the nonrobust scenario, in which the initial belief is given. We then adopt percentile optimization to enable the decision maker/controller to reduce reliance on $\mathbf{b}(0)$ and thereby make robust decisions.

3.1. A Discrete-Time Nonrobust Framework

The continuous-time Markov chain $\{\mathbf{X}^\pi(t) : t \geq 0\}$ can be converted to a discrete-time equivalent using the well-known uniformization technique (Lippman 1975). Following this method, we first select a *uniformized* exponentially distributed random variable ξ with a rate $\psi > \max_{i \in \mathcal{N}, j \in \mathcal{J}_i} \hat{\mu}_{i,j}$, which serves as our rate of observations made as follows. If the server completes service to a customer of class i within a uniformized unit of time (i.e., at the end of each period), an observation indicating the “successful” service to class i is recorded. Otherwise, if no service completion is observed within this time, an observation is recorded indicating an “incomplete” service to class i . We note that this *uniformization rate* ψ may be arbitrarily large so as to approximate continuous observations.

We let σ be the Bayesian learning operator such that $\sigma(\mathbf{b}, a, \theta)$ is an m -dimensional vector representing the updated belief after taking action a and receiving observation θ when the prior belief is \mathbf{b} . Because there are only two outcomes for observations for any given action, we let “+” signify an observed service completion (“success”) during the uniformized time period and “−” represent an incomplete service (“failure”) in that period. In this setting, we use a discrete-time dynamic program with uniformized parameters $\mu_{i,j} = \hat{\mu}_{i,j}/\psi$. For notational convenience, we let $E[\mu_i | \mathbf{b}] = \sum_{j=1}^{m_i} \mu_{i,j} b_{i,j}$ be the expected service transition probability of class $i \in \mathcal{N}$ given belief \mathbf{b} . In this way, the Bayesian learning operator updates belief \mathbf{b} with components $b_{i,j}$ to belief $\bar{\mathbf{b}} = \sigma(\mathbf{b}, a, \theta)$ with components $\bar{b}_{i,j} = \sigma(\mathbf{b}, a, \theta)_{i,j}$, where $a, i \in \mathcal{N}, j \in \mathcal{J}_i$, and

$$\sigma(\mathbf{b}, a, +)_{i,j} = \begin{cases} \frac{\mu_{a,j} b_{a,j}}{\sum_{k=1}^{m_a} \mu_{a,k} b_{a,k}} = \frac{\mu_{a,j} b_{a,j}}{E[\mu_a | \mathbf{b}]} : i = a \\ b_{i,j} : i \neq a \end{cases} \quad (2)$$

for a successful service observation and

$$\sigma(\mathbf{b}, a, -)_{i,j} = \begin{cases} \frac{(1 - \mu_{a,j}) b_{a,j}}{\sum_{k=1}^{m_a} (1 - \mu_{a,k}) b_{a,k}} = \frac{(1 - \mu_{a,j}) b_{a,j}}{(1 - E[\mu_a | \mathbf{b}])} : i = a \\ b_{i,j} : i \neq a \end{cases} \quad (3)$$

for a failed service observation. Equations (2) and (3) are established because, under realized parameter $\mu_{a,j}$, the probability of successful service in a given period is $\mu_{a,j}$ and probability of incomplete service is $(1 - \mu_{a,j})$. With this and defining a discrete-time discounting factor $\beta = \frac{\psi}{\psi + \alpha}$ and instantaneous cost $\mathbf{c} \mathbf{X}^T = \frac{\hat{\mathbf{c}} \mathbf{X}^T}{\psi + \alpha}$ where \mathbf{X} is an n -dimensional vector representing queue lengths, we can identify the nonrobust optimal policy and the associated cost via the dynamic program

$$V_{t+1}(\mathbf{X}, \mathbf{b}) = \mathbf{c} \mathbf{X}^T + \beta \left[\min_{a \in \mathcal{A}(\mathbf{X})} \left\{ E[\mu_a | \mathbf{b}] V_t(\mathbf{X} - \mathbf{e}_a, \sigma(\mathbf{b}, a, +)) \right. \right. \\ \left. \left. + (1 - E[\mu_a | \mathbf{b}]) V_t(\mathbf{X}, \sigma(\mathbf{b}, a, -)) \right\} \right], \quad (4)$$

with the terminal condition $V_0(\mathbf{X}, \mathbf{b}) = \mathbf{c} \mathbf{X}^T$, where $\mathcal{A}(\mathbf{X}) = \{i \in \mathcal{N} \cup \{0\} | X_i \neq 0\}$ is the set of admissible actions. In this setting, taking the limit as $t \rightarrow \infty$, we define $V(\mathbf{X}, \mathbf{b}) = \lim_{t \rightarrow \infty} V_t(\mathbf{X}, \mathbf{b})$ and note that $V(\mathbf{X}, \mathbf{b}) = \inf_{\pi \in \Pi} J^\pi(\mathbf{X}, \mathbf{b})$ (see Lemma 11 in Online Appendix B for a rigorous treatment), where $J^\pi(\mathbf{X}, \mathbf{b})$ is defined in (1). To account for evaluating nonoptimal policies, we let $V_{t+1}^\pi(\mathbf{X}, \mathbf{b})$ be a value function similar to that of the dynamic program (4) with the minimization operator replaced by serving the class prescribed by policy π . Likewise, we let $V^\pi(\mathbf{X}, \mathbf{b}) = \lim_{t \rightarrow \infty} V_t^\pi(\mathbf{X}, \mathbf{b})$ be the infinite-horizon dynamic program value function under policy π .

3.2. Gaining Robustness via Percentile Optimization

Because the controller is facing ambiguity with respect to the true model, the controller may distrust the initial prior on the cloud of models, $\mathbf{b}(0)$. The specification of $\mathbf{b}(0)$ is subject to model sensitivities, especially in applications in which there is little or highly variable data to perfectly quantify it. Often, the selection of a prior is a process that requires sussing out probabilities and parameter values from experts in the field, which can be a highly subjective and inaccurate task.⁵

In traditional robust optimization, one would choose a policy assuming that nature, being an antagonistic character, picks the worst-case initial belief vector $\mathbf{b}(0)$ for a chosen policy. Hence, the traditional minimax robust objective can be defined by first considering the worst-case cost under a policy $\pi \in \Pi$:

$$R^\pi(\mathbf{X}) = \max_{\mathbf{b} \in \mathcal{B}} V^\pi(\mathbf{X}, \mathbf{b}).$$

The cost under the minimax robust objective is then $R(\mathbf{X}) = \inf_{\pi \in \Pi} R^\pi(\mathbf{X})$. In this setting, the controller assumes that nature will pick the transition parameters that result in the maximum cost for any given policy and chooses a policy that minimizes the cost of this worst-case outcome.

In sharp contrast to this type of robustness, which typically yields overly pessimistic control policies, is the overly optimistic minimin objective defined by

$$N^\pi(\mathbf{X}) = \min_{\mathbf{b} \in \mathcal{B}} V^\pi(\mathbf{X}, \mathbf{b}),$$

and $N(\mathbf{X}) = \inf_{\pi \in \Pi} N^\pi(\mathbf{X})$, under which the controller chooses a policy assuming nature picks the transition parameters resulting in the best-case cost for any given policy. In what follows, we first show that both minimax and minimin optimal policies are within the well-known class of $c\mu$ policies. Thus, they (a) are fully myopic and (b) have very simple forms.

Proposition 1 (Minimax/Minimin $c\mu$ Optimal Policies). *At any state (\mathbf{X}, \mathbf{b}) , optimal policies to the minimax and minimin objectives serve classes $\arg \max_{a \in \mathcal{A}(\mathbf{X})} (\min_{j \in \mathcal{J}_a} c_a \mu_{a,j})$ and $\arg \max_{a \in \mathcal{A}(\mathbf{X})} (\max_{j \in \mathcal{J}_a} c_a \mu_{a,j})$, respectively.*

Proposition 1 establishes that optimal policies under both minimax and minimin objectives are myopic priority disciplines (known as the $c\mu$ rule) with respect to the smallest and largest transition rates within the ambiguity set for each class, respectively. However, it should be noted that such policies (a) ignore the potential for learning from the system behavior and (b) only consider the potentially unrealistic extreme best- and worst-case scenarios and can perform poorly in real-world applications. To address this deficit, we next investigate how the percentile optimization approach provides a balancing alternative between these two extreme strategies while incorporating learning about the hidden probabilities associated with the true transition parameters (i.e., service rates).

To this end, for a given $\epsilon \in [0, 1]$, we define the percentile optimization program:

$$Y^\pi(\mathbf{X}, \epsilon) = \inf_{y_\epsilon \in [N^\pi(\mathbf{X}), R^\pi(\mathbf{X})]} y_\epsilon \quad (5)$$

$$\text{s.t. } \Pr(V^\pi(\mathbf{X}, \mathbf{B}) \leq y_\epsilon) \geq 1 - \epsilon, \quad (6)$$

and let $Y(\mathbf{X}, \epsilon) = \inf_{\pi \in \Pi} Y^\pi(\mathbf{X}, \epsilon)$ represent the optimal percentile objective. In (5), we impose that $N^\pi(\mathbf{X}) \leq y_\epsilon \leq R^\pi(\mathbf{X})$ so that the value of the objective is within the most optimistic and pessimistic values attainable for any given belief in accordance with the policy, hence enforcing “realizable” expected costs. The probability operator in (6) is defined with respect to a specified probability density $\mathbb{P}_{\mathbf{B}}$ over the prior belief space,⁶ where \mathbf{B} is a random variable whose realization is \mathbf{b} . The percentile optimization program (5) and (6) allows us to find a *chance-constrained* policy: it emphasizes policy performance over a portion of the belief space. We, thus, term the policy that is the solution under the optimal percentile objective as $(1 - \epsilon)\%$ chance-constrained policy. Intuitively, the smaller the ϵ , the more protection from poor parameter settings because the proportion of the belief space that performs worse than y_ϵ becomes smaller.

It is important to note that the percentile objective acts as a bridge between nonrobust and robust objectives; expressing a manager’s optimism level is a core ambition of this type of robustness. For instance, the chance-constrained policy reduces to the minimax and minimin policies when ϵ is zero and one, respectively.

Proposition 2 (Percentile/Minimax/Minimin Relationship). *The percentile objective, minimax, and minimin policies share the following relation:*

- i. If $\epsilon = 0$ and $\mathbb{P}_{\mathbf{B}}(\mathbf{b}) > 0$ for all $\mathbf{b} \in \mathcal{B}$, then the optimal policy and cost under both minimax and percentile objectives are the same.
- ii. If $\epsilon = 1$, then the optimal policy and cost under the minimin and percentile objectives are the same.

The additional condition $\mathbb{P}_{\mathbf{B}}(\mathbf{b}) > 0$ for all $\mathbf{b} \in \mathcal{B}$ in part (i) is necessary because $\mathbb{P}_{\mathbf{B}}$ with zeros may allow the percentile objective to “ignore” certain portions of the belief space while still satisfying constraint (6). For example, if $\mathbb{P}_{\mathbf{B}}$ yields a degenerate distribution with respect to a point \mathbf{b} , $Y(\mathbf{X}, 0) = V(\mathbf{X}, \mathbf{b})$.

4. Structure of Optimal Policies Under the Percentile Objective

Analyzing program (5) and (6) is inherently complex both analytically and computationally. However, we find that the solution to this program is linked to solving the nonrobust problem. Hence, we first consider the solution of the dynamic program (4), identify important characteristics of these solutions over the belief space, establish the link between nonrobust and robust policies, and finally work to characterize optimal percentile policies. In Section 6, we develop an easy-to-use heuristic based on these insights to facilitate tractable solutions.

As the observation rate increases, tending toward continuous observations, the nonrobust problem can be transferred to a multiarmed bandit (MAB) problem by noting that, (a) under any action, only the belief about transition parameters and number of customers in the served class (the “arms” of the MAB) change and (b) the “discounted cost” can be reinterpreted as “discounted savings” of the MAB because of our clearing system environment (for further discussion, see Lemma 3 in Online Appendix B). MAB problems are typically solved by indexing policies related to the expected savings in cost experienced through exclusively serving one class over time.

