# On the Estimation of Buffer Overflow Probabilities from Measurements

Ioannis Ch. Paschalidis, *Member, IEEE,* and Spyridon Vassilaras

*Abstract*—We propose estimators of the buffer overflow probability in queues fed by a Markov-modulated input process and serviced by an autocorrelated service process. These estimators are based on large-deviations asymptotics for the overflow probability. We demonstrate that the proposed estimators are less likely to underestimate the overflow probability than the estimator obtained by certainty equivalence. As such, they are appropriate in situations where the overflow probability is associated with quality of service (QoS) and we need to provide firm QoS guarantees. We also show that as the number of observations increases to infinity the proposed estimators converge with probability one to the appropriate target, and thus, do not lead to underutilization of the system in this limit.

*Index Terms*—Empirical measure, estimation, large deviations, Markov-modulated processes, Sanov's theorem.

## I. INTRODUCTION

ESTIMATING small buffer overflow probabilities has found important applications in a variety of application areas. Two prominent examples include the control of *communication networks* and *manufacturing systems.*

In modern high-speed communication networks, congestion manifests itself as buffer overflows; the *quality of service (QoS)* faced by various connections can be quantified by the buffer overflow probability. To provide QoS guarantees the so-called *effective bandwidth* admission control mechanism has been proposed [1]–[6]. Briefly, effective bandwidth is a number between the peak and average rate of a connection chosen in such a way that when connections are allocated their effective bandwidth in an appropriately dimensioned buffer, the buffer overflow probability stays below a given small level (say, on the order of $10^{-7}$). Real-time applications can tolerate such small frequencies of congestion phenomena.

In make-to-stock manufacturing systems and supply chains, on the other hand, the objective is to control the inventory in order to avoid stockouts (see [7], [28], [8], [9]). In such systems, demand is met from a finished goods inventory, and it is backordered if inventory is not available. The stockout probability quantifies the QoS encountered by customers. It can be shown that this probability is equal to a buffer overflow probability in a corresponding make-to-order system [7], [28], [9]. Thus, the problem of estimating the stockout probability can be transformed into one of estimating a buffer overflow probability.

In both these applications, we are interested in estimating buffer overflow probabilities that are very small. Moreover, arrival and service processes are typically autocorrelated (to model bursty traffic in communication networks, and to accommodate realistic demand scenarios and model failure-prone production facilities in supply chains). As a result, obtaining exact analytic expressions is intractable; it is, therefore, natural to focus on asymptotic regimes. To that end, *large deviations* theory (see [10]) has been an important analytical tool. Large deviations techniques have recently been applied to a variety of problems in telecommunications (see [11], [12], [6] and references therein) and supply chains (see [7], [28], [9]).

All of this large deviations work assumes detailed knowledge of models for the arrival and service processes. In practice, however, such traffic models are not known *a priori* and have to be estimated from real observations. Consequently, one approach for estimating buffer overflow probabilities is to assume a certain traffic model, estimate the parameters of the model from real observations, and then calculate the overflow probability using large deviations techniques. Consider, for example, the case of a *deterministic Markov-modulated* source (i.e., a source which has a constant rate at every state of the underlying Markov chain) and let $g(\Xi)$ be the overflow probability when this source is fed to a certain buffer, where $\Xi$ denotes the transition probability matrix of the underlying Markov chain. (We assume that the rate at every state and the characteristics of the buffer are given, thus, we do not explicitly denote the dependence of the overflow probability on these quantities). Let $\hat{\Xi}$ be an unbiased estimator of the transition probability matrix $\Xi$, that is, $\boldsymbol{E}[\hat{\Xi}] = \Xi$. Suppose that we use $g(\hat{\Xi})$ as an estimator of the overflow probability. An important observation is that due to the nonlinearity of $g(\cdot)$, $\boldsymbol{E}[g(\hat{\Xi})]$ is not necessarily equal to $g(\boldsymbol{E}[\hat{\Xi}]) = g(\Xi)$. That is, a *certainty equivalence* approach can lead to an erroneous estimate.

Measurement-based, model-free admission control has received some attention recently [13]–[15]. In particular, the authors in [13] and [14] consider measurement-based admission control schemes in bufferless multiplexers. They investigate how the overflow probability is affected by estimation errors relying on approximations based on the central limit theorem. Duffield [15] extends these results by considering large-deviations asymptotics in a multiplexer with nonzero buffer. The

I. Ch. Paschalidis is with the Department of Manufacturing Engineering, Boston University, Boston, MA 02215 USA (e-mail: yannisp@bu.edu; url: http://ionix.bu.edu).

S. Vassilaras is with the College of Engineering, Boston University, Boston, MA 02215 USA (e-mail: svassila@bu.edu).

approach in [15] relies on some earlier work in [16] which develops an estimator of the large-deviations rate function of the arrival process directly from measurements, without assuming any particular stochastic model for that process. Finally, the authors in [17] argue, as we do, that certainty equivalence can lead to erroneous estimates for rare event probabilities and propose a Bayesian approach in a simpler [independent and identically distributed (i.i.d.)] context than ours.

In this paper, we will focus on queues fed by Markov-modulated arrival processes and develop new estimates of overflow probabilities that take estimation errors into account. To that end, we will establish an inverse of Sanov's theorem (see [10]) for Markov chains in the large-deviations regime. Intuitively, the proposed estimators suggest that we should quote a quantity that is larger than $g(\hat{\boldsymbol{\Xi}})$ to guard against estimation errors. We will provide analytical and numerical evidence that the proposed estimators are "safe" even when based on relatively few observations, meaning that they do not lead to substantial underestimation of the overflow probability that can compromise QoS. Still, they are *consistent* in the sense that they converge to $g(\boldsymbol{\Xi})$ with probability one (w.p. 1) as the number of observations tends to infinity. Among our main contributions we consider the following.

- The fact that the proposed estimators are "safe" for relatively few observations, but eventually (as the number of observations tends to infinity) "learn" the true overflow probability. Providing "safe" estimates from relatively few observations is crucial since it allows the estimator to quickly adapt to level shifts in the input, and thus, accommodate nonstationary scenarios. One of the proposed estimators has the same structure as the estimator suggested in [15]. In [15], however, the estimation process is different; a limiting *log*-moment-generating function is estimated directly from the data. In our work, we impose additional structure by considering Markov-modulated processes (which is appropriate in some applications, e.g., inventory control in supply chains [7], [28], [9]). Moreover, we derive additional estimators based on higher moments of a loss measure. These estimators are shown to be more appropriate ("safer," yet still consistent) for our purposes.

- The development of efficient algorithms to compute the proposed estimators. Computing the estimators amounts to solving a nonlinear programming problem; efficient computation is an issue, especially in high-dimensional instances. We study the structure of this problem and develop algorithms that perform well in practice.

- An inverse of Sanov's theorem for Markov chains that we establish. This generalizes a result in [18] which deals with i.i.d. processes. It is also related to the result in [19], which derives a large deviations rate function of the posterior distribution for the transition probabilities of a Markov chain with parameter the "true" transition probabilities. We provide an independent proof based on first principles, where in our version the limit of the empirical measure enters the rate function; this is motivated by the application where the "true" transition probabilities are not known.

Moreover, our result applies to Markov chains with zeros in the transition probability matrix, a case which is not discussed in [19].

On the organization of this paper, we start in Section II with some preliminaries on large deviations and Sanov's theorem. In Section III, we discuss analytical large-deviations expressions for buffer overflow probabilities and formulate the estimation problem. In Section IV, we analyze the estimation process in the large-deviations regime and establish an inverse of Sanov's theorem for Markov chains. We use this result in Section V to propose our estimators of buffer overflow probabilities. We discuss issues related to the computation of the proposed estimators in Section VI. Comparisons between various estimators and some illustrative numerical results are in Section VII. Conclusions are in Section VIII.

## II. PRELIMINARIES

In the form of background on large deviations and to establish some of our notation, we first review some basic results. Consider a sequence of i.i.d. random variables $X_i$, $i \geq 1$, with mean $\boldsymbol{E}[X_1] = \overline{X}$. The strong law of large numbers asserts that $\frac{\sum_{i=1}^n X_i}{n}$ converges to $\overline{X}$, as $n \to \infty$, w.p. 1. Thus, for large $n$, the event $\sum_{i=1}^n X_i \geq na$, where $a > \overline{X}$ (or $\sum_{i=1}^n X_i \leq na$, for $a < \overline{X}$) is a rare event. In particular, its probability behaves as $e^{-nr(a)}$, as $n \to \infty$, where the function $r(\cdot)$ determines the rate at which the probability of this event is diminishing. Cramér's theorem [20] determines $r(\cdot)$, and is considered the first large deviations statement.

Consider next a sequence $\{S_1, S_2, \ldots\}$ of random variables, with values in $\mathbb{R}$ and define

$$\Lambda_n(\theta) \triangleq \frac{1}{n} \log \boldsymbol{E}[e^{\theta S_n}].$$

For the applications we have in mind, $S_n$ is a partial sum process. Namely, $S_n = \sum_{i=1}^n X_i$, where $X_i, i \geq 1$, are identically distributed, possibly *dependent* random variables. We will be making the following assumption.

*Assumption A:*
1) The limit

$$\Lambda(\theta) \triangleq \lim_{n \to \infty} \Lambda_n(\theta) = \lim_{n \to \infty} \frac{1}{n} \log \boldsymbol{E}[e^{\theta S_n}]$$

exists for all $\theta$, where $\pm\infty$ are allowed both as elements of the sequence $\Lambda_n(\theta)$ and as limit points.
2) The origin is in the interior of the domain

$$D_\Lambda \triangleq \{\theta | \Lambda(\theta) < \infty\}$$

of $\Lambda(\theta)$.
3) $\Lambda(\theta)$ is differentiable in the interior of $D_\Lambda$ and the derivative tends to infinity as $\theta$ approaches the boundary of $D_\Lambda$.
4) $\Lambda(\theta)$ is lower semicontinuous, i.e.,

$$\liminf_{\theta_n \to \theta} \Lambda(\theta_n) \geq \Lambda(\theta)$$

for all $\theta$.

Let us define

$$\Lambda^*(a) \triangleq \sup_\theta (\theta a - \Lambda(\theta)) \tag{1}$$

which is the Legendre transform of $\Lambda(\cdot)$. $\Lambda(\cdot)$ and $\Lambda^*(\cdot)$ are convex duals (see [21]), namely, along with (1), it also holds that

$$\Lambda(\theta) = \sup_a (\theta a - \Lambda^*(a)). \tag{2}$$

The function $\Lambda^*(\cdot)$ is convex and lower semicontinuous (see [10]).

Cramér's theorem has been extended to cover autocorrelated processes. In particular, under Assumption A, the Gärtner–Ellis theorem (see [10]) establishes that $\{S_n\}$ satisfies a *large deviations principle (LDP)* with *rate function* $\Lambda^*(\cdot)$. More specifically, this theorem intuitively asserts that for large enough $n$ and for small $\epsilon > 0$

$$\boldsymbol{P}[S_n \in (na - n\epsilon,\, na + n\epsilon)] \sim e^{-n\Lambda^*(a)}.$$

The set of processes satisfying Assumption A (and hence Gärtner–Ellis theorem) is large enough to include renewal, Markov-modulated, and stationary processes with mild mixing conditions. Such processes can model "burstiness" and are commonly used in modeling the input traffic to communication networks. They have recently being used in modeling demand and the production process in manufacturing systems [7], [28], [9].

For discrete random variables, a stronger result than Cramér's theorem is Sanov's theorem which deals with large deviations of the empirical measure. Consider a sequence $\boldsymbol{Y} = Y_1, Y_2, \ldots, Y_n$ of i.i.d. random variables taking values from a finite alphabet $\mathcal{A} = \{a_1, \ldots, a_M\}$ with $M$ elements. Let $\boldsymbol{\mu} = (\mu(a_1), \ldots, \mu(a_M))$ denote the probability law for $Y_1$, where $\boldsymbol{\mu}$ is an element of the standard $M$-dimensional probability simplex (i.e., $\boldsymbol{\mu} \in M_1(\mathcal{A}) \triangleq \{\boldsymbol{r} \in \mathbb{R}^M | r_i \in [0, 1], i = 1, \ldots, M, \sum_{i=1}^{M} r_i = 1\}$). The empirical measure of the sequence $Y_1, Y_2, \ldots, Y_n$ is given by

$$\boldsymbol{\mathcal{E}}_n^{\boldsymbol{Y}} = (\mathcal{E}_n^{\boldsymbol{Y}}(a_1), \ldots, \mathcal{E}_n^{\boldsymbol{Y}}(a_M))$$

where

$$\mathcal{E}_n^{\boldsymbol{Y}}(a_i) = \frac{1}{n} \sum_{j=1}^{n} \mathbf{1}_{a_i}(Y_j), \qquad i = 1, \ldots, M$$

and where $\mathbf{1}_{a_i}(Y_j)$ is the indicator function of $Y_j$ being equal to $a_i$. Thus, $\mathcal{E}_n^{\boldsymbol{Y}}(a_i)$ is equal to the fraction of occurrences of $a_i$ in the sequence.

Sanov's theorem (see [10]) establishes that $\boldsymbol{\mathcal{E}}_n^{\boldsymbol{Y}}$ satisfies an LDP with rate function given by the relative entropy

$$H(\boldsymbol{\nu}|\boldsymbol{\mu}) \triangleq \sum_{i=1}^{M} \nu(a_i) \log \frac{\nu(a_i)}{\mu(a_i)}. \tag{3}$$

In particular, we have the following.

*Theorem II.1 (Sanov's):* For every set $\Gamma$ of probability vectors in $M_1(\mathcal{A})$

$$- \inf_{\boldsymbol{\nu} \in \Gamma^\circ} H(\boldsymbol{\nu}|\boldsymbol{\mu}) \le \liminf_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}\left[\boldsymbol{\mathcal{E}}_n^{\boldsymbol{Y}} \in \Gamma\right]$$
$$\le \limsup_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}\left[\boldsymbol{\mathcal{E}}_n^{\boldsymbol{Y}} \in \Gamma\right]$$
$$\le - \inf_{\boldsymbol{\nu} \in \Gamma} H(\boldsymbol{\nu}|\boldsymbol{\mu})$$

where $\Gamma^\circ$ denotes the interior of $\Gamma$.

Intuitively, the empirical measure is "close" to $\boldsymbol{\nu}$ with probability behaving as $e^{-nH(\boldsymbol{\nu}|\boldsymbol{\mu})}$. Notice that this probability is diminishing as $\boldsymbol{\nu}$ is "further away" from the *a priori* measure $\boldsymbol{\mu}$ since the relative entropy can be interpreted as distance[1] (Kullback–Leibler distance, see [22]).

On a notational remark, in the rest of the paper we will be denoting by $\Lambda_X(\cdot)$ and $\Lambda_X^*(\cdot)$ the limiting log-moment-generating function and the large-deviations rate function, respectively, of the process $X$. Moreover, vectors will be denoted by bold characters and will be assumed to be column vectors, unless otherwise specified.

## III. LARGE DEVIATIONS OF A G/G/1 QUEUE AND PROBLEM FORMULATION

In this section, we review work on the large deviations behavior of a G/G/1 queue under a Markov-modulated arrival process. In particular, we state an asymptotic expression for the probability that the queue length process exceeds a certain threshold $U$. This probability is the quantity we wish to estimate from measurements.

Consider a single class queue. We will be using a discrete-time model, where time is divided into time slots. We let $A_i$, $i \in \mathbb{Z}$, denote the aggregate number of customers that enter the queue at time $i$. The queue has an infinite buffer and is serviced according to a general service process which can clear up to $B_i$ customers during the time interval $[i, i + 1]$. We assume that the stochastic processes $\{A_i; i \in \mathbb{Z}\}$ and $\{B_i; i \in \mathbb{Z}\}$ are stationary, possibly autocorrelated, and mutually independent processes that satisfy Assumption A.

We denote by $L_i$ the queue length at time $i$ (without counting arrivals at time $i$). We assume that the server uses a work-conserving policy (i.e., the server never stays idle when there is work in the system) and that

$$\boldsymbol{E}[A_1] < \boldsymbol{E}[B_1] \tag{4}$$

which by stationarity carries over to all $i$. We also assume that the queue length process $\{L_i, i \in \mathbb{Z}\}$ is stationary. To simplify the analysis, we consider a discrete-time "fluid" model, meaning that we will be treating $A_i$, $L_i$, and $B_i$ as nonnegative real numbers (the amount of fluid entering, in queue, or served).

An LDP for the queue length process has been established in [23], [11] and is given in the next proposition. In preparation for the result, consider a convex function $g(u)$ with the property $g(0) = 0$. We define the *largest root* of $g(u)$ to be the solution of

---

[1]The relative entropy is nonnegative and equal to zero if and only if $\boldsymbol{\mu} = \boldsymbol{\nu}$. However, it is not a true metric since it is not symmetric and does not satisfy the triangle inequality.

the optimization problem $\sup_{u:\,g(u)<0} u$. If $g(\cdot)$ has a negative derivative at $u=0$, there are two cases: either $g(\cdot)$ has a single positive root or it stays below the horizontal axis $u=0$, for all $u>0$. In the latter case, we will say that $g(\cdot)$ has a root at $u=\infty$.

*Proposition III.1:* The steady-state queue length process $L_i$ satisfies

$$\lim_{U\to\infty}\frac{1}{U}\log\boldsymbol{P}[L_i\geq U]=-\theta^* \qquad (5)$$

where $\theta^*>0$ is the largest root of the equation

$$\Lambda_A(\theta)+\Lambda_B(-\theta)=0. \qquad (6)$$

More intuitively, for large enough $U$ we have

$$\boldsymbol{P}[L_i\geq U]\sim e^{-U\theta^*}.$$

This expression can be used to estimate the overflow probability in a queue with a finite buffer of size $U$. Kelly [5] establishes that the latter probability has the same asymptotic decay rate (same exponent) as $\boldsymbol{P}[L_i\geq U]$.

### A. Markov-Modulated Arrivals

Assume next that the arrival process is Markov-modulated. More specifically, consider an irreducible Markov chain with $M$ states $1, 2, \ldots, M$ and transition probability matrix $\boldsymbol{\Xi}=\{p(k,j)\}_{k,j=1}^M$. We will be using the notation $\boldsymbol{p}'=(\boldsymbol{p}_1',\ldots,\boldsymbol{p}_M')$, where $\boldsymbol{p}_i$ is the $i$th row of $\boldsymbol{\Xi}$ and prime denotes transpose. The Markov chain makes one transition per time slot; let $Y_i$ be the state at time $i$. The number of arrivals at time $i$ is a deterministic function of the state, i.e., $A_i=f_A(Y_i)$.

In [10, Sec. 3.1.1] it is established that the limiting log-moment-generating function of the arrival process $A$ is given by

$$\Lambda_A(\theta,\boldsymbol{p})=\log\rho(\boldsymbol{\Pi}_{\theta,\boldsymbol{p}}^A) \qquad (7)$$

where $\rho(\boldsymbol{\Pi}_{\theta,\boldsymbol{p}}^A)$ denotes the Perron–Frobenius eigenvalue of the $M\times M$ matrix $\boldsymbol{\Pi}_{\theta,\boldsymbol{p}}^A$ which has elements

$$\pi_{\theta,\boldsymbol{p}}^A(k,j)=p(k,j)e^{\theta f_A(j)}, \qquad k,j=1,\ldots,M. \qquad (8)$$

(In this Markov-modulated case, we are using notation that explicitly denotes the dependence of $\Lambda_A(\theta,\boldsymbol{p})$ and $\boldsymbol{\Pi}_{\theta,\boldsymbol{p}}^A$ on the transition probabilities $\boldsymbol{p}$.) Notice that because the quantities $e^{\theta f_A(j)}$ are always positive the irreducibility of $\boldsymbol{\Xi}$ implies that $\boldsymbol{\Pi}_{\theta,\boldsymbol{p}}^A$ is irreducible. This calculation of $\Lambda_A(\theta,\boldsymbol{p})$ can be easily extended to the case where $f_A(Y_i)$ is a random function of the state $Y_i$ (see [10]).

### B. Estimating the Overflow Probability

We are interested in estimating the overflow probability $\boldsymbol{P}[L_i\geq U]$ from measurements. In particular, we will be assuming that we have perfect knowledge of the service process $B$, and that we can observe the states of the Markov chain associated with the arrival process. That is, we do know $M$ and the function $f_A(\cdot)$, but the transition probability matrix $\boldsymbol{\Xi}$ is unknown.

One way of estimating the overflow probability $\boldsymbol{P}[L_i\geq U]$ is to form an estimate $\hat{\boldsymbol{\Xi}}$ of the transition probability matrix, obtain the Perron–Frobenius eigenvalue of the corresponding matrix $\boldsymbol{\Pi}_{\theta,\hat{\boldsymbol{p}}}^A$, and use $\rho(\boldsymbol{\Pi}_{\theta,\hat{\boldsymbol{p}}}^A)$ in Proposition III.1 to obtain an estimate

$\hat{\theta}^*$ of the decay rate. This procedure yields an estimate $e^{-U\hat{\theta}^*}$ of the overflow probability. A problem with this approach is that even if $\hat{\boldsymbol{\Xi}}$ is an unbiased estimate of $\boldsymbol{\Xi}$, $e^{-U\hat{\theta}^*}$ is not necessarily unbiased since it is a nonlinear function of the transition probabilities in $\hat{\boldsymbol{\Xi}}$. That is, a *certainty equivalence* approach leads to an erroneous estimate. Our objective in this paper is to develop alternative and more appropriate estimates of the overflow probability than $e^{-U\hat{\theta}^*}$.

## IV. LARGE-DEVIATIONS ANALYSIS OF THE ESTIMATION PROCESS

We start the analysis by studying the large-deviations behavior of the estimation for a Markov-modulated process. In particular, we determine the large-deviations rate function of the *a posteriori* measure for the underlying Markov chain, given the observed empirical measure.