To take advantage of the aforementioned connection, we term the myopic policy that serves the class $a \in \mathcal{A}(\mathbf{X})$ with largest value of $c_a E[\mu_a | \mathbf{b}]$ the “ $Ec\mu$ ” policy. Thus, we denote $\pi^{c\mu}$ that serves $\arg \max_{a \in \mathcal{A}(\mathbf{X}(t))} c_a E[\mu_a | \mathbf{b}(t)]$ as the $Ec\mu$ policy. This policy can be viewed as an extension of the traditional $c\mu$ policy (often seen in the literature surrounding control of multiclass queueing systems) for queueing systems with ambiguous parameters. The expectation operator in this policy dynamically combines all the possible $c\mu$ values for each class based on the belief at time t . Because the nonrobust problem learns via belief state transitions in discrete increments based on the observation rate, the myopic $Ec\mu$ policy is not optimal in general. However, as this rate increases, the belief state transitions become smooth, leading to its asymptotic optimality as shown in the following theorem.

Theorem 1 (*$Ec\mu$ Asymptotic Optimality*). *The $Ec\mu$ policy $\pi^{c\mu}$ is asymptotically optimal for the nonrobust problem: $\lim_{\psi \rightarrow \infty} V^{\pi^{c\mu}}(\mathbf{X}, \mathbf{b}) = \lim_{\psi \rightarrow \infty} V(\mathbf{X}, \mathbf{b})$ for all $\mathbf{X} \in \mathbb{Z}_+^n$ and $\mathbf{b} \in \mathcal{B}$.*

Theorem 1 is surprising in its simplicity because problems based on POMDP formulations typically do not yield closed-form results. In contrast to the usual complexities, the asymptotic optimality of the $Ec\mu$ policy implies that the only information necessary to make decisions is the expected transition rates among nonempty queues. Therefore, transition rate variability only has an effect on parameter learning and not on the policy. Similarly, queue lengths are essentially irrelevant to the decision maker. However, the $Ec\mu$ policy features a momentum property; if the current action a prescribed by the policy yields enough successes so that $c_a E[\mu_a | \mathbf{b}]$ does not fall below the threshold defined by $c_{\hat{a}} E[\mu_{\hat{a}} | \mathbf{b}]$ of the next highest available class \hat{a} , the $Ec\mu$ policy will continue to serve class a regardless of the state of other classes. In turn, this means that the policy will not attempt to serve a class with smaller $c_a E[\mu_a | \mathbf{b}]$ until other classes with larger values have experienced a sufficient number of service failures or have cleared their queue. This property may run counterintuitive to the exploration-minded individual; even if a class has the potential to be endowed with a very large $c_a \mu_{a,j}$ value (under the realization of system parameters), this potential is only rated on the basis of its contribution to the expected service rate.

Another important property of the $Ec\mu$ policy is that, under mild conditions, $V^{\pi^{c\mu}}(\mathbf{X}, \mathbf{b})$ is piecewise-linear over the belief space (excluding beliefs near edges and faces of \mathcal{B}).⁷

Proposition 3 (*Piecewise-Linearity of the Approximate Nonrobust Value Function*). *Let \mathcal{B}' be any closed subset of \mathcal{B} such that for any $\mathbf{b} \in \mathcal{B}'$, $b_{i,j} > 0$ for all $i \in \mathcal{N}$, $j \in \mathcal{J}_i$. If $\min_{j \in \mathcal{J}_i} c_i \mu_{i,j} \neq \min_{j \in \mathcal{J}_k} c_k \mu_{k,j}$ for any distinct pair $i, k \in \mathcal{N}$, then $V^{\pi^{c\mu}}(\mathbf{X}, \mathbf{b})$ is piecewise-linear on \mathcal{B}' .*

This result is related to two facts: (i) for any given initial prior $\mathbf{b} \in \mathcal{B}'$ (and $\mathbf{X} \in \mathbb{Z}_+^n$), the $Ec\mu$ policy is unique unless \mathbf{b} lies on the break points of the piecewise-linear function $V^{\pi^{c\mu}}(\mathbf{X}, \mathbf{b})$ (see Lemmas 7 and 3 in Online Appendix B), and (ii) policies can be evaluated as linear functions of the belief in any POMDP. Therefore, with respect to closed, nonzero portions of the belief space, the value function $V^{\pi^{c\mu}}(\mathbf{X}, \mathbf{b})$ is differentiable (except at break points). As we show in Theorem 2, the differentiability of the value function strongly enhances the relationship between optimal policies of the nonrobust problem and those under the robust percentile optimization program (5) and (6). Thus, identifying an asymptotically optimal policy that exhibits this property enables us to solve the robust percentile optimization program in an efficient way. This is an important insight to our search for robust chance-constrained

policies, especially because, as Zhang (2010) states, there are no known general conditions over which a POMDP value function is differentiable on its entire belief space.

To the purpose of finding robust chance-constrained policies, we introduce the following set of policies. Fix the initial \mathbf{X} and let $\mathcal{H}_{\mathbf{b}} = \{\pi_{\mathbf{b}}^1, \pi_{\mathbf{b}}^2, \dots, \pi_{\mathbf{b}}^k\}$ be any finite set of optimal policies to the nonrobust problem when the initial prior is \mathbf{b} and $\mathbf{p} = (p_1, p_2, \dots, p_k)$ be an associated distribution such that $\sum_{i=1}^k p_i = 1$. We define a policy $\pi_{\mathcal{H}_{\mathbf{b}}}^{\mathbf{p}}$ to be a *randomized policy* if, at time 0, an element of $\mathcal{H}_{\mathbf{b}}, \pi_{\mathbf{b}}^i$, is chosen with probability p_i , which will dictate all current and future decisions.⁸

Interestingly, similar to other nonlearning robust problems (see, e.g., Bertsimas and Thiele 2006), we find that there exists a randomized policy that forms an optimal solution to the robust percentile problem. Because all policies are evaluated as linear functions over \mathcal{B} , we find that the hyperplane that generates a lower bound to the percentile objective is achievable via a convex combination of optimal nonrobust policies. This means that there exists an optimal robust policy that randomizes between optimal nonrobust policies obtained for a single belief point $\mathbf{b} \in \mathcal{B}$. Furthermore, we shed light on conditions (associated with the differentiability of $V(\mathbf{X}, \mathbf{b})$ with respect to the belief space) such that a *deterministic* nonrobust policy is optimal even for the robust percentile problem.

Theorem 2 (Chance-Constrained Policy). *For any given $\epsilon \geq 0$, there exists a $\mathbf{b}^* \in \mathcal{B}$ and a distribution \mathbf{p}^* forming a randomized policy $\pi_{\mathcal{H}_{\mathbf{b}^*}}^{\mathbf{p}^*}$ that is optimal under the percentile optimization program (5) and (6): $Y^{\pi_{\mathcal{H}_{\mathbf{b}^*}}^{\mathbf{p}^*}}(\mathbf{X}, \epsilon) = Y(\mathbf{X}, \epsilon) = V(\mathbf{X}, \mathbf{b}^*)$. Furthermore, if $V^{\pi_{\mathbf{b}}}(\mathbf{X}, \mathbf{b})$ is differentiable at \mathbf{b}^* , then $\mathcal{H}_{\mathbf{b}^*}$ consists of a single policy, and hence, $\pi_{\mathcal{H}_{\mathbf{b}^*}}^{\mathbf{p}^*}$ is deterministic.*

This result significantly reduces the complexity of the search for optimal robust policies. Importantly, it implies that we can combine policies associated with the function $V(\mathbf{X}, \mathbf{b}^*)$ to find chance-constrained policies. In this way, we no longer need to look at the general space of policies but, rather, can focus on the class of nonrobust optimal policies. Moreover, Proposition 3 shows that the differentiability condition of Theorem 2 can be met by a surface that converges to the value function. If \mathbf{b}^* lies on a linear segment of the value function that is not a break point, $\mathcal{H}_{\mathbf{b}^*}$ can be composed of a single policy yielding a deterministic chance-constrained policy. Hence, under this assumption, one need not be concerned with finding \mathbf{p}^* .

However, Theorem 2 leaves us with an important question: what belief, \mathbf{b}^* , should be used to form the chance-constrained policy $\pi_{\mathcal{H}_{\mathbf{b}^*}}^{\mathbf{p}^*}$ for a given percentile problem? If such a \mathbf{b}^* is characterized, then the solution to the percentile problem can easily be found by a randomization of nonrobust policies associated with \mathbf{b}^* . The answer to this question turns out to be closely related to the geometrical concept of the convex floating body first discussed by Dupin (1822) and later used in robust literature to generate ambiguity sets that guarantee performance for policies evaluated within these sets (see, e.g., Lagoa et al. 2005 and Bertsimas et al. 2018). However, we utilize the convex floating body to characterize \mathbf{b}^* , which generates a policy satisfying the chance-constrained objective.

Definition 1 (Convex Floating Body). Let $\mathcal{W}_{\epsilon} = \{(\mathbf{w}, w) \in \mathbb{R}^m \times \mathbb{R} : \Pr(\mathbf{B}\mathbf{w}^T \geq w) \leq \epsilon\}$ be the set of all half spaces that “cut off” ϵ or less volume of the belief space \mathcal{B} with respect to $\mathbb{P}_{\mathbf{B}}$. An ϵ -based convex floating body on \mathcal{B} is $\mathcal{L}_{\epsilon} = \bigcap_{(\mathbf{w}, w) \in \mathcal{W}_{\epsilon}} \{\mathbf{b} \in \mathcal{B} : \mathbf{b}\mathbf{w}^T \leq w\}$. We let $\delta\mathcal{L}_{\epsilon}$ be the boundary of \mathcal{L}_{ϵ} .⁹

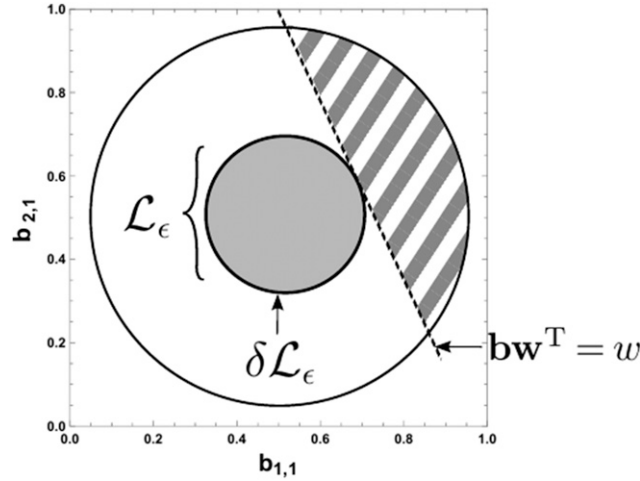
Based on this definition, a convex floating body is the region left from hyperplanes cutting off a specified volume (ϵ) from an object. For every $\mathbf{b} \in \delta\mathcal{L}_{\epsilon}$, there exists a hyperplane that divides \mathcal{B} into two pieces: one that has volume less than or equal to ϵ . Figure 2 illustrates the convex floating body of a sphere with uniform density, which is either the empty set or another sphere. We study convex floating bodies with respect to density $\mathbb{P}_{\mathbf{B}}$ on the belief space of our priors to characterize \mathbf{b}^* and thereby find optimal chance-constrained policies as discussed in Theorem 2.

For the purposes of characterizing \mathbf{b}^* , it is important that \mathcal{L}_{ϵ} is nonempty. Fortunately, Fresen (2013) states that when $\mathbb{P}_{\mathbf{B}}$ yields a log-concave probability distribution, \mathcal{L}_{ϵ} exists so long as $\epsilon \leq e^{-1}$. Hence, for many robust applications that tend toward pessimism (where ϵ is small), under common distributions, the convex floating body is guaranteed to exist.¹⁰ If \mathcal{L}_{ϵ} is nonempty, we find that \mathbf{b}^* (defined in Theorem 2) is found at the largest value of the nonrobust problem on the boundary of the convex floating body.