In the simpler case where we are observing the realization of a sequence of i.i.d. random variables, [18] provides the result. In particular, let $\boldsymbol{Y}$ denote a sequence $Y_1, Y_2, \ldots$ of i.i.d. random variables taking values from a finite alphabet $\mathcal{A}$, with probability law $\boldsymbol{\mu}\in M_1(\mathcal{A})$. The probability law $\boldsymbol{\mu}$ is not known; we have instead a prior distribution $\phi_{\boldsymbol{\mu}^-}\in M_1(M_1(\mathcal{A}))$. Given $n$ observations $Y_1, Y_2, \ldots, Y_n$ the posterior distribution is a function of the empirical measure; we will be denoting by $\phi_{\boldsymbol{\mu}_n^+(\boldsymbol{\mathcal{E}}_n^Y)}$ the posterior distribution and by $\boldsymbol{\mu}_n^+(\boldsymbol{\mathcal{E}}_n^Y)$ the corresponding random variable. Assume that as $n\to\infty$, $\boldsymbol{\mathcal{E}}_n^Y$ converges weakly to some $\boldsymbol{\nu}$ in the support of $\phi_{\boldsymbol{\mu}}^-$. In [18], it is established that for any prior $\phi_{\boldsymbol{\mu}}^-$, $\boldsymbol{\mu}_n^+(\boldsymbol{\mathcal{E}}_n^Y)$ satisfies an LDP with rate function

$$I_1(\boldsymbol{\mu})=\begin{cases} H(\boldsymbol{\nu}|\boldsymbol{\mu}), & \text{if }\boldsymbol{\mu}\text{ is in the support of }\phi_{\boldsymbol{\mu}}^- \\ \infty, & \text{otherwise.}\end{cases}$$

This result can be interpreted as an "inverse" of Sanov's theorem. Intuitively, given $n$ observations $\boldsymbol{\mu}_n^+(\boldsymbol{\mathcal{E}}_n^Y)$ is "close" to $\boldsymbol{\mu}$ with probability behaving as $e^{-nI_1(\boldsymbol{\mu})}$. Recall, that Sanov's theorem is dealing with deviations of the empirical measure from its typical value and considers the law $\boldsymbol{\mu}$ as given.

We will develop an equivalent of this result in the Markov case. Before we proceed with this agenda, it is instructive to state an equivalent of Sanov's theorem in the Markov case. Consider an irreducible Markov chain with $M$ states $1, 2, \ldots, M$ and the transition probability matrix $\boldsymbol{\Xi}=\{p(i,j)\}_{i,j=1}^M$. Let $\boldsymbol{p}$ denote the vector consisting of the rows of $\boldsymbol{\Xi}$. Let $\boldsymbol{Y}$ denote a sequence $Y_1, Y_2, \ldots, Y_n$ of states that the Markov chain visits with the initial state being $Y_0=\sigma$, and consider the empirical measures

$$\mathcal{E}_{n,2}^Y(\boldsymbol{y})=\frac{1}{n}\sum_{k=1}^n\mathbf{1}_{\boldsymbol{y}}(Y_{k-1}Y_k),$$

where

$$\boldsymbol{y}\in\mathcal{A}^2\triangleq\{1,\ldots,M\}\times\{1,\ldots,M\}.$$

Note that when $\boldsymbol{y}=(i,j)\in\mathcal{A}^2$, the empirical measure $\mathcal{E}_{n,2}^Y(\boldsymbol{y})$ denotes the fraction of times that the Markov chain makes transitions from $i$ to $j$ in the sequence $\boldsymbol{Y}$. Let now

$$\mathcal{A}_{\boldsymbol{p}}^2\triangleq\{(i,j)\in\mathcal{A}^2|p(i,j)>0\}$$

denote the set of pairs of states that can appear in the sequence $Y_1, Y_2, \ldots, Y_n$ and denote by $M_1(\mathcal{A}_{\boldsymbol{p}}^2)$ the standard $|\mathcal{A}_{\boldsymbol{p}}^2|$-di-

mensional probability simplex, where $|\mathcal{A}_{\boldsymbol{p}}^2|$ denotes the cardinality of $\mathcal{A}_{\boldsymbol{p}}^2$. Note that the vector of $\mathcal{E}_{n,2}^{\boldsymbol{Y}}(\boldsymbol{y})$'s denoted by

$$\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = (\mathcal{E}_{n,2}^{\boldsymbol{Y}}(\boldsymbol{y}); y \in \mathcal{A}_{\boldsymbol{p}}^2)$$

is an element of $M_1(\mathcal{A}_{\boldsymbol{p}}^2)$. For any $\boldsymbol{q} \in M_1(\mathcal{A}_{\boldsymbol{p}}^2)$, let

$$q_1(i) \triangleq \sum_{j=1}^{M} q(i, j) \quad \text{and} \quad q_2(i) = \sum_{j=1}^{M} q(j, i) \qquad (9)$$

be its marginals. Whenever $q_1(i) > 0$, let

$$q_f(j|i) \triangleq q(i, j)/q_1(i).$$

We will be using the notation

$$\boldsymbol{q}_f = (q_f(1|1), \ldots, q_f(M|1), q_f(1|2), \ldots, q_f(M|M)).$$

We will say that a probability measure $\boldsymbol{q} \in M_1(\mathcal{A}_{\boldsymbol{p}}^2)$ is *shift-invariant* if both its marginals are identical, i.e., $q_1(i) = q_2(i)$ for all $i$. An LDP for $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}$ is established in the next theorem and is proven in [10, Sec. 3.1.3].

*Theorem IV.1 ([10]):* For every $\boldsymbol{q} \in M_1(\mathcal{A}_{\boldsymbol{p}}^2)$ let

$$I_2(\boldsymbol{q}) = \begin{cases} \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)|p(i, \cdot)), & \text{if } \boldsymbol{q} \text{ is shift invariant} \\ \infty, & \text{otherwise} \end{cases}$$

where $H(q_f(\cdot|i)|p(i, \cdot))$ is the relative entropy defined in (3), that is,

$$H(q_f(\cdot|i)|p(i, \cdot)) = \sum_{j=1}^{M} q_f(j|i) \log \frac{q_f(j|i)}{p(i, j)}.$$

Then, for any set $\Gamma$ of probability vectors in $M_1(\mathcal{A}_{\boldsymbol{p}}^2)$

$$- \inf_{\boldsymbol{q} \in \Gamma^{\circ}} I_2(\boldsymbol{q}) \leq \lim \inf_{n \to \infty} \frac{1}{n} \log \boldsymbol{P} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} \in \Gamma \right]$$

$$\leq \lim \sup_{n \to \infty} \frac{1}{n} \log \boldsymbol{P} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} \in \Gamma \right]$$

$$\leq - \inf_{\boldsymbol{q} \in \Gamma} I_2(\boldsymbol{q})$$

where $\Gamma^{\circ}$ denotes the interior of $\Gamma$.

We next develop the main result of this section, which can be viewed as the "inverse" of Theorem IV.1. In particular, we observe a certain sequence $\boldsymbol{Y}$ of states $Y_1, \ldots, Y_n$ in the Markov chain with the initial state being $Y_0 = \sigma$. Let again $\boldsymbol{\Xi}$ be the transition probability matrix of the Markov chain and $\boldsymbol{p}$ the vector of transition probabilities. Note that $\boldsymbol{p}$ is an element of $(M_1(\mathcal{A}))^M$ (i.e., the $M$-times Cartesian product of $M_1(\mathcal{A})$). We assume that the transition probabilities $\boldsymbol{p}$ are not known; instead, we have a prior $\phi_{\boldsymbol{p}^-} \in (M_1(M_1(\mathcal{A})))^M$, which assigns probability mass only to $\boldsymbol{p}$'s corresponding to irreducible Markov chains. Let $\underline{\boldsymbol{p}}^-$ denote the support of $\phi_{\boldsymbol{p}^-}$. Given the sequence $\boldsymbol{Y}$, the posterior distribution is a function of the empirical measure $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}$; we will write $\phi_{\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})}$ for the posterior distribution and $\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})$ for the corresponding random variable.

Before we proceed with the main result, we establish the following lemma. We will be using the notation $\boldsymbol{P}_{\sigma}^{\boldsymbol{p}}[\cdot]$ to denote probabilities when the initial state is $\sigma$ and the transition probabilities are equal to $\boldsymbol{p}$.

*Lemma IV.2:* Suppose that as $n \to \infty$ the empirical measure $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}$ converges weakly to some shift-invariant $\boldsymbol{q}$ which satisfies $q_1(i) > 0$, $i = 1, \ldots, M$, and $\boldsymbol{q}_f \in \underline{\boldsymbol{p}}^-$. Suppose, also, that $q(i, j) = 0$ implies $\mathcal{E}_{n,2}^{\boldsymbol{Y}}(i, j) = 0$ for all $n$. Then, for any given

$\boldsymbol{q}_n \in M_1(\mathcal{A}_{\boldsymbol{p}}^2)$ that for all $n$ satisfies $q_n(i, j) = 0$ whenever $q(i, j) = 0$, we have

$$\lim_{n \to \infty} \frac{1}{n} \log \int_{(M_1(\mathcal{A}))^M} \boldsymbol{P}_{\sigma}^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$= \sum_{i=1}^{M} \sum_{j=1}^{M} q(i, j) \log q_f(j|i)$$

$$+ \sum_{i=1}^{M} q_1(i) \sum_{j=1}^{M} q_f(j|i) \log q_f(j|i).$$

*Proof:* We have

$$\boldsymbol{P}_{\sigma}^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] = \sum_{\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{y}} = \boldsymbol{q}_n} \boldsymbol{P}_{\sigma}^{\boldsymbol{p}} [Y_1 = y_1, \ldots, Y_n = y_n]$$

$$= \sum_{\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{y}} = \boldsymbol{q}_n} \prod_{i=1}^{M} \prod_{j=1}^{M} (p(i, j))^{n\mathcal{E}_{n,2}^{\boldsymbol{y}}(i, j)}$$

$$= |T_n(\boldsymbol{q}_n)| \prod_{i=1}^{M} \prod_{j=1}^{M} (p(i, j))^{nq_n(i, j)} \qquad (10)$$

where $|T_n(\boldsymbol{q}_n)|$ denotes the size of the type class of $\boldsymbol{q}_n$, i.e., the number of sequences $\boldsymbol{y}$ for which $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{y}} = \boldsymbol{q}_n$. For the first equality above note that given that the Markov chain is in state $i$, transitions "out of $i$" are i.i.d. with probability $p(i, j)$ for a transition to state $j$. Furthermore, by the definition of the empirical measure $n\mathcal{E}_{n,2}^{\boldsymbol{y}}(i, j)$ is equal to the number of transitions from $i$ to $j$ in the sequence $\boldsymbol{y}$. Therefore, letting

$$q_{n1}(i) = \sum_{j=1}^{M} q_n(i, j)$$

[cf. (9)] and $q_{nf}(j|i) = q_n(i, j)/q_{n1}(i)$ we obtain

$$\log \boldsymbol{P}_{\sigma}^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right]$$

$$= \log |T_n(\boldsymbol{q}_n)| + \sum_{i=1}^{M} \sum_{j=1}^{M} nq_n(i, j) \log p(i, j)$$

$$= \log |T_n(\boldsymbol{q}_n)| + \sum_{i=1}^{M} \sum_{j=1}^{M} nq_{n1}(i) q_{nf}(j|i) \log p(i, j)$$

$$= \log |T_n(\boldsymbol{q}_n)| + \sum_{i=1}^{M} nq_{n1}(i)$$

$$\cdot \left( \sum_{j=1}^{M} q_{nf}(j|i) \log q_{nf}(j|i) \right.$$

$$\left. - \sum_{j=1}^{M} q_{nf}(j|i) \log \frac{q_{nf}(j|i)}{p(i, j)} \right)$$

$$= \log |T_n(\boldsymbol{q}_n)| + \sum_{i=1}^{M} nq_{n1}(i)$$

$$\cdot (-H(q_{nf}(\cdot|i)) - H(q_{nf}(\cdot|i)|p(i, \cdot))) \qquad (11)$$

where

$$H(q_{nf}(\cdot|i)) \triangleq - \sum_{j=1}^{M} q_{nf}(j|i) \log q_{nf}(j|i)$$

is the entropy of the probability vector $(q_{nf}(1|i), \ldots, q_{nf}(M|i))$.