Proposition 4 (Characterizing $\mathcal{H}_{\mathbf{b}^*}$). *If \mathcal{L}_{ϵ} is nonempty, then $\mathbf{b}^* = \arg\max_{\mathbf{b} \in \delta\mathcal{L}_{\epsilon}} V(\mathbf{X}, \mathbf{b})$, where \mathbf{b}^* satisfies $Y(\mathbf{X}, \epsilon) = V(\mathbf{X}, \mathbf{b}^*)$.*

Interestingly, Proposition 4 relates percentile optimization to a minimax objective: one can search for a worst-case belief within a specified set. Because $V(\mathbf{X}, \mathbf{b})$ is concave in \mathbf{b} (by the convexity results of Sondik 1971 and Smallwood and Sondik 1973), if $\delta\mathcal{L}_{\epsilon}$ is easily characterized, we can apply gradient-based optimization to solve the problem rather than evaluating the entire surface, which is computationally intractable. Although Theorem 2 states that $\mathcal{H}_{\mathbf{b}^*}$ is a singleton when the value function is differentiable at \mathbf{b}^* , the differentiability is not always guaranteed. To this end, in the proof of Proposition 4 (see Online Appendix B), we characterize \mathbf{p}^* . We find that the distribution \mathbf{p}^* such that the contour $\{\mathbf{b} \in \mathcal{B} | V^{\pi_{\mathcal{H}_{\mathbf{b}^*}}^{\mathbf{p}^*}}(\mathbf{X}, \mathbf{b}) = V(\mathbf{X}, \mathbf{b}^*)\}$ is a subgradient hyperplane to \mathcal{L}_{ϵ} .

Figure 2. A Convex Floating Body \mathcal{L}_ϵ When $\mathbb{P}_\mathbf{B}$ Is Uniform Within the Circle and Is Zero Elsewhere



Notes. It is generated from the intersection of half spaces $(\mathbf{w}, w) \in \mathcal{W}_\epsilon$, and the striped area must contain less than or equal to ϵ volume. ($n = 2, m_1, m_2 = 2$).

In general, because nonrobust policies are only partially characterized (they converge to $Ec\mu$ policies asymptotically), it is important to connect the $Ec\mu$ policies to the percentile optimization objective. The following corollary is similar to Proposition 4 and shows that there exists a finite randomization of $Ec\mu$ policies that are asymptotically optimal to the percentile objective as $\psi \rightarrow \infty$.

Corollary 1 (Robust $Ec\mu$ Optimality). *If \mathcal{L}_ϵ is nonempty, then there exists a policy π that is a finite randomization of $Ec\mu$ policies such that $Y^\pi(\mathbf{X}, \epsilon) - Y(\mathbf{X}, \epsilon) \leq V^{\pi^{cu}}(\mathbf{X}, \hat{\mathbf{b}}) - V(\mathbf{X}, \mathbf{b}^*)$, where $\hat{\mathbf{b}} = \arg \max_{\mathbf{b} \in \delta \mathcal{L}_\epsilon} V^{\pi^{cu}}(\mathbf{X}, \mathbf{b})$ and \mathbf{b}^* is defined in Theorem 2.*

This corollary holds despite the fact that $V^{\pi^{cu}}(\mathbf{X}, \mathbf{b})$ is not guaranteed to be concave in \mathbf{b} . In fact, if it is concave in \mathbf{b} , the randomized policy π can be directly built from nonrobust policies. However, if $V^{\pi^{cu}}(\mathbf{X}, \mathbf{b})$ is not concave in \mathbf{b} , we can still form the appropriate randomized policy satisfying Corollary 1 via a randomization of policies that satisfy minimax solutions within the set of $Ec\mu$ policies on the boundary of the convex floating body, namely, $\min_{\mathbf{b}^1 \in \mathcal{B}} \max_{\mathbf{b}^2 \in \delta \mathcal{L}_\epsilon} V^{\pi^{cu}}(\mathbf{X}, \mathbf{b}^2)$.

With respect to optimal solutions to the percentile objective, additional results can further confine $\mathcal{H}_{\mathbf{b}^*}$ (of Theorem 2) by noting that \mathbf{b}^* must lie near the extreme belief state with worst-case transition parameters. We denote this worst-case belief state by \mathbf{b}_0 and note that it is composed of components

$$b_{i,j}^0 = \begin{cases} 1 & \text{if } \mu_{i,j} = \min_{k \in \mathcal{J}_i} \mu_{i,k}, \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

It can be shown (see the proof of Proposition 5) that, for any policy, \mathbf{b}_0 is the worst-case (most expensive) belief state for the system. To further characterize \mathbf{b}^* , we define the concept of *visibility* (adopted from geometry literature but repurposed for our needs).

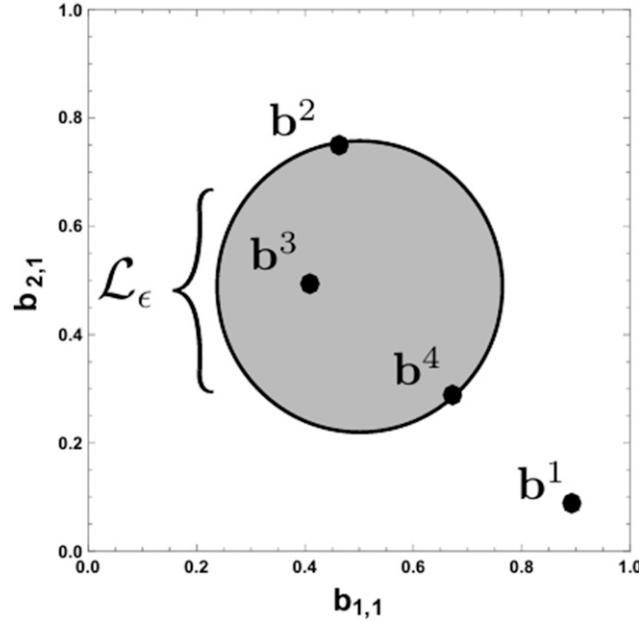
Definition 2 (Visibility). A belief point $\mathbf{b} \in \mathcal{L}_\epsilon$ is said to be visible from a reference belief $\mathbf{b}^1 \in \mathcal{B}$ if $\{\mathbf{b}^2 \in \mathcal{B} : \mathbf{b}^2 = \eta \mathbf{b} + (1 - \eta) \mathbf{b}^1, \eta \in [0, 1]\} \cap \mathcal{L}_\epsilon = \mathbf{b}$.

As demonstrated in Figure 3, a belief \mathbf{b} in the convex floating body is visible from a reference belief \mathbf{b}^1 if, on the line segment connecting these points, only \mathbf{b} lies within the convex floating body. This implies that, if the reference belief point \mathbf{b}^1 is distinct from \mathbf{b} and \mathbf{b} is visible from \mathbf{b}^1 , then \mathbf{b} must lie in the boundary ($\mathbf{b} \in \delta \mathcal{L}_\epsilon$). However, not every point on $\delta \mathcal{L}_\epsilon$ is visible from a reference point \mathbf{b}^1 . In the following proposition, we show that the belief \mathbf{b}^* (introduced in Theorem 2) must be visible from the worst-case belief state \mathbf{b}_0 .

Proposition 5 (Visibility of \mathbf{b}^*). *If \mathcal{L}_ϵ is nonempty, then there exists a \mathbf{b}^* visible from the worst-case belief \mathbf{b}_0 .*

Proposition 5 significantly helps us find \mathbf{b}^* (of Theorem 2): we only need to search part of $\delta \mathcal{L}_\epsilon$, which is visible from \mathbf{b}_0 . Proposition 5 also can facilitate establishing effective heuristics that circumvent the calculation of the nonrobust problem. For instance, Figure 4 demonstrates the implications of Proposition 5 for a uniform type $\mathbb{P}_\mathbf{B}$: \mathbf{b}^* lies somewhere on the dashed line.

Figure 3. Belief Points \mathbf{b}^2 and \mathbf{b}^3 Are Not Visible from Reference Belief \mathbf{b}^1 , Whereas \mathbf{b}^4 Is Visible from Reference Belief \mathbf{b}^1 ($n = 2, m_1, m_2 = 2$)

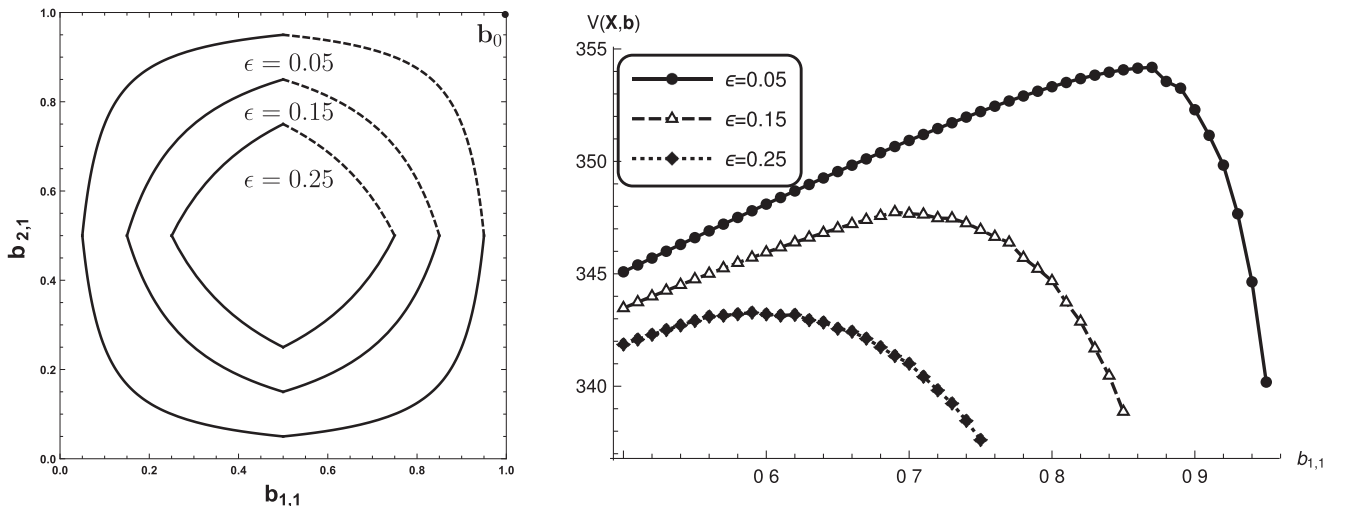


5. Asymptotically Tight Bounds

Although we have characterized the optimal policies of the nonrobust and percentile problems, evaluating the nonrobust value function $V(\mathbf{X}, \mathbf{b})$ is still a computationally complex problem (see, e.g., Littman et al. 1998, Mundhenk et al. 2000, and Papadimitriou and Tsitsiklis 1987 for an in-depth discussion regarding the complexity of POMDP programs). If the value function $V(\mathbf{X}, \mathbf{b})$ and the convex floating body's boundary $\delta\mathcal{L}_\epsilon$ are known, the solution to the percentile optimization is easily characterizable (Theorem 2 and Propositions 4 and 5). Therefore, we provide computationally tractable bounds to the nonrobust problem that can be evaluated in closed form to facilitate the computability of chance-constrained policies.

The bounds we form are based on the performance of (a) queues under no model ambiguity with fixed rate parameters equal to $E[\mu_i|\mathbf{b}]$ and (b) following a particular server allocation priority rule based on the initial parameter belief. These imply that our bounds rely only on the valuation of fixed priority-based policies

Figure 4. On the Left, Convex Floating Bodies \mathcal{L}_ϵ for $\epsilon = 0.05, 0.15, 0.25$ with $n = 2, m_1, m_2 = 2$ and Uniform \mathbb{P}_B



Notes. To be visible from \mathbf{b}_0 , belief \mathbf{b}^* associated with Proposition 4 must lie on the dashed lines assuming $\mu_{1,1} < \mu_{1,2}$ and $\mu_{2,1} < \mu_{2,2}$ (Proposition 5). On the right, $V((10, 10), \mathbf{b})$ is evaluated on these boundaries when $\mu_{1,1} = 0.1, \mu_{1,2} = 0.2, \mu_{2,1} = 0.05, \mu_{2,2} = 0.25$. Belief \mathbf{b}^* lies at the peak of these curves.

that do not change with dynamic observations, significantly reducing the computational complexity of the problem.