It can be shown that (see [10, Exercise 3.1.21])

$$\log |T_n(\boldsymbol{q}_n)| \leq n \sum_{i=1}^{M} \sum_{j=1}^{M} q_n(i,j) \log q_{nf}(j|i). \qquad (12)$$

Using (12), (11), and the fact that the relative entropy is non-negative, we obtain

$$\log \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] \leq n \sum_{i=1}^{M} \sum_{j=1}^{M} q_n(i,j) \log q_{nf}(j|i)$$

$$- \sum_{i=1}^{M} n q_{n1}(i) H(q_{nf}(\cdot|i)).$$

Taking the limit as $n \to \infty$ and since the prior distribution should integrate to one in $(M_1(\mathcal{A}))^M$ we obtain

$$\limsup_{n\to\infty} \frac{1}{n} \log \int_{(M_1(\mathcal{A}))^M} \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\leq \sum_{i=1}^{M} \sum_{j=1}^{M} q(i,j) \log q_f(j|i) - \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)). \quad (13)$$

We are left with proving a matching lower bound. To this end, consider an "$\epsilon$-neighborhood" of $\boldsymbol{q}_f$ defined as follows:

$$\mathcal{N}_\epsilon(\boldsymbol{q}_f) \triangleq \left\{ \boldsymbol{r} \in (M_1(\mathcal{A}))^M \,\middle|\, \frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{M} |q_f(j|i) - r(i,j)| < \epsilon \right\} \qquad (14)$$

and note that when $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n$, for all $\epsilon > 0$, $\boldsymbol{q}_{nf} \in \mathcal{N}_\epsilon(\boldsymbol{q}_f)$ for sufficiently large $n$. We have

$$\liminf_{n\to\infty} \frac{1}{n} \log \int_{(M_1(\mathcal{A}))^M} \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\geq \lim \inf_{n\to\infty} \frac{1}{n} \log \int_{\mathcal{N}_\epsilon(\boldsymbol{q}_f)} \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\geq \lim \inf_{n\to\infty} \frac{1}{n} \log \int_{\mathcal{N}_\epsilon(\boldsymbol{q}_f)} d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$+ \lim \inf_{n\to\infty} \frac{1}{n} \log |T_n(\boldsymbol{q}_n)|$$

$$- \lim \sup_{n\to\infty} \sum_{i=1}^{M} q_{n1}(i) H(q_{nf}(\cdot|i))$$

$$- \lim \sup_{n\to\infty} \sup_{\boldsymbol{p}\in\mathcal{N}_\epsilon(\boldsymbol{q}_f)} \sum_{i=1}^{M} q_{n1}(i) H(q_{nf}(\cdot|i)|p(i,\cdot)) \quad (15)$$

where we have used (11) in the last inequality above. The first term on the right-hand side of the above is zero, since $\boldsymbol{q}_f \in \underline{\boldsymbol{p}}^-$ implies that the integral is strictly positive. For the second term, it can be shown (see [10, eq. 3.1.22]) that

$$\liminf_{n\to\infty} \frac{1}{n} \log |T_n(\boldsymbol{q}_n)| \geq \sum_{i=1}^{M} \sum_{j=1}^{M} q(i,j) \log q_f(j|i). \quad (16)$$

The third term on the right-hand side of (15) equals

$$- \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)).$$

Next, note that the fourth term on the right-hand side of (15) is not infinity since $q(i,j) = 0$ implies $q_n(i,j) = 0$ for all $n$.

Furthermore, it can be made negligible by taking $\epsilon \to 0$. We conclude that

$$\lim \inf_{n\to\infty} \frac{1}{n} \log \int_{(M_1(\mathcal{A}))^M} \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\geq \sum_{i=1}^{M} \sum_{j=1}^{M} q(i,j) \log q_f(j|i) - \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)).$$

This along with (13) establishes the desired result. $\qquad \square$

The main result of this section is the following theorem.

*Theorem IV.3:* Suppose that as $n \to \infty$ the empirical measure $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}$ converges weakly to some shift-invariant $\boldsymbol{q}$ which satisfies $q_1(i) > 0$, $i = 1, \ldots, M$, and $\boldsymbol{q}_f \in \underline{\boldsymbol{p}}^-$. Suppose, also, that $q(i,j) = 0$ implies $\mathcal{E}_{n,2}^{\boldsymbol{Y}}(i,j) = 0$ for all $n$. Then for any prior $\phi_{\boldsymbol{p}^-}$ and for every set $\Gamma$ in $(M_1(\mathcal{A}))^M$ we have

$$- \inf_{\boldsymbol{p}\in\Gamma^\circ} I_3(\boldsymbol{p}) \leq \liminf_{n\to\infty} \frac{1}{n} \log \boldsymbol{P} \left[ \boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}) \in \Gamma \right]$$

$$\leq \limsup_{n\to\infty} \frac{1}{n} \log \boldsymbol{P} \left[ \boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}) \in \Gamma \right]$$

$$\leq - \inf_{\boldsymbol{p}\in\Gamma} I_3(\boldsymbol{p}) \qquad (17)$$

where $\Gamma^\circ$ denotes the interior of $\Gamma$ and the rate function is given by

$$I_3(\boldsymbol{p}) = \begin{cases} \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)|p(i,\cdot)), & \text{if } \boldsymbol{p} \in \underline{\boldsymbol{p}}^- \\ \infty, & \text{otherwise.} \end{cases} \quad (18)$$

*Proof:* Conditional on the empirical measure being equal to $\boldsymbol{q}_n$, and using Bayes' theorem we have

$$\boldsymbol{P} \left[ \boldsymbol{p}_n^+(\boldsymbol{q}_n) \in \Gamma \right]$$

$$= \int_\Gamma \frac{\boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right]}{\int_{(M_1(\mathcal{A}))^M} \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})} d\phi_{\boldsymbol{p}^-}(\boldsymbol{p}). \quad (19)$$

We assume here that the integrand is well-defined, that is, its denominator is not equal to zero.[2] If this is the case, then the integrand is equal to the Radon–Nikodym derivative $\frac{d\phi_{\boldsymbol{p}_n^+(\boldsymbol{q}_n)}}{d\phi_{\boldsymbol{p}^-}}$ evaluated at $\boldsymbol{p}$.

To obtain an upper bound on $\boldsymbol{P}[\boldsymbol{p}_n^+(\boldsymbol{q}_n) \in \Gamma]$ consider first the numerator in the right-hand side of (19). We have

$$\limsup_{n\to\infty} \frac{1}{n} \log \int_\Gamma \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\leq \limsup_{n\to\infty} \frac{1}{n} \log \int_\Gamma d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$+ \limsup_{n\to\infty} \frac{1}{n} \log \sup_{\boldsymbol{p}\in\Gamma\cap\underline{\boldsymbol{p}}^-} \boldsymbol{P}_\sigma^{\boldsymbol{p}} \left[ \boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n \right]$$

$$= \limsup_{n\to\infty} \sup_{\boldsymbol{p}\in\Gamma\cap\underline{\boldsymbol{p}}^-} \left[ \sum_{i=1}^{M} \sum_{j=1}^{M} q_n(i,j) \log q_{nf}(j|i) \right.$$

$$- \sum_{i=1}^{M} q_{n1}(i) \left( H(q_{nf}(\cdot|i)) \right.$$

$$\left. + H(q_{nf}(\cdot|i)|p(i,\cdot)) \right] \qquad (20)$$

[2] For example, it is equal to zero when $n$ is large (i.e., $\boldsymbol{q}_n$ is close to $\boldsymbol{q}$) and $\boldsymbol{q}_f \notin \underline{\boldsymbol{p}}^-$.

where we have used (11) and (12) in the last equality. Here, we are assuming that $\int_\Gamma d\phi_{\boldsymbol{p}^-}(\boldsymbol{p}) > 0$; otherwise, $\boldsymbol{P}[\boldsymbol{p}_n^+(\boldsymbol{q}_n) \in \Gamma] = 0$ and the upper bound still holds. Since as $n \to \infty$, $\boldsymbol{q}_n$ converges to $\boldsymbol{q}$, we combine (19), (20), and the result of Lemma IV.2 to conclude

$$\limsup_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}\left[\boldsymbol{p}_n^+(\boldsymbol{q}_n) \in \Gamma\right]$$

$$\leq - \inf_{\boldsymbol{p} \in \Gamma \cap \underline{\boldsymbol{p}}^-} \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)|p(i,\cdot)) \quad (21)$$

which establishes the upper bound in (17).

Next, note that if $\underline{\boldsymbol{p}}^- \cap \Gamma^\circ$ is empty then by the definition of $I_3(\cdot)$ we have

$$\inf_{\boldsymbol{p} \in \Gamma^\circ} I_3(\boldsymbol{p}) = \inf_{\boldsymbol{p} \in \Gamma^\circ \cap \underline{\boldsymbol{p}}^-} I_3(\boldsymbol{p}) = \infty$$

and the lower bound in (17) holds trivially. Otherwise, for any $\boldsymbol{r} \in \underline{\boldsymbol{p}}^- \cap \Gamma^\circ$ there exists sufficiently small $\epsilon > 0$ such that the "$\epsilon$-neighborhood" of $\boldsymbol{r}$, $\mathcal{N}_\epsilon(\boldsymbol{r})$, is contained in $\underline{\boldsymbol{p}}^- \cap \Gamma^\circ$. By considering the numerator in the right-hand side of (19) we obtain

$$\liminf_{n \to \infty} \frac{1}{n} \log \int_\Gamma \boldsymbol{P}_\sigma^{\boldsymbol{p}}\left[\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n\right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\geq \liminf_{n \to \infty} \frac{1}{n} \log \int_{\mathcal{N}_\epsilon(\boldsymbol{r})} \boldsymbol{P}_\sigma^{\boldsymbol{p}}\left[\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n\right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\geq \liminf_{n \to \infty} \frac{1}{n} \log \int_{\mathcal{N}_\epsilon(\boldsymbol{r})} d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$+ \liminf_{n \to \infty} \frac{1}{n} \log \inf_{\boldsymbol{p} \in \mathcal{N}_\epsilon(\boldsymbol{r})} \boldsymbol{P}_\sigma^{\boldsymbol{p}}\left[\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n\right]$$

$$\geq \lim_{n \to \infty} \inf \frac{1}{n} \log |T_n(\boldsymbol{q}_n)|$$

$$- \lim_{n \to \infty} \sup \sum_{i=1}^{M} q_{n1}(i) H(q_{nf}(\cdot|i))$$

$$- \lim_{n \to \infty} \sup \sup_{\boldsymbol{p} \in \mathcal{N}_\epsilon(\boldsymbol{r})} \sum_{i=1}^{M} q_{n1}(i) H(q_{nf}(\cdot|i)|p(i,\cdot))$$

$$\geq \sum_{i=1}^{M} \sum_{j=1}^{M} q(i,j) \log q_f(j|i) - \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i))$$

$$- \sup_{\boldsymbol{p} \in \mathcal{N}_\epsilon(\boldsymbol{r})} \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)|p(i,\cdot)). \quad (22)$$