For a given belief $\hat{\mathbf{b}} \in \mathcal{B}$, consider a counterpart system identical to our original setting with the exception of the ambiguity sets being $\hat{\mathcal{M}}_i = \{E[\mu_i|\hat{\mathbf{b}}]\}$ (analogous to the original ambiguity sets \mathcal{M}_i). That is, the counterpart queueing system has fully known service rates that are calculated based on taking an expectation of service rates in \mathcal{M}_i over belief $\hat{\mathbf{b}}$. Obviously, the optimal policy for this system is the traditional $c\mu$ rule because all of its parameters are fully known. Let $\pi_{\hat{\mathbf{b}}}$ denote this $c\mu$ rule and $\bar{V}^{\pi_{\hat{\mathbf{b}}}}(\mathbf{X}, \hat{\mathbf{b}})$ be the associated infinite-horizon cost of the counterpart system under $\pi_{\hat{\mathbf{b}}}$. It is important to emphasize that $\pi_{\hat{\mathbf{b}}}$ exhaustively serves class $\arg \max_{a \in \mathcal{A}(\mathbf{X})} c_a E[\mu_a|\hat{\mathbf{b}}]$ until no customer of that class remains in the system and acts only as a function of the queue state, not of belief, even when $\pi_{\hat{\mathbf{b}}}$ is implemented in the original system. When $\pi_{\hat{\mathbf{b}}}$ is implemented in the original system, we denote the infinite-horizon cost by $V^{\pi_{\hat{\mathbf{b}}}}(\mathbf{X}, \hat{\mathbf{b}})$. Using the counterpart system's cost and its associated policy, we can bound the nonrobust cost (which is needed to calculate the robust cost; see Theorem 2 and Proposition 4) using the following proposition.

Proposition 6 (Asymptotically Tight Bounds). *For any state $(\mathbf{X}, \hat{\mathbf{b}})$, the nonrobust cost $V(\mathbf{X}, \hat{\mathbf{b}})$ is bounded as $\bar{V}^{\pi_{\hat{\mathbf{b}}}}(\mathbf{X}, \hat{\mathbf{b}}) \leq V(\mathbf{X}, \hat{\mathbf{b}}) \leq V^{\pi_{\hat{\mathbf{b}}}}(\mathbf{X}, \hat{\mathbf{b}})$. Furthermore,*

- i. *The gap between the upper and lower bound costs decrease to zero as queue length X_i increases to infinity, where $i = \arg \max_{a \in \mathcal{A}(\mathbf{X})} c_a E[\mu_a|\hat{\mathbf{b}}]$.*
- ii. *The gap between the upper and lower bound costs monotonically decrease to zero as $\text{Var}[\mu_i|\hat{\mathbf{b}}]$ decrease to zero (for all $i \in \mathcal{N}$).*

Both the upper and lower bounds of Proposition 6 are easily calculable (see Online Appendix B). Furthermore, under these conditions, these bounds become arbitrarily close approximations, which adds computational tractability to the problem as well as analytical insight to the relationship between our nonrobust and traditional $c\mu$ policies. In particular, part (ii) of Proposition 6 supports the intuition that gathering more data on unknown service parameters can provide more accurate bound information. Part (i) of Proposition 6 provides conditions under which the myopic, nonlearning policy's cost converges to that of the optimal policy.

Remark 1. Because the percentile objective relies on the computation of the nonrobust problem, the bound results can be easily applied to the percentile formulation as well. For instance, one can refine the search for $\arg \max_{\mathbf{b} \in \delta\mathcal{L}_\epsilon} V(\mathbf{X}, \mathbf{b})$ as in Proposition 4: if the upper bound for a $\mathbf{b} \in \delta\mathcal{L}_\epsilon$ is less than the lower bound for $\mathbf{b}' \in \delta\mathcal{L}_\epsilon$, \mathbf{b} must not be the belief point \mathbf{b}^* . Because most infinite-horizon POMDPs are calculated by finite-horizon approximations, a second application of the bounds is to use them as the terminal cost used in the finite-horizon dynamic program. That is, when evaluating the finite-horizon approximation, one can replace $V_0(\mathbf{X}, \mathbf{b})$ by lower and upper bounds $\bar{V}^{\pi_{\hat{\mathbf{b}}}}(\mathbf{X}, \hat{\mathbf{b}})$ and $V^{\pi_{\hat{\mathbf{b}}}}(\mathbf{X}, \hat{\mathbf{b}})$, respectively. This can provide very tight bounds on the POMDP because after a certain number of “learning periods,” where the POMDP is explicitly evaluated, the controller might have collected enough information to have enough confidence in the true transition parameters.

6. An Analytically Rooted Heuristic Policy

Chance-constrained policies are inherently difficult to calculate even given the analytical results established in the previous section. To circumvent complexity arising from (a) the PSPACE-hard problem of evaluating a POMDP over a belief space with high dimensionality and (b) finding the shape of the convex floating body, which requires high-dimensional polytope approximations, we now introduce an effective heuristic policy. This heuristic policy operates by simply choosing the $Ec\mu$ policy associated with the belief point on the convex floating body's boundary $\delta\mathcal{L}_\epsilon$ that minimizes the distance from \mathbf{b}_0 (the worst-case parameter settings for each class characterized in (7)). This is typically an easy-to-perform task, especially in the cases of uniform and spherical-type distributions on the belief space, allowing for managers to benefit from our approach without requiring demanding computations. Moreover, as we show in Section 7, this heuristic performs extremely well on both randomly generated and real-world data that we have collected from a leading U.S. hospital.

We term the $Ec\mu$ policy with expectation taken based on belief point $\arg \min_{\mathbf{b} \in \delta\mathcal{L}_\epsilon} \|\mathbf{b}_0 - \mathbf{b}\|$, where $\|\cdot\|$ is the l^2 norm, as the $(1 - \epsilon)Ec\mu$ heuristic policy. This heuristic policy takes advantage of three main structural results of the chance-constrained policy (that we established in the previous section) while providing a much simpler version of it:

1. It assumes that the true optimal policies of the nonrobust problem are $Ec\mu$, a fact supported by Theorem 1, which shows the asymptotic relationship of the optimal policies to $Ec\mu$.
2. It locates belief $\arg \min_{\mathbf{b} \in \delta\mathcal{L}_\epsilon} \|\mathbf{b}_0 - \mathbf{b}\|$ to be near \mathbf{b}^* (of Theorem 2) based on Proposition 4. The worst-case (most expensive) belief state is \mathbf{b}_0 , and through the proof of Proposition 5 (see Online Appendix B), the value function is

nonincreasing in λ with respect to belief $\lambda \mathbf{b} + (1 - \lambda) \mathbf{b}_0$ for $\lambda \in [0, 1]$. Thus, $\arg \max_{\mathbf{b} \in \delta \mathcal{L}_e} V(\mathbf{X}, \mathbf{b})$ is expected to be near \mathbf{b}_0 .¹¹

3. It takes advantage of the fact that $\arg \min_{\mathbf{b} \in \delta \mathcal{L}_e} \|\mathbf{b}_0 - \mathbf{b}\|$ satisfies Proposition 5 (because this belief is visible from \mathbf{b}_0).

7. Numerical Experiments

We now perform various numerical experiments to (a) identify the advantages of chance-constrained policies in a variety of environments under model ambiguity, (b) demonstrate the sensitivities of the underlying queueing models, (c) study the effectiveness of the proposed $Ec\mu$ heuristic in mimicking the optimal chance-constrained policies, and (d) demonstrate the implications of our results in real-world applications. To pursue these goals, we present our analyses in five parts: we (a) establish the sensitivities in initial prior selection; (b) investigate how our policies perform over a large parameter suite but in a relatively small queueing system; (c) evaluate our proposed heuristic alongside percentile, minimax, and minimin policies in a larger system; (d) demonstrate the gap between the $Ec\mu$ and optimal (nonrobust) policies; and (e) apply the $Ec\mu$ heuristic to a hospital ED setting using real-world data and discuss its significant implications on improving the current patient-flow policies.

To help establish the necessity of our robust percentile formulation, it is first important to establish the sensitivities of the nonrobust value function under small perturbations in belief. To this end, we evaluate the expected cost under a variety of parameter settings when $n = 2, m_1 = 2$, and $m_2 = 2$ with respect to a “central prior” $\bar{\mathbf{b}} = (0.5, 0.5, 0.5, 0.5)$, which assumes a uniform distribution on parameters, a slightly pessimistic $\bar{\mathbf{b}}^p = (0.6, 0.4, 0.6, 0.4)$, and a slightly optimistic prior $\bar{\mathbf{b}}^o = (0.4, 0.6, 0.4, 0.6)$. Table 1 displays the results from comparing the percentage difference between nonrobust value functions evaluated at these priors (for various parameter configurations) via the expression

$$\frac{|V(\mathbf{X}, \mathbf{b}) - V(\mathbf{X}, \hat{\mathbf{b}})|}{(V(\mathbf{X}, \mathbf{b}) + V(\mathbf{X}, \hat{\mathbf{b}}))/2} \%$$

for two distinct priors $\mathbf{b}, \hat{\mathbf{b}} \in \mathcal{B}$.

Even with relatively small perturbations to the selection of the prior, as can be seen from Table 1, differences in value function are substantial. Thus, we make the following:

Observation 1 (Sensitivity to Prior Specification). The expected cost of the nonrobust problem is sensitive to the choice of prior.

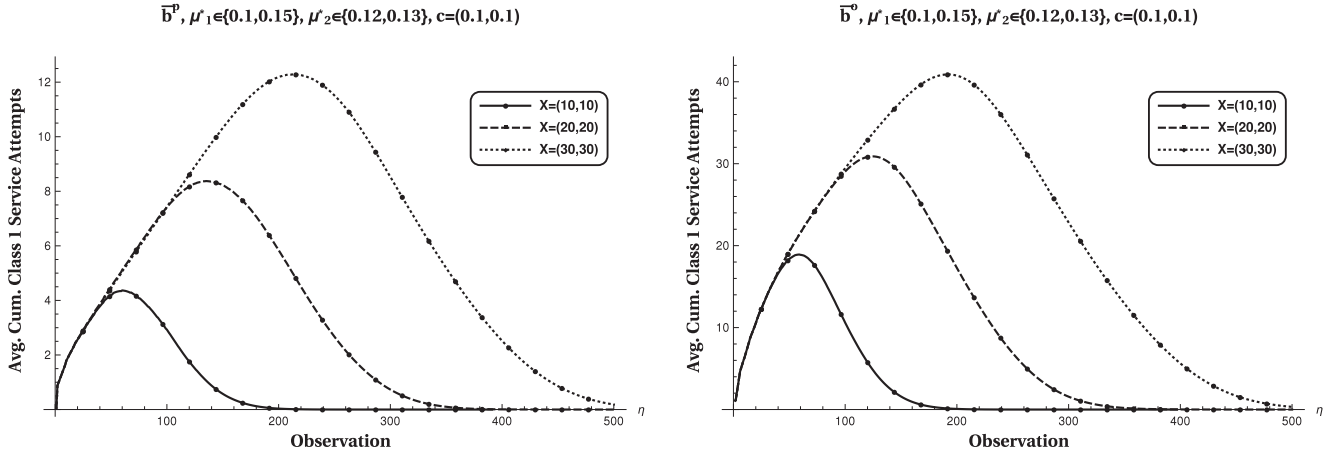
Can slight perturbations in the prior also cause significant differences in policies obtained from the nonrobust framework? To answer this, we again turn to our $n = 2, m_1 = 2, m_2 = 2$ environment and investigate the differences nonrobust policies $Ec\mu$ experience as their prior changes from $\bar{\mathbf{b}}, \bar{\mathbf{b}}^p$, and $\bar{\mathbf{b}}^o$. We run simulations in which the true parameter settings are selected according to $\bar{\mathbf{b}}$. To identify differences between the policies at different initial priors, we track the cumulative number of attempts to serve class 1 by time t under each policy and depict the results in Figure 5.

Figure 5 shows that policies experience an extended period of time in which they disagree on the class to serve. This is especially evident when a large number of customers are in the system, indicating that policies only begin to converge after having finished the service of the class. Furthermore, as discussed earlier, the $Ec\mu$ policy experiences

Table 1. Percentage Gaps for $\bar{\mathbf{b}} = (0.5, 0.5, 0.5, 0.5)$, $\bar{\mathbf{b}}^p = (0.6, 0.4, 0.6, 0.4)$, and $\bar{\mathbf{b}}^o = (0.4, 0.6, 0.4, 0.6)$, Where $n = 2, m_1 = 2, m_2 = 2$

X	Percentage differences						X	Percentage differences					
	$(\mu_{1,1}, \mu_{1,2})$	$(\mu_{2,1}, \mu_{2,2})$	(c_1, c_2)	$\bar{\mathbf{b}} \text{ vs } \bar{\mathbf{b}}^p$	$\bar{\mathbf{b}} \text{ vs } \bar{\mathbf{b}}^o$	$\bar{\mathbf{b}}^p \text{ vs } \bar{\mathbf{b}}^o$		$(\mu_{1,1}, \mu_{1,2})$	$(\mu_{2,1}, \mu_{2,2})$	(c_1, c_2)	$\bar{\mathbf{b}} \text{ vs } \bar{\mathbf{b}}^p$	$\bar{\mathbf{b}} \text{ vs } \bar{\mathbf{b}}^o$	$\bar{\mathbf{b}}^p \text{ vs } \bar{\mathbf{b}}^o$
(5, 5)	(0.1, 0.2)	(0.15, 0.3)	(0.1, 0.1)	5.24	5.52	10.76	(5, 5)	(0.05, 0.15)	(0.1, 0.2)	(0.15, 0.2)	5.31	5.8	11.1
(10, 10)				4.25	4.98	9.22	(10, 10)				4.15	4.47	8.62
(15, 15)				3.62	4.29	7.9	(15, 15)				3.39	3.37	6.75
(5, 5)	(0.05, 0.15)	(0.1, 0.2)	(0.1, 0.1)	5.21	6.1	11.3	(5, 5)	(0.1, 0.2)	(0.15, 0.3)	(0.2, 0.15)	5.33	5.8	11.12
(10, 10)				4.11	4.63	8.74	(10, 10)				4.7	4.77	9.46
(15, 15)				3.19	3.7	6.89	(15, 15)				4.02	3.89	7.91
(5, 5)	(0.05, 0.1)	(0.06, 0.08)	(0.1, 0.1)	3.03	3.23	6.26	(5, 5)	(0.05, 0.15)	(0.1, 0.2)	(0.2, 0.15)	6.32	6.0	12.32
(10, 10)				2.24	2.23	4.47	(10, 10)				4.75	4.6	9.35
(15, 15)				1.58	1.63	3.21	(15, 15)				3.64	3.81	7.45
(5, 5)	(0.1, 0.2)	(0.1, 0.2)	(0.15, 0.3)	5.46	5.27	10.72	(5, 5)	(0.05, 0.1)	(0.06, 0.08)	(0.2, 0.15)	3.44	3.95	7.39
(10, 10)				4.39	4.89	9.27	(10, 10)				2.66	2.63	5.3
(15, 15)				3.97	4.05	8.02	(15, 15)				2.04	1.96	4.0
Average percentage											3.72	3.93	7.64

Figure 5. Comparison of Two Nonrobust Policies Under Slight Perturbations of the Initial Prior \bar{b}
 ($\mu_{1,1} = 0.1, \mu_{1,2} = 0.15, \mu_{2,1} = 0.12, \mu_{2,2} = 0.13$)

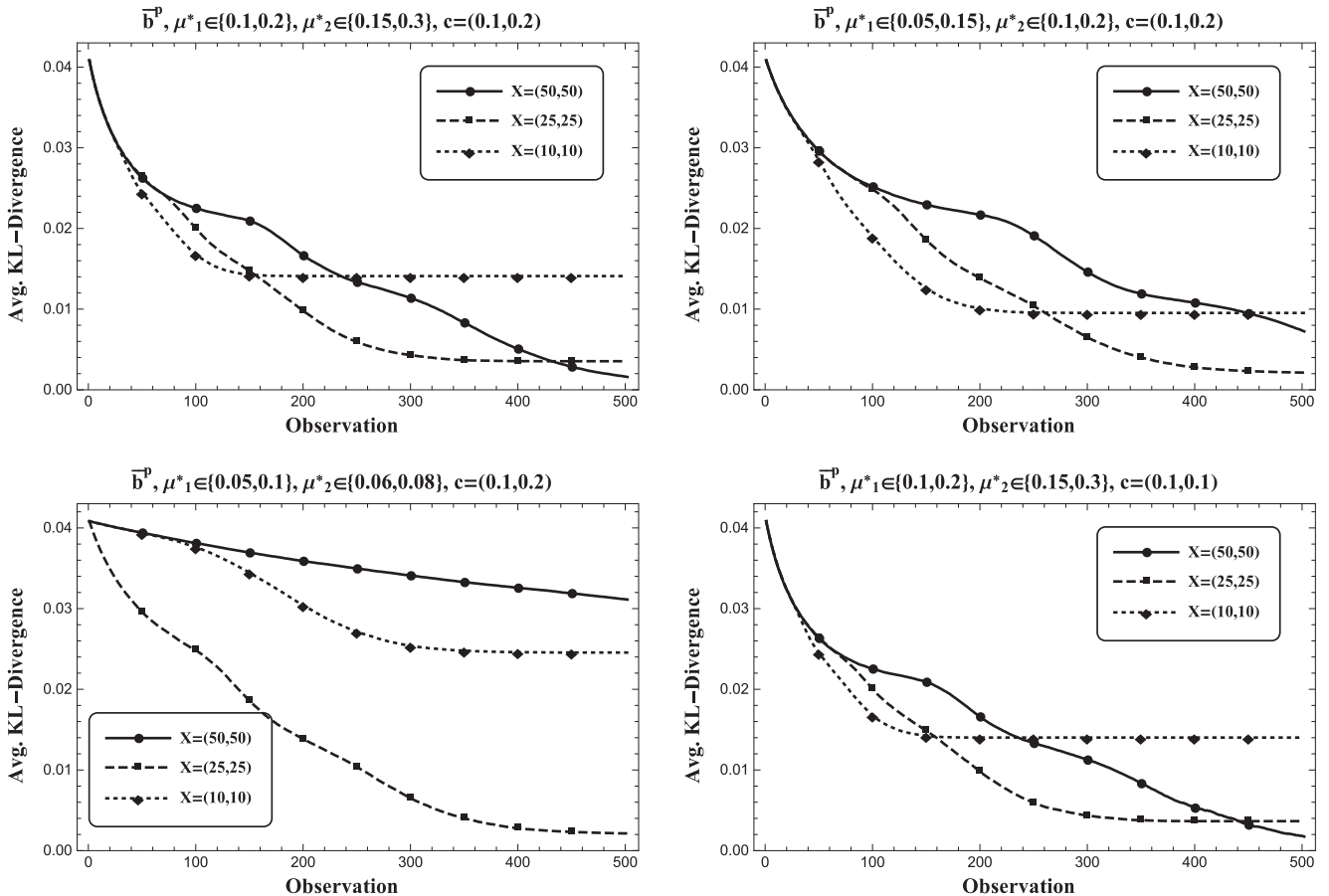


momentum toward serving a customer after a successful service. Therefore, as can be seen from Figure 5, two policies with slightly different starting beliefs (i.e., initial priors) may experience very different action profiles.

Thus, not only does the value function experience sensitivity among different selections of the prior, but these differences also correspond to policy changes. Thus, we make the following:

Observation 2 (Sensitivity in Policy). The policies generated from the nonrobust problem are sensitive to the choice of prior.

Figure 6. Comparison of the Average KL Divergence Between Two Policies' Beliefs When the True Prior Is \bar{b}



If the duration of time in which the optimal policy experiences learning is relatively small, the choice of the initial prior becomes inconsequential because the difference between initial priors will be quickly “washed out” by the incoming data. To test whether the differences between initial priors are long-lasting, in Figure 6 we compare the KL divergence¹² between two beliefs after each observation under their associated policies.

From Figure 6, it is evident that the beliefs converge to one another given that there are enough customers to serve. In general, the learning is faster for smaller queue states until all of the customers have been served because the policies are, in effect, closer to one another. However, even in the smaller queue settings, the learning rate is not fast enough to disregard the choice of the initial prior. Thus, we make the following:

Observation 3 (Slow Convergence in Belief). The differences between the beliefs about the correct model with differing initial priors are long-lasting.

To better understand the relative performance of our robust percentile policies, we start by considering a large parameter suite including more than 1,000 parameter settings in an $n = 2, m_1 = 2$, and $m_2 = 2$ setting with four different \mathbb{P}_B at their 95% chance-constrained policy. We name these \mathbb{P}_B densities f_1, f_2, f_3 , and f_4 , respectively: f_1, f_2 , and f_3 are truncated multivariate normal with means $\mu_1 = (0.5, 0.5)$, $\mu_2 = (0.4, 0.4)$, $\mu_3 = (0.6, 0.6)$ and covariance matrices $\Sigma_1 = \begin{pmatrix} 1.5 & 0.0 \\ 0.0 & 1.5 \end{pmatrix}$, $\Sigma_2 = \begin{pmatrix} 0.5 & 0.0 \\ 0.0 & 0.5 \end{pmatrix}$, and $\Sigma_3 = \begin{pmatrix} 0.5 & 0.0 \\ 0.0 & 0.5 \end{pmatrix}$, respectively. Finally, f_4 is uniform.

We include two nonlearning robust policies (minimin and minimax) as benchmarks for the performance of our robust percentile policies and compare the policies by evaluating their total cost when each model (i.e., parameter configuration) is equally likely. That is, we assume that the true (but unknown) prior of our system is $\bar{\mathbf{b}} = (0.5, 0.5, 0.5, 0.5)$, and we evaluate the total cost under 95% chance-constrained, minimax, and minimin policies. Furthermore, we assume $c_1 = c_2$. In every problem instance, we assume $\mu_{2,1} < \mu_{1,1}$ and $\mu_{1,2} < \mu_{2,2}$ so that the policy is not uniform throughout the belief space, which provides incentive for gaining additional knowledge. Further detail on this parameter suite is presented in Online Appendix A.1.

We next compare our proposed policies with other nonlearning robust policies (minimax and minimin). In Table 2, we present the results of this comparison expressed by the average (among all models) optimality gap percentage under various policies. The optimality gap percentage for policy π at \mathbf{b} is defined as

$$\frac{V^\pi(\mathbf{X}, \mathbf{b}) - V(\mathbf{X}, \mathbf{b})}{V(\mathbf{X}, \mathbf{b})} \%.$$

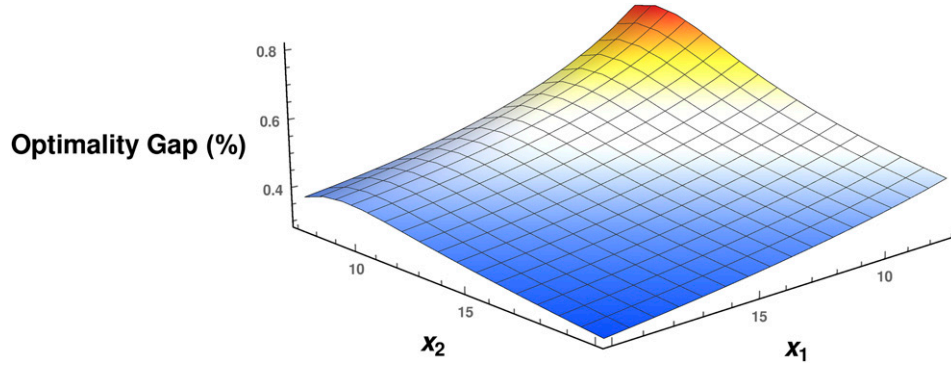
From Table 2, we observe that, on average, our proposed chance-constrained policies perform much better than the other nonlearning policies. Because there is equal chance of every parameter configuration, nonlearning policies serve the wrong class for a realized set of parameters 50% of the time, which results in poor performance.