In the third inequality above we have used (11). In the last inequality above we have used the fact that $\boldsymbol{q}_n$ converges to $\boldsymbol{q}$ and (16). Taking now $\epsilon \to 0$ we conclude that

$$\liminf_{n \to \infty} \frac{1}{n} \log \int_\Gamma \boldsymbol{P}_\sigma^{\boldsymbol{p}}\left[\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n\right] d\phi_{\boldsymbol{p}^-}(\boldsymbol{p})$$

$$\geq \sum_{i=1}^{M} \sum_{j=1}^{M} q(i,j) \log q_f(j|i) - \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i))$$

$$- \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)|r(i,\cdot)). \quad (23)$$

Using the result of Lemma IV.2 and (19) we obtain

$$\liminf_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}\left[\boldsymbol{p}_n^+(\boldsymbol{q}_n) \in \Gamma\right] \geq - \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)|r(i,\cdot)).$$

$$(24)$$

Recall now that the inequality above holds for all $\boldsymbol{r} \in \underline{\boldsymbol{p}}^- \cap \Gamma^\circ$. Optimizing over such $\boldsymbol{r}$ to obtain the tightest bound we have

$$\liminf_{n \to \infty} \frac{1}{n} \log \boldsymbol{P}\left[\boldsymbol{p}_n^+(\boldsymbol{q}_n) \in \Gamma\right]$$

$$\geq - \inf_{\boldsymbol{r} \in \underline{\boldsymbol{p}}^- \cap \Gamma^\circ} \sum_{i=1}^{M} q_1(i) H(q_f(\cdot|i)|r(i,\cdot))$$

$$= - \inf_{\boldsymbol{r} \in \Gamma^\circ} I_3(\boldsymbol{r})$$

which establishes the lower bound in (17). $\qquad \square$

## V. ESTIMATES OF THE OVERFLOW PROBABILITY

In this section we develop estimates of the overflow probability in a G/G/1 queue and discuss confidence measures.

Consider the model of the G/G/1 queue, introduced in Section III, under a Markov-modulated arrival process. We assume that we have perfect knowledge of the service process $B$, of the number of states $M$ that characterize the arrival process, and of the function $f_A(\cdot)$ that maps Markov states to amount of arriving fluid. The transition probability matrix $\boldsymbol{\Xi}$ of the underlying Markov chain of the arrival process is unknown. Suppose we observe a sequence $\boldsymbol{Y} = Y_1, Y_2, \ldots, Y_n$ of states that this Markov chain visits with the initial state being $Y_0 = \sigma$. With the notation of Section IV, and letting the empirical measure $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}$ be equal to $\boldsymbol{q}_n$, a maximum-likelihood estimator of the transition probabilities is given by

$$\hat{p}_n(i, j) = q_{nf}(j|i) = \frac{q_n(i,j)}{q_{n1}(i)}, \qquad i, j = 1, \ldots, M \quad (25)$$

where we assume that $n$ is large enough to have $q_{n1}(i) > 0$ for all $i$. (To avoid overburdening the notation, in (25) and the sequel we suppress the dependence of the estimators on the sequence $\boldsymbol{Y}$.) Let $\hat{\boldsymbol{p}}_n$ denote the vector of these estimates. We can now construct a matrix $\boldsymbol{\Pi}_{\theta, \hat{\boldsymbol{p}}_n}^A$ with elements [cf. (8)]

$$\pi_{\theta, \hat{\boldsymbol{p}}_n}^A(k, j) = \hat{p}_n(k, j) e^{\theta f_A(j)}, \qquad k, j = 1, \ldots, M$$

and obtain an estimate $\Lambda_A(\theta, \hat{\boldsymbol{p}}_n)$ of the limiting log-moment-generating function for the arrival process by computing the Perron–Frobenius eigenvalue of $\boldsymbol{\Pi}_{\theta, \hat{\boldsymbol{p}}_n}^A$. Applying the result of Proposition III.1, for large values of $U$ we obtain the following estimate of the overflow probability $\boldsymbol{P}[L_i \geq U]$:

$$\mathcal{P}_n^I \triangleq e^{-U\theta^*(\hat{\boldsymbol{p}}_n)} \quad (26)$$

where $\theta^*(\hat{\boldsymbol{p}}_n)$ is the largest root of the equation $\Lambda_A(\theta, \hat{\boldsymbol{p}}_n) + \Lambda_B(-\theta) = 0$. If we knew the *true* transition probabilities, say $\boldsymbol{p}$, then by Proposition III.1 the ersatz $e^{-U\theta^*(\boldsymbol{p})}$ approximates the overflow probability; in (26) we use the estimate $\hat{\boldsymbol{p}}_n$ in the place of the unknown $\boldsymbol{p}$. Notice, however, that as we outlined in Section III, even though $\hat{\boldsymbol{p}}_n$ is an unbiased estimate of the transition probabilities, $\mathcal{P}_n^I$ is not necessarily unbiased. Due to the nonlinearity of $e^{-U\theta^*(\cdot)}$, estimation errors in $\hat{\boldsymbol{p}}_n$ can be severely amplified.

To address this drawback, we next propose an alternative estimator. In particular, taking the Bayesian approach of Section IV, we consider the estimator

$$\boldsymbol{E}\left[e^{-U\theta^*(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}))}\right] \quad (27)$$

where the expectation is taken with respect to the posterior probability distribution of the transition probabilities. Recall that we

have assumed the prior to assign probability mass only to measures corresponding to irreducible Markov chains. As a result, the posterior is also irreducible and $\theta^*(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}))$ is well-defined. Calculating the expectation in (27) exactly is intractable; we will resort to asymptotics. In particular, we let $U = ns$, where $s$ is a scalar, and compute the dominant exponent in the next proposition.

*Proposition V.1:* It holds

$$\lim_{n \to \infty} \frac{1}{n} \log \boldsymbol{E}\left[e^{-ns\theta^*(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}))}\right]$$
$$= -\inf_{\boldsymbol{p} \in (M_1(\mathcal{A}))^M}\left[s\theta^*(\boldsymbol{p}) + I_3(\boldsymbol{p})\right].$$

*Proof:* The result follows by direct application of Varadhan's integral lemma [10, Theorem 4.3.1]. It suffices to verify that the necessary assumptions are satisfied. Indeed, the moment condition

$$\limsup_{n \to \infty} \frac{1}{n} \log \boldsymbol{E}\left[e^{-\gamma ns\theta^*(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}))}\right] < \infty$$

is trivially satisfied for any $\gamma > 1$ since $-\gamma ns\theta^*(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})) \leq 0$, which implies that

$$\boldsymbol{E}\left[e^{-\gamma ns\theta^*(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}))}\right] \leq 1.$$

Moreover, as it is established in the lemma that follows this proof, $\theta^*(\cdot)$ is a continuous function.  □

*Lemma V.2:* Assume that $\boldsymbol{p}$ is the vector of transition probabilities corresponding to an irreducible Markov chain. Then, $\theta^*(\boldsymbol{p})$ is a continuous function of $\boldsymbol{p}$.

*Proof:* It suffices to show that

$$\lim_{\boldsymbol{r} \to \boldsymbol{0}} \theta^*(\boldsymbol{p} + \boldsymbol{r}) = \theta^*(\boldsymbol{p}). \tag{28}$$

Recall that $\Lambda_A(\theta, \boldsymbol{p})$ is the Perron–Frobenius eigenvalue of the matrix $\boldsymbol{\Pi}_{\theta,\boldsymbol{p}}^A$. As such, it is real, nonnegative, and continuous in $\boldsymbol{p}$ and $\theta$ [24]. Therefore, $\Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta)$ is continuous in $\boldsymbol{p}$ and $\theta$ as well, which implies

$$\lim_{\boldsymbol{r} \to \boldsymbol{0}}[\Lambda_A(\theta, \boldsymbol{p} + \boldsymbol{r}) + \Lambda_B(-\theta)] = \Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta)$$
$$\forall \theta > 0.$$

We next distinguish between three cases: $0 < \theta^*(\boldsymbol{p}) < \infty$, $\theta^*(\boldsymbol{p}) = 0$, and $\theta^*(\boldsymbol{p}) = \infty$.

If $0 < \theta^*(\boldsymbol{p}) < \infty$, $\Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta)$ has a single positive root, which by continuity implies that for small enough $||\boldsymbol{r}||$ the perturbed equation

$$\Lambda_A(\theta, \boldsymbol{p} + \boldsymbol{r}) + \Lambda_B(-\theta) = 0$$

has a single root $\theta^*(\boldsymbol{p} + \boldsymbol{r})$ as well. Furthermore,

$$\lim_{\boldsymbol{r} \to \boldsymbol{0}}[\Lambda_A(\theta^*(\boldsymbol{p}), \boldsymbol{p} + \boldsymbol{r}) + \Lambda_B(-\theta^*(\boldsymbol{p}))]$$
$$= \Lambda_A(\theta^*(\boldsymbol{p}), \boldsymbol{p}) + \Lambda_B(-\theta^*(\boldsymbol{p})) = 0$$

and for all $\theta \neq \theta^*(\boldsymbol{p})$

$$\lim_{\boldsymbol{r} \to \boldsymbol{0}}[\Lambda_A(\theta, \boldsymbol{p} + \boldsymbol{r}) + \Lambda_B(-\theta)] = \Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta) \neq 0.$$

Therefore, $\lim_{\boldsymbol{r} \to \boldsymbol{0}} \theta^*(\boldsymbol{p} + \boldsymbol{r}) = \theta^*(\boldsymbol{p})$.

If now $\theta^*(\boldsymbol{p}) = 0$, $\Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta) > 0$ for all $\theta > 0$, which implies

$$\lim_{\boldsymbol{r} \to \boldsymbol{0}}[\Lambda_A(\theta, \boldsymbol{p} + \boldsymbol{r}) + \Lambda_B(-\theta)] = \Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta) > 0$$
$$\forall \theta > 0.$$

Consequently, $\lim_{\boldsymbol{r} \to \boldsymbol{0}} \theta^*(\boldsymbol{p} + \boldsymbol{r}) = 0$.

Finally, if $\theta^*(\boldsymbol{p}) = \infty$, $\Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta) < 0$ for all $\theta > 0$, which implies

$$\lim_{\boldsymbol{r} \to \boldsymbol{0}}[\Lambda_A(\theta, \boldsymbol{p} + \boldsymbol{r}) + \Lambda_B(-\theta)] = \Lambda_A(\theta, \boldsymbol{p}) + \Lambda_B(-\theta) < 0$$
$$\forall \theta > 0.$$

Consequently, $\lim_{\boldsymbol{r} \to \boldsymbol{0}} \theta^*(\boldsymbol{p} + \boldsymbol{r}) = \infty$.