Comparing the chance-constrained policies under f_1, f_2, f_3 , and f_4 in Table 2 reveals yet another interesting insight: they exhibit similar performance. The reason behind this is threefold: (a) as a property of Proposition 5, because we used 95% chance-constrained policies, each \mathbf{b}^* tends to be near \mathbf{b}_0 ; (b) even though the distributions f_1, f_2, f_3 , and f_4 are different (e.g., they have differing covariance structures and are centered at different beliefs), their convex floating bodies are quite similar; and (c) the chance-constrained policies we propose exhibit learning. Hence, we can make the following:

Table 2. Performance of Various Robust Policies over the Test Suite ($n = 2, m_1 = 2, m_2 = 2$)

X	Optimality gap (%)					
	Minimax	Minimin	95% chance constrained f_1	95% chance constrained f_2	95% chance constrained f_3	95% chance constrained f_4
(2, 2)	3.17	15.51	1.84	2.07	2.2	1.97
(2, 5)	2.52	13.58	0.85	0.81	0.86	0.86
(2, 10)	1.35	8.21	0.52	0.51	0.54	0.49
(5, 2)	4.48	8.73	2.65	2.74	2.33	2.21
(5, 5)	5.01	10.3	0.85	0.81	0.75	0.79
(5, 10)	3.37	7.56	0.61	0.58	0.48	0.57
(10, 2)	4.14	4.05	1.76	1.93	1.32	1.35
(10, 5)	5.49	5.79	0.53	0.51	0.54	0.55
(10, 10)	4.34	5.15	0.33	0.33	0.35	0.35
Average	3.76	8.76	1.10	1.14	1.04	1.02

Figure 7. (Color online) The Optimality Gap (%) of $Ec\mu$ Policy When Evaluated on the Central Prior $\bar{\mathbf{b}}$ ($\mu_{1,1} = 0.6, \mu_{1,2} = 0.7, \mu_{2,1} = 0.5, \mu_{2,2} = 0.8$)



Observation 4 (Sensitivity). The performance of chance-constrained policies is not sensitive to the choice of $\mathbb{P}_{\mathbf{B}}$.

In Section 6, we introduced the $Ec\mu$ heuristic as an easy-to-implement policy that mimics the performance of robust, optimal, chance-constrained policies. To demonstrate the validity of the first assumption underlying this heuristic—that the optimal policies of the nonrobust problem are $Ec\mu$ —in Figure 7, we depict the percent optimality gap of the $Ec\mu$ heuristic policy by comparing its cost to that of the optimal nonrobust policies in a situation in which ψ is small. Because we know that $Ec\mu$ becomes optimal as ψ becomes large (Theorem 1), this poses a worst-case scenario for the performance of the $Ec\mu$ policies. From Figure 7, we can make the following:

Observation 5 (Near Optimality of $Ec\mu$). Even when ψ is small, the $Ec\mu$ performance is close to the nonrobust optimal policy, especially when the system is highly congested.

Observation 5 confirms that the myopic $Ec\mu$ policy provides us with a good approximation of the optimal POMDP value function (as we would expect given its asymptotic relationship to the chance-constrained policy; see Theorem 1). However, using such a rule to find the explicit surface of the POMDP value function is computationally challenging even though the $Ec\mu$ policy is simple. This is because policy evaluation (even when a policy is known) in POMDPs is PSPACE-complete (see, e.g., Mundhenk et al. 2000). Hence, the ideal task of searching for the max of the convex floating body as in Proposition 4, even with the help of Proposition 5, is highly difficult even in moderate problem instances in which $n > 3$ and $m > 6$. Furthermore, oftentimes the shape of \mathcal{L}_ϵ is difficult to determine explicitly as is the case even in the simple uniform distributions in more than two dimensions, which further complicates our search. Hence, for implementation in real applications, we turn to our robust heuristic policy.

To gain deeper insights into the performance of our heuristic, we simulate systems with $m_1 = m_2 = m_3 = 3$ with uniform $\mathbb{P}_{\mathbf{B}}$ in the largest inscribed sphere of the belief space. To also evaluate the robustness of our proposed heuristic vis-à-vis the optimal percentile policy as well as minimin and minimax policies, we use conditional value at risk, $CVar(q)$, which is the average cost within the most costly $q\%$ of our simulated runs. Therefore, if $\mathcal{S} = \{s_1, \dots, s_r\}$ is the set of the costs from a simulation of r runs ordered from most costly to least costly, then

$$CVar(q) = \frac{\sum_{i=1}^{\lceil (1-q)(r-1)+1 \rceil} s_i}{\lceil (1-q)(r-1)+1 \rceil}.$$

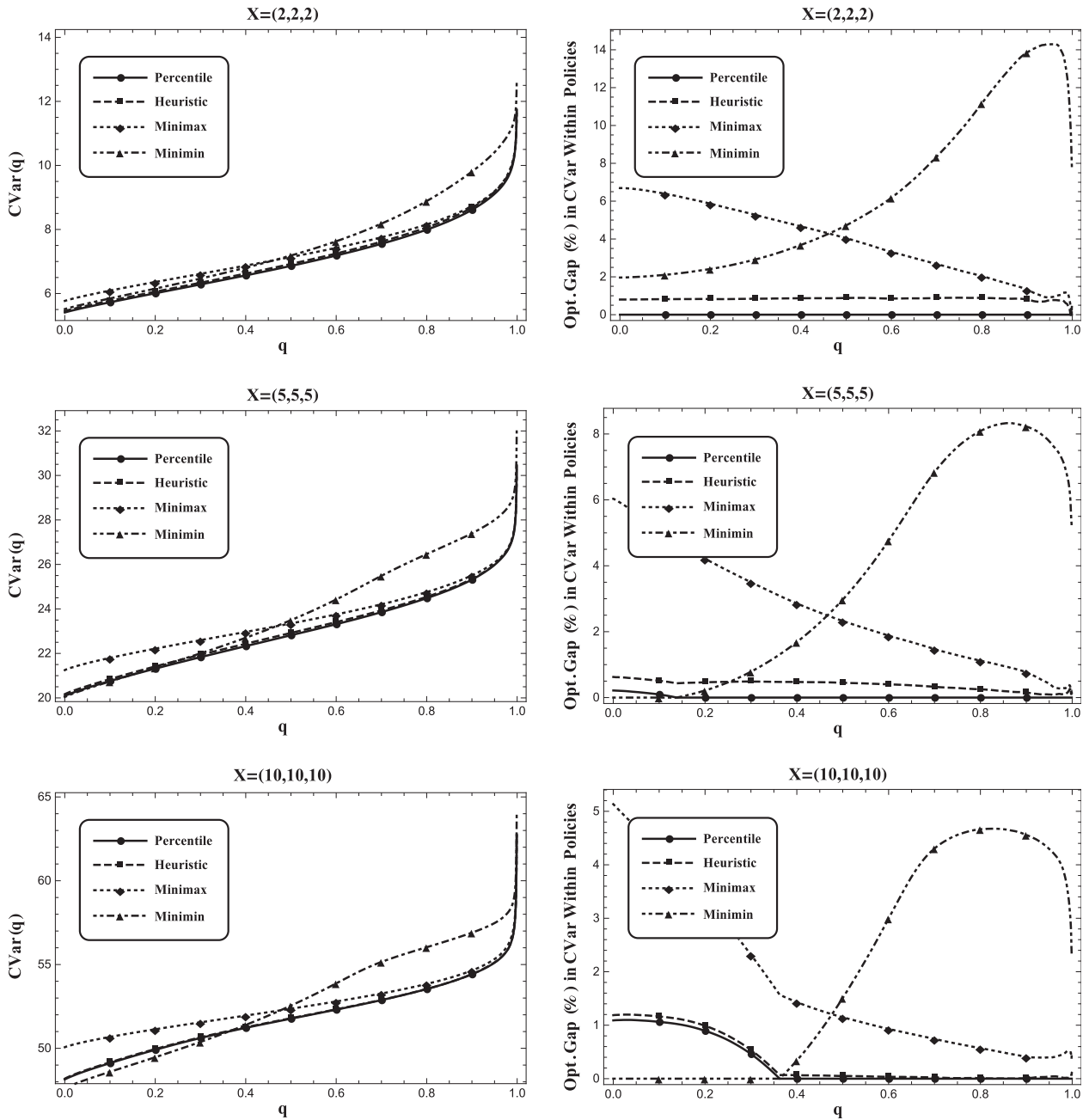
This statistic may roughly be seen as a function that increases in pessimism because we use fewer low-cost data points in the expectation as q increases.¹³

Using a 95% chance-constrained policy, the $Ec\mu$ heuristic, minimin, and minimax policies, Figure 8 illustrates performance over 20,000 simulation runs.¹⁴ The leftmost subfigures display the raw $CVar$ values. However, we direct our attention to the rightmost figures, which display the percentage gap (of $CVar$ s) between the four selected policies and “best” policy at a given q . From Figure 8, we observe the following:

Observation 6 (Heuristic Performance). The $Ec\mu$ heuristic performs nearly identically to the chance-constrained policy with a diminishing difference as the system becomes more congested.

We note that percentile optimization is not concerned about the worst-case scenarios and, rather, optimizes based on a proportion of the belief space. Hence, being a statistic concerned with the tail performance of the distribution of costs, $CVar$ (as compared with the expected cost) provides us with a more accurate representation of the value of robustness that percentile optimization offers. Furthermore, Figure 8 demonstrates that the proposed

Figure 8. Comparison of Policies with Respect to CVar (20,000 Simulated Runs and a Uniform \mathbb{P}_B on the Largest Inscribed Sphere of the Belief Space)



heuristic captures the essence of the chance-constrained policy in that it lies near the optimal policy, mirroring its performance in each simulated run. Overall, our goal to provide an alternative to the overconservatism and overoptimism of the minimax and minimin policies seems to be met by our percentile optimization technique, which is consistent with established robust optimization literature (see, e.g., Bertsimas and Sim 2004). Moreover, although our policies are generated from a fixed pessimism level (i.e., 95% chance-constrained), they perform well throughout the spectrum of optimism/pessimism in the CVar statistic.

Even in cases in which the chance-constrained policy is inferior to other policies with regard to the CVar statistic (e.g., the fourth row of Figure 8 with $X = (10, 10, 10)$, where the minimin policy is seen to perform best with regard to $\text{CVar}(0)$), we can see that fixed priority policies (e.g., those obtained under the minimin objective) miss out on the advantages of robustness that the chance-constrained policy offers throughout the optimism spectrum. Furthermore, percentile optimization is flexible: by modifying ϵ , we can change our policy's focus to be more or less

optimistic to the point of becoming a minimax and minimin policy itself (Proposition 2). A similar advantage is also gained in the APOMDP framework of Saghaian (2018), where α -maximin expected utility (α -MEU) preferences are used.

7.1. Real-World Application: ED Patient Prioritization

In most hospital EDs in the United States, patients upon arrival are sorted by means of an urgency-based triage system into one of (typically) five classes known as emergency severity index (ESI) levels. These ESI levels classify patients in descending order of urgency so that a patient of ESI 1, being in dire condition, is immediately treated, whereas patients of levels 4 and 5 are sent to a “fast-track” area to be treated. Therefore, the classes served by the main section of the ED (the majority of arrivals) are those with ESI levels 2 and 3 (see, e.g., Saghaian et al. 2012, 2014, and the references therein). We denote ESI 2 and 3 patients by “urgent” and “nonurgent” patients, respectively.

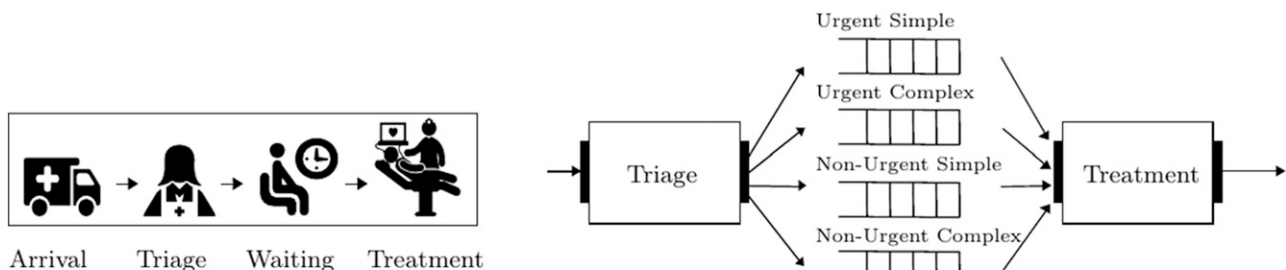
As patients wait to receive treatment, their condition may worsen over time and lead to adverse medical events. Sprivulis et al. (2006) and Plunkett et al. (2011) show that higher patient mortality is associated with longer waiting times prior to seeing a physician. Other research (e.g., an extremely large study on data of nearly 14 million patients by Guttman et al. 2011) indicate that the risk of adverse events (ROAE) for patients increases with higher waiting times leading to higher mortality and hospital admission rates. Therefore, with the objective of increasing patient safety, we consider the goal of minimizing average ROAE for ED patients and investigate optimal prioritization policies. To do so, we assume adverse events occur based on a Poisson process with a higher rate for urgent patients and note that ROAEs in this setting play the role of holding cost parameters in our multiclass queueing model introduced earlier. The same approach is used in Saghaian et al. (2014), in which the benefits of further stratifying these levels in terms of a patient’s *complexity* is discussed. Simple patients are those who experience only a single interaction with the physician and, thus, are more quickly treated by the ED than complex patients, whose treatment necessitates several interactions with the physician interspersed with various tests (CT scans, MRI, etc.).

Figure 9 (left) illustrates a schematic flow of patients as a multiclass queueing system. To analyze the multiclass queueing system of Figure 9 (right) in a traditional way, one needs to obtain point estimates of various parameters (e.g., service/treatment rates for each class), a task that is subject to inevitable errors.¹⁵ Furthermore, triaged urgency and complexity levels are subject to misclassifications, which further confuse the true parameter settings of the system. Although misclassifications can be included in the analysis when all of the parameters of the system are known, the misclassification probabilities themselves are also hard to quantify. These create parameter ambiguity, and one needs to use robust analyses to hedge against them. However, current ED patient prioritization policies are based on analyses that ignore such ambiguities.

To demonstrate the benefits of our percentile optimization approach, we now focus on two questions: How should EDs prioritize their patients given that they are faced with parameter ambiguity? How much benefit can they get by taking ambiguities into consideration? To answer these questions, we first model the ED from a broad perspective with nonstationary Poisson process arrivals and known service rates for all four classes: urgent simple (US), urgent complex (UC), nonurgent simple (NS), and nonurgent complex (NC) patients. In this way, we model the ED as a single “superserver” (i.e., with a pooled capacity that we estimate from our data set so as to match the input–output process of the ED as a whole). This allows us to gain insights into the questions we raised by noting that the ED queueing model of Figure 9 (right) is essentially a special case of our general model depicted in Figure 1 with $n = 4$.

Patient arrivals in an ED fluctuate throughout a given day, so we model these arrivals with a nonstationary Poisson process with hourly rates shown in Figure 16 in Online Appendix A, which depicts the actual time-dependent arrival rates to the ED based on our data set. Furthermore, because patient length of stay (LOS) in our data has a lognormal distribution, we fit lognormal service distributions to match the LOS of patients for each class of patients. Next, we design our cloud of models by perturbing the fitted rate parameters such that, for each class i with fitted rate $\hat{\mu}_{i,3}$, we incorporate four additional possible rate parameters so $\hat{\mu}_{i,1} < \hat{\mu}_{i,2} < \hat{\mu}_{i,3} < \hat{\mu}_{i,4} < \hat{\mu}_{i,5}$. Because

Figure 9. Patient Flow in Hospital Emergency Departments (Left: the Overall Flow; Right: the Multiclass Flow)



patients become fairly stable upon seeing a physician, we focus on adverse events in the waiting area of EDs and assume ROAE drops to zero once the treatment stage begins. Our model is nonpreemptive, which is a reflection of physicians' behavior in EDs: upon initiating treatment to a patient, they rarely pause treatment to serve a different patient. Because there is a possibility that the ROAE for simple patients differs from that of complex patients, we also consider a variety of such "cost" structures in our study.

Although this model allows for dynamic arrivals (unlike our model introduced in Section 3), we can still incorporate chance-constrained policies through the use of our heuristic and compare its performance to the complexity-based prioritization policy that serves classes US, UC, NS, and NC in descending priority (demonstrated to be optimal for EDs in Saghafian et al. 2014 when ambiguity is ignored), minimax, and minimin policies. To do so, we simply modify the Bayesian belief to also incorporate arrival data. We simulate these policies and track the nondiscounted ROAE by assuming that $\mathbb{P}_{\mathbf{B}}$ is uniform. The result of 20,000 simulated days expressed in terms of the CVar statistic is reported in Figure 10 (see Online Appendix A.7 for four additional ROAE settings and in-depth discussions).

A widely discussed topic in the literature surrounding EDs is the overcrowding issue (see e.g., Derlet and Richards 2000, Derlet et al. 2001, and Trzeciak and Rivers 2003) that stems from high arrival rates and limited resources (such as capacity, physicians, equipment, etc.). Overcrowding in EDs results in high ROAE that endangers patients. The third row of Figure 10 demonstrates how policies perform in overcrowded EDs by considering an ambiguity set with smaller service rates (in comparison with the other ambiguity sets). We note that percentile optimization, in comparison with other policies, is especially suited for studying patient prioritization in overcrowded EDs. This is because, under heavy congestion, chance-constrained policies learn faster because more classes are available to serve at any given time. Furthermore, as we show in Corollary 4 in Online Appendix B, the E_{μ} policy becomes asymptotically optimal when arrivals occur during intense bursts followed by lull periods. Because hospital EDs typically experience long periods of heavy traffic in the afternoon followed by little traffic after midnight (see the actual arrival pattern depicted in Figure 16 in Online Appendix A), this further establishes our approach in hospital ED applications. Using these results, we can make the following:

Observation 7 (High Traffic). Our percentile optimization approach performs well for prioritizing patients in EDs, especially in highly congested ones (e.g., those in busy research hospitals).

Also, Figure 10 shows that, once again, the chance-constrained policies nearly dominate the entire spectrum of the CVar statistic because they explicitly incorporate both learning and robustness. Hence, even though our stylized environment is less detailed than those ED flow models in studies such as Huang et al. (2015) and Saghafian et al. (2012, 2014, 2015) and the references therein (which feature patient feedback), these experiments indicate a performance advantage over complexity-based prioritization, which suggests implementation regardless of optimism/pessimism levels. Hence, to establish the potential benefits percentile optimization can offer to EDs over the current status quo, we make the following:

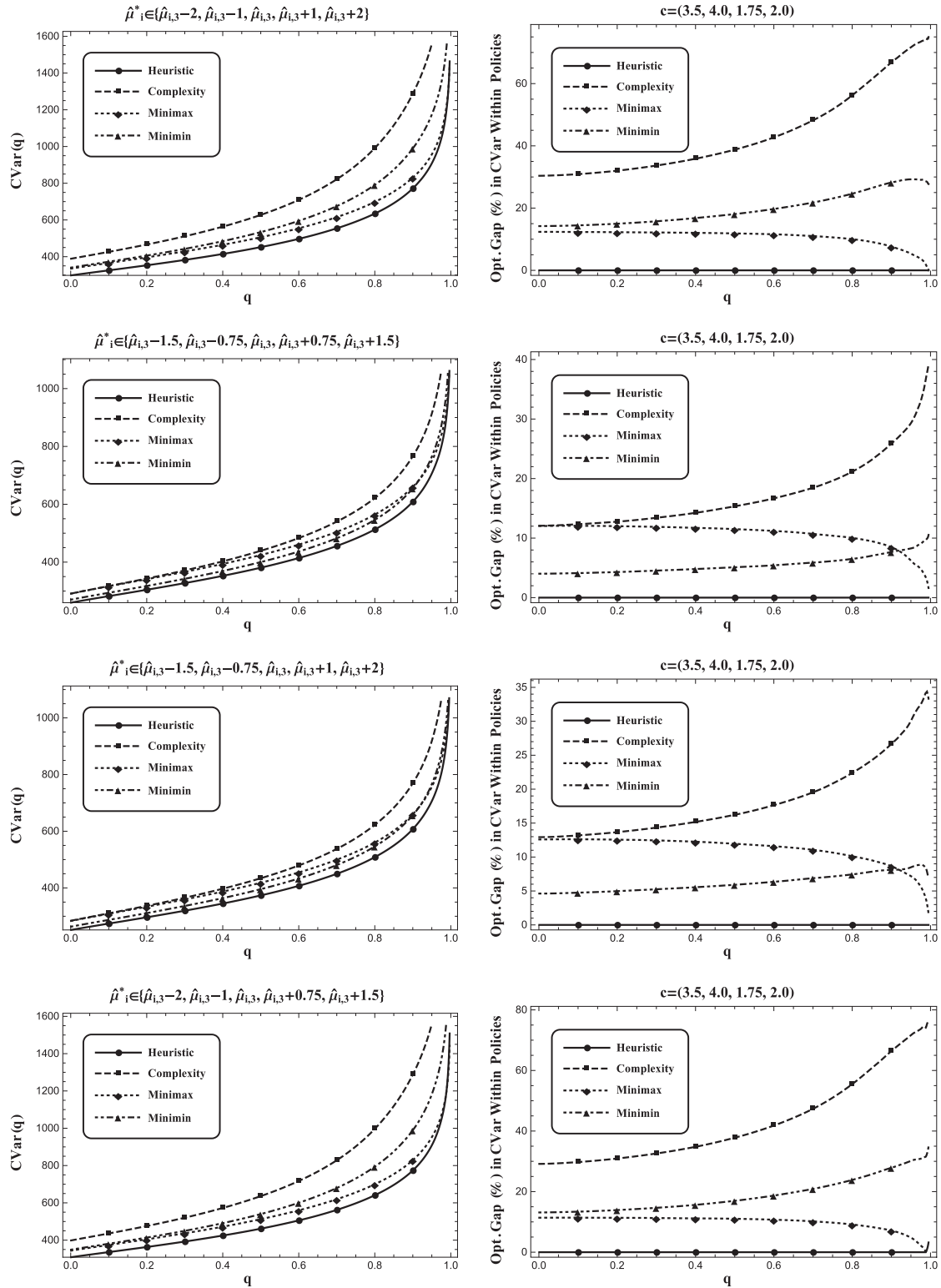
Observation 8 (Improved System Performance). Percentile optimization can improve the performance of EDs regardless of a manager's disposition.

In systems with high traffic, learning may occur at an advanced rate because it has available customers from each class a majority of the time the system is online. Hence, although static priority policies continue to serve the "wrong" classes (because of the underlying parameter ambiguity), the chance-constrained policy quickly identifies the optimal $c\mu$ priority using the observed values. This enhances the quality the robustness percentile optimization offers, especially because one is typically more concerned with overcrowded/busy systems (EDs with low traffic have short patient LOS naturally and are not in significant need of optimization).

Furthermore, our clearing system is a model often used to study queues undergoing overcrowded situations. Therefore, a more congested ED is a better fit to our original model, and in considering dynamic arrivals, we can reconfirm all the previous insights generated in the clearing environment. This further confirms the results of Section A.4, in which we show that most of the main insights gained from the clearing system holds for systems with dynamic arrivals.

Finally, we note that, in communities with unstable patient population characteristics, in which ED service rates or misclassification probabilities are more ambiguous, ED managers can incorporate percentile optimization to effectively hedge against such ambiguities. Moreover, percentile optimization is well suited to high levels of ambiguity. In our simulations, this is captured through modifying our cloud of models to incorporate larger differences in the fitted parameters (see the first row of Figure 10 and compare it with the second row). Hence, when patient population characteristics are unstable, percentile optimization stands out as a method that protects from negative consequences of focusing only on extreme outcomes while simultaneously learning from incoming data. This results in the following:

Figure 10. 20,000 Simulated Days in the ED for the Complexity-Based Prioritization, 95% E_{μ} Heuristic, Minimax, and Minimin Policies, When $\mathbb{P}_{\mathbf{B}}$ Is Uniform, and the Cloud of Models Perturbs the Fitted Service Rate $\hat{\mu}_{i,3}$ in Terms of Two-Hour Time Increments with $\mathbf{c} = (3.5, 4.0, 1.75, 2.0)$



Note. Triage levels US, UC, NS, and NC are denoted 1, 2, 3, and 4, respectively.

Observation 9 (Uncertain Population Characteristics). Percentile optimization can significantly help EDs that are placed in geographical areas with unstable or unknown patient population characteristics to better prioritize their patients.

8. Conclusion

Multiclass queues are versatile structures widely used in operations management that see a large variety of applications in both service and manufacturing sectors. In such environments, often exact parameter specification is rife with estimation errors that (if ignored) can cause system managers to implement wrong policies. We identify and implement a novel data-driven percentile optimization framework for use in POMDPs. Our method layers chance-constrained optimization on a nonrobust learning model, effectively enabling learning of the true system state parameters and allowing the manager to set an optimism level indicating the extent of protection against poor parameter scenarios the manager desires. We characterize the optimal policies to both the nonrobust and percentile problems and find that chance-constrained policies can be established via the nonrobust problem.