We conclude that (28) holds in all three cases.  □

Based on (27) and the result of Proposition V.1, for large values of $n$ we obtain the following estimator of the overflow probability $\boldsymbol{P}[L_i \geq ns]$:

$$\mathcal{P}_n^{II} \triangleq \exp\left\{-n \inf_{\boldsymbol{p} \in (M_1(\mathcal{A}))^M}[s\theta^*(\boldsymbol{p}) + I_3(\boldsymbol{p})]\right\}. \tag{29}$$

A few remarks on this estimator are in order. Consider a G/G/1 queue serviced by an independent copy of the service process $B$. We draw a sample vector of transition probabilities according to the posterior distribution $\phi_{\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})}$ and feed the queue with an independent copy of the arrival process $A$ with transition probabilities equal to this sample vector. Let $L(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}))$ denote the queue length in this queue. For large enough $n$ we have

$$\boldsymbol{P}\left[L\left(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})\right) > ns\right]$$
$$= \int \boldsymbol{P}\left[L(\boldsymbol{q}) > ns|\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})\right] d\phi_{\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})}(\boldsymbol{q})$$
$$\sim \int e^{-ns\theta^*(\boldsymbol{q})} d\phi_{\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}})}(\boldsymbol{q})$$
$$= \boldsymbol{E}\left[e^{-ns\theta^*(\boldsymbol{p}_n^+(\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}))}\right]$$
$$\sim \exp\left\{-n \inf_{\boldsymbol{p} \in (M_1(\mathcal{A}))^M}[s\theta^*(\boldsymbol{p}) + I_2(\boldsymbol{p})]\right\}. \tag{30}$$

Essentially, the estimator $\mathcal{P}_n^{II}$ captures estimation errors by computing what the loss probability would have been if the arrival process were not exactly known, but characterized by some posterior probability distribution. This interpretation of $\mathcal{P}_n^{II}$ is similar to the result in [15, Theorem 8]; the difference is that [15] does not consider the Markov-modulated case and estimates the limiting log-moment-generating function directly from measurements.

An interesting observation is that $\mathcal{P}_n^{II}$ can be decomposed into two terms. In particular, let $\boldsymbol{p}^*$ be the optimal solution of the optimization problem appearing in the definition of $\mathcal{P}_n^{II}$. We have

$$\mathcal{P}_n^{II} = e^{-ns\theta^*(\boldsymbol{p}^*)}e^{-nI_3(\boldsymbol{p}^*)}. \tag{31}$$

The first term in the right-hand side of the above approximates the loss probability in a G/G/1 queue with buffer size $U = ns$ and arrival process governed by transition probabilities equal to $\boldsymbol{p}^*$. The second term in the right-hand side of the above approximates the probability that the posterior transition probabilities of the arrival process are equal to $\boldsymbol{p}^*$. The vector $\boldsymbol{p}^*$ can be interpreted as the *most likely* value of the transition probabilities of the arrival process, given the observation $\boldsymbol{Y}$, that leads to an overflow. In view of (30), $\boldsymbol{p}^*$ can be thought of as the "worst" transition probability vector in $\boldsymbol{p}^-$, i.e., the one leading to larger overflow probability $\boldsymbol{P}[L(\boldsymbol{p}_n^+(\bar{\boldsymbol{\mathcal{E}}}_{n,2}^{\boldsymbol{Y}})) > ns]$.

To understand on a more intuitive level the estimator $\mathcal{P}_n^{II}$, note that when the empirical measure $\boldsymbol{\mathcal{E}}_{n,2}^{\boldsymbol{Y}}$ converges to some

$\boldsymbol{q}$, $I_3(\boldsymbol{q}_f) = 0$. Furthermore, we have assumed in Theorem IV.3 that $\boldsymbol{q}_f \in \underline{\boldsymbol{p}}^-$. Consequently,

$$\mathcal{P}_n^{II} = \exp\left\{-n \inf_{\boldsymbol{r} \in (M_1(\mathcal{A}))^M} [s\theta^*(\boldsymbol{r}) + I_3(\boldsymbol{r})]\right\} \geq e^{-ns\theta^*(\boldsymbol{q}_f)}. \tag{32}$$

The right-hand side of the above is the large-deviations asymptotic for the overflow probability, if we believe the transition probabilities to be given by $\boldsymbol{q}_f$ (their estimated values). Instead, we quote $\mathcal{P}_n^{II}$ which is larger. We can say that we pay an "estimation premium" to guard against estimation errors. Taking into account the interpretation in (30), $e^{-ns\theta^*(\boldsymbol{q}_f)}$ underestimates the overflow probability, i.e., *certainty equivalence fails*.

If we view the estimator $\mathcal{P}_n^{II}$ as an expectation of the loss measure $e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^{\boldsymbol{Y}}))}$ [cf. (27)], we can also compute higher order moments. More specifically, consider the quantities

$$\boldsymbol{E}\left[e^{-kU\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^{\boldsymbol{Y}}))}\right], \qquad k = 1, 2, \ldots. \tag{33}$$

Using exactly the same reasoning as in the proof of Proposition V.1 we can compute the dominant exponent of these quantities. We can thus define estimators based on higher moments of the loss measure at hand. In particular, for $k = 1, 2, \ldots$ we define

$$\mathcal{P}_n^{III}(k) \triangleq \exp\left\{-n \inf_{\boldsymbol{p} \in (M_1(\mathcal{A}))^M} [ks\theta^*(\boldsymbol{p}) + I_3(\boldsymbol{p})]\right\}. \tag{34}$$

Interestingly enough, we can also asymptotically characterize the distribution of the loss measure $e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^{\boldsymbol{Y}}))}$. The next proposition describes the result.

*Proposition V.3:* We have

$$\lim_{n\to\infty} \frac{1}{n} \log \boldsymbol{P}\left[e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^{\boldsymbol{Y}}))} \geq \epsilon\right] = -\inf_{\substack{\theta^*(\boldsymbol{p}) \leq -\log \epsilon/U \\ \boldsymbol{p} \in (M_1(\mathcal{A}))^M}} I_3(\boldsymbol{p}). \tag{35}$$

*Proof:* The result follows by direct application of the contraction principle [10, Theorem 4.2.1]. In particular, $e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^{\boldsymbol{Y}}))}$ satisfies an LDP in $\mathbb{R}$ with rate function

$$I_4(z) = \inf_{\substack{\theta^*(\boldsymbol{r})=-\log z/U \\ \boldsymbol{r} \in (M_1(\mathcal{A}))^M}} I_3(\boldsymbol{r}).$$

Therefore,

$$\lim_{n\to\infty} \frac{1}{n} \log \boldsymbol{P}\left[e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^{\boldsymbol{Y}}))} \geq \epsilon\right] = -\inf_{z \geq \epsilon} I_4(z)$$

which is equivalent with (35). $\qquad\square$

## VI. COMPUTING THE ESTIMATORS

Computing the estimators in (29) and (34) requires solving nonlinear optimization problems. In this section, we examine the structure of these problems and devise efficient algorithms for their solution.

The optimization problems in (29) and (34) have the following form:

$$\begin{aligned} \text{minimize} \quad & ks\theta^*(\boldsymbol{p}) + I_3(\boldsymbol{p}) \\ \text{s.t.} \quad & p(i, j) \geq 0, \qquad i, j = 1, \ldots, M \\ & \sum_{j=1}^M p(i, j) = 1, \qquad i = 1, \ldots, M \end{aligned} \tag{36}$$

where $k$ is some positive scalar and $I_3(\boldsymbol{p})$ is given by [cf. (18)]

$$I_3(\boldsymbol{p}) = \begin{cases} \sum_{i=1}^M q_1(i) \sum_{j=1}^M q_f(j|i) \log \frac{q_f(j|i)}{p(i,j)}, & \text{if } \boldsymbol{p} \in \underline{\boldsymbol{p}}^- \\ \infty, & \text{otherwise.} \end{cases} \tag{37}$$

We are interested in solving relatively large instances of (36); typically, $M$ would be in the range of 5–100 which brings the number of decision variables in the range of 25–10 000. As a result, computational efficiency is critical and special-purpose algorithms that exploit the special structure of interest.

On the structure of the objective function, a first observation is that $I_3(\boldsymbol{p})$ is a convex function if $\underline{\boldsymbol{p}}^-$ is a convex set. This can be easily established directly from the definition of convexity. Furthermore, $I_3(\boldsymbol{p})$ is strictly convex in $\underline{\boldsymbol{p}}^-$ (where it is finite). In Section IV, we have assumed that the prior assigns probability mass only to $\boldsymbol{p}$'s corresponding to irreducible Markov chains. Let $\mathcal{I}$ be the set of $\boldsymbol{p}$'s corresponding to irreducible Markov chains. The next Proposition shows that $\mathcal{I}$ is a convex set.

*Proposition VI.1:* The set $\mathcal{I}$ of transition probability vectors $\boldsymbol{p}$ corresponding to irreducible Markov chains is a convex set.

*Proof:* Consider two irreducible Markov chains $\mathcal{M}_1$ and $\mathcal{M}_2$ with transition probability vectors $\boldsymbol{p}_1$ and $\boldsymbol{p}_2$, respectively. For some scalar $\lambda \in [0, 1]$ form a new Markov chain $\mathcal{M}_3$ with transition probabilities $\lambda\boldsymbol{p}_1 + (1-\lambda)\boldsymbol{p}_2$. Clearly, for either $\lambda = 0$ or $\lambda = 1$, $\mathcal{M}_3$ is irreducible. We will use contradiction to show that this is also the case for $\lambda \in (0, 1)$. Assume otherwise. This implies that $\mathcal{M}_3$ has two states $i$ and $j$ such that $j$ is not accessible from $i$. Since $\mathcal{M}_1$ is irreducible, there exists a path from $i$ to $j$ with intermediate states $k_1, k_2, \ldots, k_m$, i.e.,

$$p_1(i, k_1), p_1(k_1, k_2), \ldots, p_1(k_m, j) > 0.$$

Consequently, for any $\lambda \in (0, 1)$

$$\begin{aligned} \lambda p_1(i, k_1) + (1-\lambda)p_2(i, k_1), \ldots, \lambda p_1(k_m, j) \\ + (1-\lambda)p_2(k_m, j) > 0. \end{aligned}$$

which contradicts our initial assumption. $\qquad\square$

Based on this proposition, and assuming that the support of the prior $\underline{\boldsymbol{p}}^-$ is a convex subset of $\mathcal{I}$, $I_3(\boldsymbol{p})$ is a convex function. Henceforth, and in the absence of more information on the true transition probabilities of the Markov chain we wish to estimate, we will be making the following assumption.

*Assumption B:* The support of the prior $\underline{\boldsymbol{p}}^-$ is the set $\mathcal{I}$ of transition probability vectors $\boldsymbol{p}$ corresponding to irreducible Markov chains.

Recall next from the statement of Theorem IV.3 that the limit of the empirical measure $\mathcal{E}_{n,2}^{\boldsymbol{Y}}$, $\boldsymbol{q}$, which appears in the definition of $I_3(\boldsymbol{p})$ [cf. (37)], satisfies $\boldsymbol{q}_f \in \underline{\boldsymbol{p}}^-$. As a result, $\boldsymbol{q}_f$ is the transition probability vector of an irreducible Markov chain. Consider now the following optimization problem:

$$\begin{aligned} \text{minimize} \quad & ks\theta^*(\boldsymbol{p}) + \hat{I}_3(\boldsymbol{p}) \\ \text{s.t.} \quad & p(i, j) \geq 0, \qquad i, j = 1, \ldots, M \\ & \sum_{j=1}^M p(i, j) = 1, \qquad i = 1, \ldots, M \end{aligned} \tag{38}$$

where

$$\hat{I}_3(\boldsymbol{p}) = \sum_{i=1}^{M} q_1(i) \sum_{j=1}^{M} q_f(j|i) \log \frac{q_f(j|i)}{p(i,j)}. \qquad (39)$$

The next lemma establishes a property of the optimal solution.

*Lemma VI.2:* Let $\boldsymbol{p}^*$ be an optimal solution of the optimization problem in (38). Then $\boldsymbol{p}^*$ is the transition probability vector of an irreducible Markov chain.

*Proof:* Recall that $\boldsymbol{q}_f$ is the transition probability vector of an irreducible Markov chain, to be denoted by $\mathcal{M}_1$. Then for any pair of states $i$, $j$ of $\mathcal{M}_1$, $j$ is accessible from $i$. More specifically, if $k_1, \ldots, k_m$ are the intermediate states, it holds that

$$q_f(k_1|i), q_f(k_2|k_1), \ldots, q_f(j|k_m) > 0.$$

This implies that

$$p^*(i, k_1), p^*(k_1, k_2), \ldots, p^*(k_m, j) > 0$$

otherwise, $\hat{I}_3(\boldsymbol{p}^*) = \infty$. Thus, the Markov chain with transition probabilities $\boldsymbol{p}^*$ is irreducible. $\qquad \square$

The result of this lemma suggests that under Assumption B the optimization problem in (38) is equivalent to the problem in (36). Hence, we will focus on solving (38).

A very important issue for any nonlinear programming problem is the form of the objective function. A convex objective function, especially under polyhedral constraints (as in (38)), can lead to more efficient algorithms. Let

$$\overline{\theta}(\boldsymbol{p}) \triangleq ks\theta^*(\boldsymbol{p}) + \hat{I}_3(\boldsymbol{p})$$

denote the objective function in (38). Notice that it is the weighted sum of a strictly convex function $\hat{I}_3(\boldsymbol{p})$, and $\theta^*(\boldsymbol{p})$, which is not necessarily convex [see Fig. 1 for a case where $\theta^*(\boldsymbol{p})$ is not convex]. Consequently, $\overline{\theta}(\boldsymbol{p})$ is not convex in general. Nevertheless, since $\hat{I}_3(\boldsymbol{p})$ is strictly convex, it can be seen that $\overline{\theta}(\boldsymbol{p})$ will be convex for small enough values of $s$. Recalling that $s = U/n$ and assuming that the buffer size $U$ is given, we will be dealing with a convex objective function if we can afford a large number $n$ of measurements.

### A. A Heuristic Algorithm

To solve the problem in (38) we have developed a heuristic algorithm which performs very well in practice, typically giving first decimal digit approximations to the optimal value by the third iteration and third decimal digit by the fourth.

To describe the heuristic, note that the constraint set is the Cartesian product of simplices, and that $q_f(j|i) > 0$ implies $p(i, j) > 0$ for all $i, j = 1, \ldots, M$. This feasible set implies the following optimality conditions [25, pp. 178–179]:

$$\frac{\partial \overline{\theta}(\boldsymbol{p})}{\partial p(i,m)} = \frac{\partial \overline{\theta}(\boldsymbol{p})}{\partial p(i,l)} \qquad \forall i, m, l = 1, \ldots, M,$$

$$\text{satisfying } q_f(m|i), q_f(l|i) > 0. \quad (40)$$

The heuristic algorithm iterates as follows.

1) Initialize with $\boldsymbol{p}^{(0)} = \boldsymbol{q}_f$.

2) Let $\boldsymbol{p}^{(m)}$ be the outcome of the $m$th iteration. Stop if $\overline{\theta}(\boldsymbol{p}^{(m-1)}) - \overline{\theta}(\boldsymbol{p}^{(m)}) < \epsilon$, where $\epsilon$ is the desired accuracy.
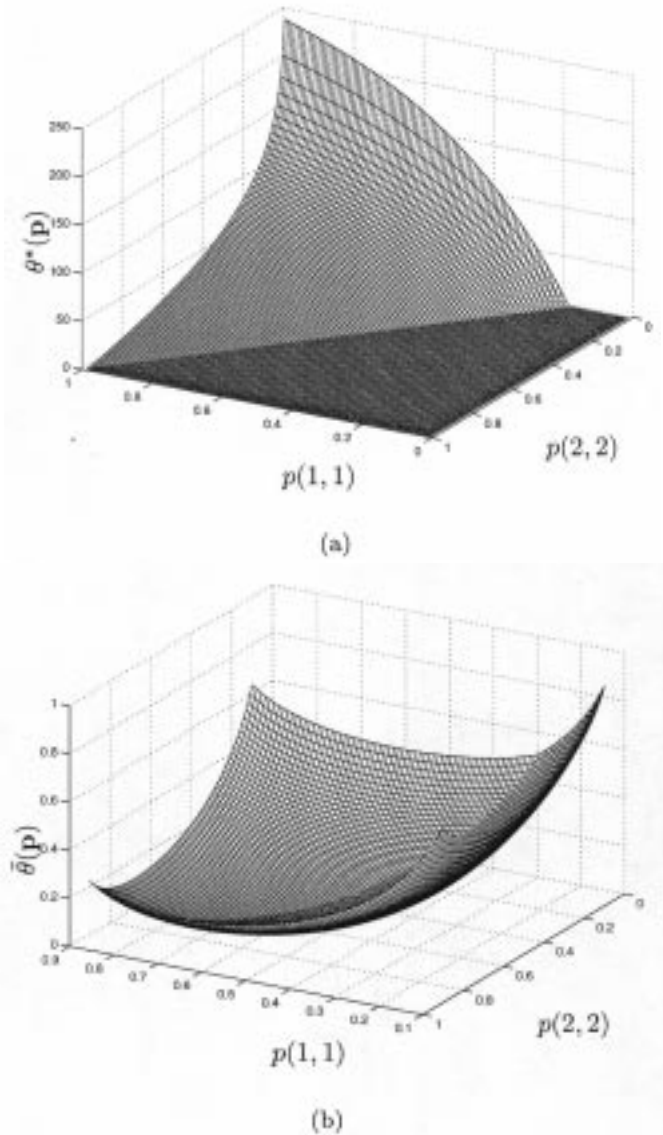


(a)



(b)

Fig. 1. We consider an example where $k = 1$ and the arrival process is the superposition of 36 two-state Markov-modulated processes with $(f_A(1), f_A(2)) = (0.042, 0.077)$. This aggregate process is fed into a buffer with capacity $c = 50/24$. The parameter $s$ was set equal to $0.002$. In (a) we plot $\theta^*(\boldsymbol{p})$ and in (b) $\overline{\theta}(\boldsymbol{p})$ versus the two (arbitrarily chosen as) independent decision variables $p(1, 1)$ and $p(2, 2)$. It can be seen that although $\theta^*(\boldsymbol{p})$ is not convex (it is convex only along some directions), $\overline{\theta}(\boldsymbol{p})$ is convex.

3) Form an approximation of $\overline{\theta}(\boldsymbol{p})$

$$\overline{\theta}_{app}^{(m)}(\boldsymbol{p}) \triangleq ks\theta^* \left( \boldsymbol{p}^{(m)} \right)$$
$$+ ks\nabla\theta^* \left( \boldsymbol{p}^{(m)} \right)' \left( \boldsymbol{p} - \boldsymbol{p}^{(m)} \right) + \hat{I}_3(\boldsymbol{p}). \quad (41)$$

Minimize the expression in (41) subject to the constraints of problem (38) to obtain an optimal solution $\boldsymbol{p}_{min}^{(m)}$. (This minimization can be done by solving the optimality conditions in (40). More specifically, the system of equations in (40) reduces to a scalar nonlinear equation which can be solved numerically.)

4) Set $\boldsymbol{p}^{(m+1)} := \boldsymbol{p}_{min}^{(m)}$ and return to Step 2.

The intuition behind this algorithm is that as we get closer to the real minimum $\boldsymbol{p}^*$ of problem (38), the approximations $\overline{\theta}_{app}^{(m)}(\boldsymbol{p})$

of $\overline{\theta}(\boldsymbol{p})$ improve and yield $\boldsymbol{p}^{(m+1)}$ that are closer to $\boldsymbol{p}^*$. Because at each step we solve the approximate problem exactly, the algorithm typically needs much fewer iterations than a standard gradient-based algorithm.

Although, this heuristic performs well in practice, it does not guarantee convergence. This is the case because the approximations $\overline{\theta}_{\mathrm{app}}^{(m)}(\boldsymbol{p})$ are good in a small region around the expansion point but may be well off away from that point. This might lead to unpleasant situations where $\overline{\theta}(\boldsymbol{p}^{(m)}) > \overline{\theta}(\boldsymbol{p}^{(m-1)})$. An algorithm that does guarantee convergence to a stationary point (i.e., a local minimum) is the *conditional gradient method* (see [25, Sec. 2.2]). This is a *feasible direction* method that iterates as follows:

$$\boldsymbol{p}^{(m+1)} := \boldsymbol{p}^{(m)} + a^{(m)}\left(\overline{\boldsymbol{p}}^{(m)} - \boldsymbol{p}^{(m)}\right) \qquad (42)$$

where $a^{(m)}$ is the step size and $\overline{\boldsymbol{p}}^{(m)}$ is the optimal solution of the subproblem

minimize   $\nabla\overline{\theta}\left(\boldsymbol{p}^{(m)}\right)'(\boldsymbol{p} - \boldsymbol{p}^{(m)})$

s.t.   $p(i, j) \ge 0$,                   $i, j = 1, \ldots, M$

$\sum_{j=1}^{M} p(i, j) = 1$,                   $i = 1, \ldots, M.$   (43)

Notice that this subproblem is a linear programming problem and, thus, can be solved very fast with sophisticated solvers (e.g., such as CPLEX).

To obtain the excellent performance of the heuristic in practice and still be able to guarantee convergence, one can implement a hybrid algorithm that uses the heuristic initially (when far away from the local minimum) and switches to the conditional gradient method when the iterate approaches a local minimum. Both algorithms we examined above, require the gradient of $\overline{\theta}(\boldsymbol{p})$. The gradient of $\hat{I}_3(\boldsymbol{p})$ can be easily obtained. The calculation of the gradient of $\theta^*(\boldsymbol{p})$ is a bit more involved and is given in the Appendix.

Concluding this section, we note that the optimization problem in (35) can be handled similarly. That is, we can dualize the nonlinear constraint and apply the *method of multipliers* [25, Sec. 4.2]. As a result, at each iteration we will be dealing with a problem identical to (38) which makes the discussion above applicable.

## VII. COMPARISONS AND NUMERICAL RESULTS

In this section, we compare the various estimators discussed so far and present some illustrative numerical results.

We consider a queue with a certain Markov-modulated arrival process with transition probability vector $\boldsymbol{p}$. As in Section III, we denote by $L_i$ the queue length and by $U$ the buffer size. The result of Proposition III.1 suggests approximating the overflow probability $\boldsymbol{P}[L_i \ge U]$ with $e^{-U\theta^*(\boldsymbol{p})}$ (see [6], [7], [28] for related approximations and numerical results indicating that this approximation is very accurate for a wide range of $U$'s). We pretend now that we do not know $\boldsymbol{p}$ and we are interested in estimating the loss probability $\boldsymbol{P}[L_i \ge U]$ from measurements. In the sequel, when computing various estimators we will be treating the buffer size $U$ as constant; we will discuss the behavior of the estimators as the number of measurements $n$ increases.