Because percentile optimization problems are typically computationally difficult, we introduce an analytically rooted heuristic that can be used to effectively incorporate robustness in managing large and complex service or manufacturing systems. To further improve computational tractability, we find asymptotically tight bounds to the nonrobust problem, which can be used to efficiently solve the percentile optimization problem.

Finally, we demonstrate the efficacy of our methods numerically in both stylized and realistic environments. Using real-world data collected from a leading hospital, we observe that our approach provides promising results in improving current patient flow policies, especially for overcrowded EDs or those facing unknown patient population characteristics. Because ED managers typically do not fully know the service rate parameters, traditional patient-flow policies based on queueing models that assume full service rate knowledge subject patients to higher risk than chance-constrained policies. Our work is the first to take into account the inevitable ambiguities in ED operations and sheds light on the dire consequences of ignoring such ambiguities.

Endnotes

- ¹ See, for example, Saghaian et al. (2015) for a recent review of various models used to optimize patient flow and improve ED operations.
- ² Although we mainly focus on a queueing model, our approach can be used for the general class of Bayesian decision-making problems in which the decision maker faces ambiguity with respect to parameters that shape the decision maker's prior (see Corollaries 2 and 3 in Online Appendix B).
- ³ Clearing systems are typically used to model busy periods by focusing on the customers/jobs already in the system. The goal is then to clear the system with the minimum cost.
- ⁴ In Section 7, we relax the exponential distribution assumption. For instance, our data shows that service times in EDs are close to lognormal. As we show, our main insights and heuristic control procedures remain effective even when the service times are not exponential.
- ⁵ This is indeed a general criticism to Bayesianism and goes well beyond the queueing setting of this paper.
- ⁶ One may criticize the use of the percentile objective because of the potential ambiguity of \mathbb{P}_B ; however, it should be noted that this is a second-order distribution, and perturbations in \mathbb{P}_B result in very similar convex floating bodies, which is the geometric structure investigated in Section 4 that generates our optimal robust policies.
- ⁷ Although all nonrobust policies for a finite-horizon POMDP are guaranteed to evaluate as linear functions over \mathcal{B} , an infinite-horizon POMDP value function is not always guaranteed to be piecewise-linear (see, e.g., White and Harrington 1980).
- ⁸ For these randomized policies, we disallow policies that are not picked at time zero for the purpose of targeting specific contours of the value function.
- ⁹ We note that, if \mathcal{L}_ϵ is nonempty, $\delta\mathcal{L}_\epsilon$ always exists because closed, convex, and compact sets are equal to the convex hull of their boundary.
- ¹⁰ For additional discussion and examples of convex floating bodies, see Online Appendix A.6.
- ¹¹ This does not imply that $\arg \max_{\mathbf{b} \in \delta\mathcal{L}_\epsilon} V(\mathbf{X}, \mathbf{b}) = \arg \min_{\mathbf{b} \in \delta\mathcal{L}_\epsilon} \|\mathbf{b}_0 - \mathbf{b}\|$; $V(\mathbf{X}, \mathbf{b})$ is only assured to be nonincreasing on line segments connected to \mathbf{b}_0 .
- ¹² For two belief points, $\mathbf{b}, \hat{\mathbf{b}} \in \mathcal{B}$ with all positive components in the setting in which $n = 2, m_1 = 2$, and $m_2 = 2$, the KL divergence is $D_{KL}(\mathbf{b} \parallel \hat{\mathbf{b}}) = \sum_{i=1}^2 \sum_{j=1}^2 b_{1,i} b_{2,j} \log \frac{b_{1,i} b_{2,j}}{\hat{b}_{1,i} \hat{b}_{2,j}}$.
- ¹³ For instance, one would expect the minimax policy to perform well in comparison with other policies at $\text{CVar}(1)$.
- ¹⁴ The associated confidence intervals are tight, so we only show the averages.
- ¹⁵ Even after using a large data set that we have collected from a leading U.S. hospital, one which includes data about more than 18,000 patient visits, we see that our point estimates are not reliable for various reasons, including the large variation among patient characteristics as well as the need to estimate parameters for each patient class separately.

References

- Bagnell J, Ng AY, Schneider J (2001) Solving uncertain Markov decision problems. Technical Report CMU-RI-TR-01-25, Robotics Institute, Carnegie Mellon University, Pittsburgh.
- Bandi C, Bertsimas D (2012) Tractable stochastic analysis in high dimensions via robust optimization. *Math. Programming* 134(1):23–70.
- Bandi C, Bertsimas D, Youssef N (2015) Robust queueing theory. *Oper. Res.* 63(3):676–700.
- Bassamboo A, Zeevi A (2009) On a data-driven method for staffing large call centers. *Oper. Res.* 57(3):714–726.

- Bertsekas D (1995) *Dynamic Programming and Optimal Control*, vol. 1 (Athena Scientific, Belmont, MA).
- Bertsimas D, Sim M (2004) The price of robustness. *Oper. Res.* 52(1):35–53.
- Bertsimas D, Thiele A (2006) A robust optimization approach to inventory theory. *Oper. Res.* 54(1):150–168.
- Bertsimas D, Gamarnik D, Rikun AA (2011) Performance analysis of queueing networks via robust optimization. *Oper. Res.* 59(2):455–466.
- Bertsimas D, Gupta V, Kallus N (2018) Data-driven robust optimization. *Math. Programming* 167(2):235–292.
- Bertsimas D, Pachamanova D, Sim M (2004) Robust linear optimization under general norms. *Oper. Res. Lett.* 32(6):510–516.
- Buyukkoc C, Varaiya P, Walrand J (1985) The $c\mu$ rule revisited. *Adv. Appl. Probab.* 17(1):237–238.
- Chen Y, Farias V (2013) Simple policies for dynamic pricing with imperfect forecasts. *Oper. Res.* 61(3):612–624.
- Chow Y, Ghavamzadeh M, Janson L, Pavone M (2018) Risk-constrained reinforcement learning with percentile risk criteria. *J. Machine Learn. Res.* 18:1–51.
- Delage E, Mannor S (2007) Percentile optimization in uncertain Markov decision processes with application to efficient exploration. *Proc. 24th Internat. Conf. Machine Learn.* (ACM, New York), 225–232.
- Delage E, Mannor S (2010) Percentile optimization for Markov decision processes with parameter uncertainty. *Oper. Res.* 58(1):203–213.
- Derlet RW, Richards JR (2000) Overcrowding in the nation’s emergency departments: Complex causes and disturbing effects. *Ann. Emergency Medicine* 35(1):63–68.
- Derlet RW, Richards JR, Kravitz RL (2001) Frequent overcrowding in U.S. emergency departments. *Acad. Emergency Medicine* 8(2):151–155.
- Dupin C (1822) *Applications de Géométrie et de Mécanique* (Bachelier, successeur de Mme. Ve. Courcier, libraire).
- Fresen D (2013) A multivariate Gnedenko law of large numbers. *Ann. Probab.* 41(5):3051–3080.
- Guttmann A, Schull MJ, Vermeulen MJ, Stukel TA (2011) Association between waiting times and short term mortality and hospital admission after departure from emergency department: Population based cohort study from Ontario, Canada. *British Medical J.* 342:d2983.
- Hansen L, Sargent T (2007) Recursive robust estimation and control without commitment. *J. Econom. Theory* 136(1):1–27.
- Huang J, Carmeli B, Mandelbaum A (2015) Control of patient flow in emergency departments, or multiclass queues with deadlines and feedback. *Oper. Res.* 63(4):892–908.
- Iyengar GN (2005) Robust dynamic programming. *Math. Oper. Res.* 30(2):257–280.
- Jain A, Lim A, Shanthikumar JG (2010) On the optimality of threshold control in queues with model uncertainty. *Queueing Systems* 65(2):157–174.
- Lagoa CM, Li X, Szaier M (2005) Probabilistically constrained linear programs and risk-adjusted controller design. *SIAM J. Optim.* 15(3):938–951.
- Lim AEB, Shanthikumar JG, Vahn G (2012) Robust portfolio choice with learning in the framework of regret: Single-period case. *Management Sci.* 58(9):1732–1746.
- Lippman SA (1975) Applying a new device in the optimization of exponential queueing systems. *Oper. Res.* 23(4):687–710.
- Littman ML, Goldsmith J, Mundhenk M (1998) The computational complexity of probabilistic planning. *J. Artificial Intelligence Res.* 9(1):1–36.
- Mannor S, Simester D, Sun P, Tsitsiklis JN (2007) Bias and variance approximation in value function estimates. *Management Sci.* 53(2):308–322.
- Mundhenk M, Goldsmith J, Lusena C, Allender E (2000) Complexity of finite-horizon Markov decision process problems. *J. ACM* 47(4):681–720.
- Nemirovski A, Shapiro A (2006) Convex approximations of chance constrained programs. *SIAM J. Optim.* 17(4):969–996.
- Nilim A, El Ghaoui L (2005) Robust control of Markov decision processes with uncertain transition matrices. *Oper. Res.* 53(5):780–798.
- Osogami T (2015) Robust partially observable Markov decision process. *Proc. Machine Learn. Res.* 37:106–115.
- Papadimitriou CH, Tsitsiklis JN (1987) The complexity of Markov decision processes. *Math. Oper. Res.* 12(3):441–450.
- Pedarsani R, Walrand J, Zhong Y (2014) Robust scheduling and congestion control for flexible queueing networks. *2014 Internat. Conf. Comput. Networking Comm. (ICNC)* (IEEE, Piscataway, NJ), 467–471.
- Plunkett PK, Byrne DG, Breslin T, Bennett K, Silke B (2011) Increasing wait times predict increasing mortality for emergency medical admissions. *Eur. J. Emergency Medicine* 18(4):192–196.
- Ross S, Pineau J, Chaib-draa B, Kreitmann P (2011) A Bayesian approach for learning and planning in partially observable Markov decision processes. *J. Machine Learn. Res.* 12:1729–1770.
- Saghafian S (2018) Ambiguous partially observable Markov decision processes: Structural results and applications. *J. Econom. Theory* 178:1–35.
- Saghafian S, Veatch MH (2016) A $c\mu$ rule for two-tiered parallel servers. *IEEE Trans. Automatic Control* 61(4):1046–1050.
- Saghafian S, Austin G, Traub SJ (2015) Operations research/management contributions to emergency department patient flow optimization: Review and research prospects. *IIE Trans. Healthcare Systems Engrg.* 5(2):101–123.
- Saghafian S, Hopp WJ, Van Oyen MP, Desmond JS, Kronick SL (2012) Patient streaming as a mechanism to improve responsiveness in emergency departments. *Oper. Res.* 60(5):1080–1097.
- Saghafian S, Hopp WJ, Van Oyen MP, Desmond JS, Kronick SL (2014) Complexity-augmented triage: A tool for improving patient safety and operational efficiency. *Manufacturing Service Oper. Management* 16(3):329–345.
- Smallwood R, Sondik EJ (1973) The optimal control of partially observable Markov processes over a finite horizon. *Oper. Res.* 21(5):1071–1088.
- Sondik EJ (1971) The optimal control of partially observable Markov processes. Unpublished doctoral dissertation, Stanford University, Stanford, CA.
- Sprivilis PC, Da Silva JA, Jacobs IG, Frazer AR, Jelinek GA (2006) The association between hospital overcrowding and mortality among patients admitted via Western Australian emergency departments. *Medical J. Australia* 184(5):208–212.
- Su H (2006) Robust fluid control of multiclass queueing networks. Unpublished master’s thesis, Massachusetts Institute of Technology, Cambridge.
- Thrun S (1999) Monte Carlo POMDPs. *Adv. Neural Inform. Processing Systems* 12:1064–1070.
- Trzeciak S, Rivers EP (2003) Emergency department overcrowding in the United States: An emerging threat to patient safety and public health. *Emergency Medicine J.* 20(5):402–405.
- Van Mieghem JA (1995) Dynamic scheduling with convex delay costs: The generalized $c\mu$ rule. *Ann. Appl. Probab.* 5(3):809–833.
- White C, Harrington D (1980) Application of Jensen’s inequality to adaptive suboptimal design. *J. Optim. Theory Appl.* 32(1):89–99.
- Wiesemann W, Kuhn D, Rustem B (2013) Robust Markov decision processes. *Math. Oper. Res.* 38(1):153–183.
- Zhang H (2010) Partially observable Markov decision processes: A geometric technique and analysis. *Oper. Res.* 58(1):214–228.