Given that we have observed a sequence $\boldsymbol{Y}$ of $n$ transitions of the Markov chain corresponding to the arrival process, we will consider three estimators: i) $\mathcal{P}_n^I$, ii) $\mathcal{P}_n^{II}$, and iii)

$$\mathcal{P}_n^{IV}(m) \triangleq \mathcal{P}_n^{II} + m\sqrt{\mathcal{P}_n^{III}(2) - (\mathcal{P}_n^{II})^2}$$

where $\mathcal{P}_n^I$, $\mathcal{P}_n^{II}$, and $\mathcal{P}_n^{III}(\cdot)$ were defined in (26), (29), and (34), respectively, and $m$ is some scalar. Recall that in Section V we interpreted $\mathcal{P}_n^{II}$ as the expectation of the loss measure $e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^{\boldsymbol{Y}}))}$, and $\mathcal{P}_n^{III}(2)$ as the corresponding second moment. In this light, $\mathcal{P}_n^{IV}(m)$ can be interpreted as expectation plus $m$ times standard deviation. To compute these estimators we will be using the empirical measure $\mathcal{E}_{n,2}^{\boldsymbol{Y}}$ as an approximation of its limit $\boldsymbol{q}$, which appears in the definition of the rate function $I_3(\cdot)$. Maintaining the notation established so far, we let $\mathcal{E}_{n,2}^{\boldsymbol{Y}} = \boldsymbol{q}_n$, and denote by $\boldsymbol{q}_{nf}$ the corresponding vector of transition probabilities. The next proposition compares the three estimators discussed above.

*Proposition VII.1:* For any nonnegative $n, m$ it holds that

$$\mathcal{P}_n^I \le \mathcal{P}_n^{II} \le \mathcal{P}_n^{IV}(m). \qquad (44)$$

Furthermore, for all $m$ and w.p. 1.

$$\lim_{n\to\infty} \mathcal{P}_n^I = \lim_{n\to\infty} \mathcal{P}_n^{II} = \lim_{n\to\infty} \mathcal{P}_n^{IV}(m) = e^{-U\theta^*(\boldsymbol{p})}. \quad (45)$$

*Proof:* The first inequality in (44) is due to the same argument used in establishing (32). That is,

$$\mathcal{P}_n^{II} = \exp\left\{-n \inf_{\boldsymbol{r}\in(M_1(\mathcal{A}))^M}[s\theta^*(\boldsymbol{r}) + I_3(\boldsymbol{r})]\right\}$$
$$\ge e^{-ns\theta^*(\boldsymbol{q}_{nf})} = e^{-U\theta^*(\boldsymbol{q}_{nf})} = \mathcal{P}_n^I$$

since $I_3(\boldsymbol{q}_{nf}) = 0$. The second inequality in (44) is due to the fact that the standard deviation is nonnegative.

To establish the behavior of the estimators as $n \to \infty$, using renewal reward arguments it can be seen that $\boldsymbol{q}_{nf} \to \boldsymbol{p}$ w.p. 1. As a result

$$\mathcal{P}_n^I = e^{-U\theta^*(\boldsymbol{q}_{nf})} \to e^{-U\theta^*(\boldsymbol{p})}, \qquad \text{w.p. 1.}$$

Next consider $\mathcal{P}_n^{II}$ and note that as $n \to \infty$, and treating $U = ns$ as a constant, we have $s = U/n \to 0$. Observe that $\boldsymbol{q}_{nf}$ is a strict global minimum of $I_3(\cdot)$, and that both $I_3(\cdot)$ and $\theta^*(\cdot)$ are continuous in a neighborhood of $\boldsymbol{q}_{nf}$. Then it can be shown (see [25, Ch. 1, Exercise 1.9]) that the optimal solution of the optimization problem

$$\inf_{\boldsymbol{r}\in(M_1(\mathcal{A}))^M}[s\theta^*(\boldsymbol{r}) + I_3(\boldsymbol{r})]$$

is some function of $n$, say $\boldsymbol{r}^*(n)$, which tends to $\boldsymbol{p}$ as $n \to \infty$. Intuitively, for large enough $n$, $s$ becomes sufficiently small and $I_3(\boldsymbol{r})$ dominates the objective function, which implies that the optimal solution is "close" to the minimizer $\boldsymbol{q}_{nf}$ of $I_3(\boldsymbol{r})$. Namely, for large enough $n$

$$\mathcal{P}_n^{II} = \exp\left\{-n \inf_{\boldsymbol{r}\in(M_1(\mathcal{A}))^M}[s\theta^*(\boldsymbol{r}) + I_3(\boldsymbol{r})]\right\}$$
$$\approx e^{-ns\theta^*(\boldsymbol{q}_{nf})} \to e^{-U\theta^*(\boldsymbol{p})}, \qquad \text{w.p. 1.}$$

The same argument leads to

$$\mathcal{P}_n^{III}(2) = \exp\left\{-n \inf_{\boldsymbol{r} \in (M_1(\mathcal{A}))^M} [2s\theta^*(\boldsymbol{r}) + I_3(\boldsymbol{r})]\right\}$$
$$\to e^{-2U\theta^*(\boldsymbol{p})}, \qquad \text{w.p. } 1$$

which implies

$$m\sqrt{\mathcal{P}_n^{III}(2) - (\mathcal{P}_n^{II})^2} \to 0, \qquad \text{w.p. } 1.$$

We conclude that

$$\mathcal{P}_n^{IV}(m) \to e^{-U\theta^*(\boldsymbol{p})}, \qquad \text{w.p. } 1. \qquad \square$$

Essentially, this proposition establishes that the estimators considered are *consistent* since they converge to the "target" $e^{-U\theta^*(\boldsymbol{p})}$. We can view $e^{-U\theta^*(\boldsymbol{p})}$ as the target because this is what we would have quoted as the overflow probability if we knew the true transition probabilities $\boldsymbol{p}$. When we estimate an overflow probability we are more concerned with underestimating the "true" value than with overestimating it. The reason is that the overflow probability quantifies QoS (see, for example, applications in communication networks [6] and supply chains [7], [28], [9]); controlling the queue to guarantee the desired QoS is important. In this light, the estimators $\mathcal{P}_n^{II}$ and $\mathcal{P}_n^{IV}(m)$ are "safer" than $\mathcal{P}_n^I$ since due to (44) they are less likely to underestimate the overflow probability; yet they are consistent and converge to the correct target $e^{-U\theta^*(\boldsymbol{p})}$ as $n \to \infty$. The parameter $m$ in $\mathcal{P}_n^{IV}(m)$ determines how conservative we choose to be for relatively small values of $n$.

### A. An Example

We next present a numerical example. Consider a Markov-modulated arrival process $A$ driven by a Markov chain with five states and transition probability matrix

$$\Xi = \begin{bmatrix} 0.35 & 0.3 & 0.2 & 0.1 & 0.05 \\ 0.2 & 0.3 & 0.2 & 0.18 & 0.12 \\ 0.14 & 0.22 & 0.3 & 0.21 & 0.13 \\ 0.14 & 0.21 & 0.22 & 0.28 & 0.15 \\ 0.12 & 0.23 & 0.25 & 0.22 & 0.18 \end{bmatrix}.$$

The Markov chain makes one transition per time slot and the number of arrivals per time slot at every state is characterized by

$$[f_A(1), f_A(2), f_A(3), f_A(4), f_A(5)]$$
$$= [0.2, 0.4, 0.6, 0.85, 1.1].$$

We fed this arrival process into a queue with capacity $c = 0.9$ per time slot and buffer size $U = 1.728$. To improve the accuracy of the large-deviations asymptotic we used a refinement discussed in [6] and [7], [28] that introduces a constant in front of the exponential. In particular, the refined large-deviations asymptotic is

$$\boldsymbol{P}[L_i \geq U] \sim \mathcal{P}^{LD} \triangleq \alpha e^{-U\theta^*(\boldsymbol{p})} \qquad (46)$$

where $\alpha = \boldsymbol{E}[L_i]\theta^*(\boldsymbol{p})$ (see [6], [7], [28] for details). In the particular example we are considering, we have $\mathcal{P}^{LD} = 4.9 \times 10^{-7}$. For comparison, simulation yields $6.3 \times 10^{-7}$, which indicates that $\mathcal{P}^{LD}$ is a good approximation (in the sense that it captures

the order of magnitude and approximates well the first significant digit). To compare the different estimators of the overflow probability we generated 10 000 long sample paths of the arrival process. For every one of these sample paths, for all $n = 1, 2, \ldots$, and for some values of $m$, we computed a refined version of $\mathcal{P}_n^I$, $\mathcal{P}_n^{II}$, and $\mathcal{P}_n^{IV}(m)$, which included the constant $\alpha$ in front of exponentials as in (46).[3] Fig. 2 depicts the histogram of the values obtained for $n$ equal to 5000 [Fig. 2(a)–(c)] and 10 000 [Fig. 2(d)–(f)]. We were interested in the fraction of sample paths that lead to a substantial underestimation of $\mathcal{P}^{LD}$. In Fig. 3, for every $n$ we plot the fraction of sample paths that yield an estimator which does not stay above $2.5 \times 10^{-7}$ for all times after $n$. More specifically, let $\mathcal{X}_n$ be either $\mathcal{P}_n^I$, $\mathcal{P}_n^{II}$, or $\mathcal{P}_n^{IV}(m)$, for some values of $m$. For all such $\mathcal{X}_n$ and $n$ we plot the fraction of sample paths yielding $\mathcal{X}_n$ such that there exists $l \geq n$ with $\mathcal{X}_l < 2.5 \times 10^{-7}$. From both Figs. 2 and 3 it can be seen that $\mathcal{P}_n^{II}$ is "safer" (in the sense discussed above) than $\mathcal{P}_n^I$, with the difference being significant for small $n$. Moreover, $\mathcal{P}_n^{IV}(3)$ is "safe" even for very small $n$ (around 2000). To put the values of $n$ into perspective, note that estimating an overflow probability on the order of $10^{-7}$ by directly observing the occupancy of the queue (i.e., measuring the fraction of time that we have an overflow) requires more than $O(10^8)$ observations. Having consistent, yet "safe" estimators with relatively few observations is important in applications since it allows to quickly track changes in the statistics of the input and adapt to level shifts and other nonstationary phenomena. Of course, to avoid underestimation errors the price to pay is that for small $n$ and for a number of sample paths the estimators $\mathcal{P}_n^{II}$ and $\mathcal{P}_n^{IV}(3)$ overestimate the overflow probability and lead to underutilization of the system. Nevertheless, according to Proposition VII.1, as the number of measurements increases they converge to the appropriate target.

Finally, we consider the result of Proposition V.3 to assess the probability of underestimating the "true" overflow probability. For the same data reported above, we generated a sample path of the arrival process and computed the empirical measure $\boldsymbol{q}_n$ based on the first 10 000 observations. Using this empirical measure, we computed $\mathcal{P}_n^I = 3.33 \times 10^{-7}$, $\mathcal{P}_n^{II} = 3.72 \times 10^{-7}$, and $\mathcal{P}_n^{IV}(3) = 9.41 \times 10^{-7}$. In Fig. 4, we plot

$$\log \boldsymbol{P}[e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^Y))} \geq \epsilon]$$

for various values of $\epsilon$. It can be seen that calculations of $\log \boldsymbol{P}[e^{-U\theta^*(\boldsymbol{p}_n^+(\mathcal{E}_{n,2}^Y))} \geq \epsilon]$ can be used to assess the level of confidence in the "safety" of the proposed estimators. Alternatively, this computation can be used to determine the appropriate overflow probability $\epsilon$ to quote for a given confidence level.

### VIII. CONCLUSION

We have proposed estimators of buffer overflow probabilities in buffers fed by Markov-modulated inputs. These estimators are based on analytical large-deviations asymptotics for these probabilities. In particular, we considered a situation where the transition probabilities of the Markov chain characterizing the

---

[3]To avoid overburdening the notation, we will use the same symbols, $\mathcal{P}_n^I$, $\mathcal{P}_n^{II}$, and $\mathcal{P}_n^{IV}(m)$, to denote the refined estimators.
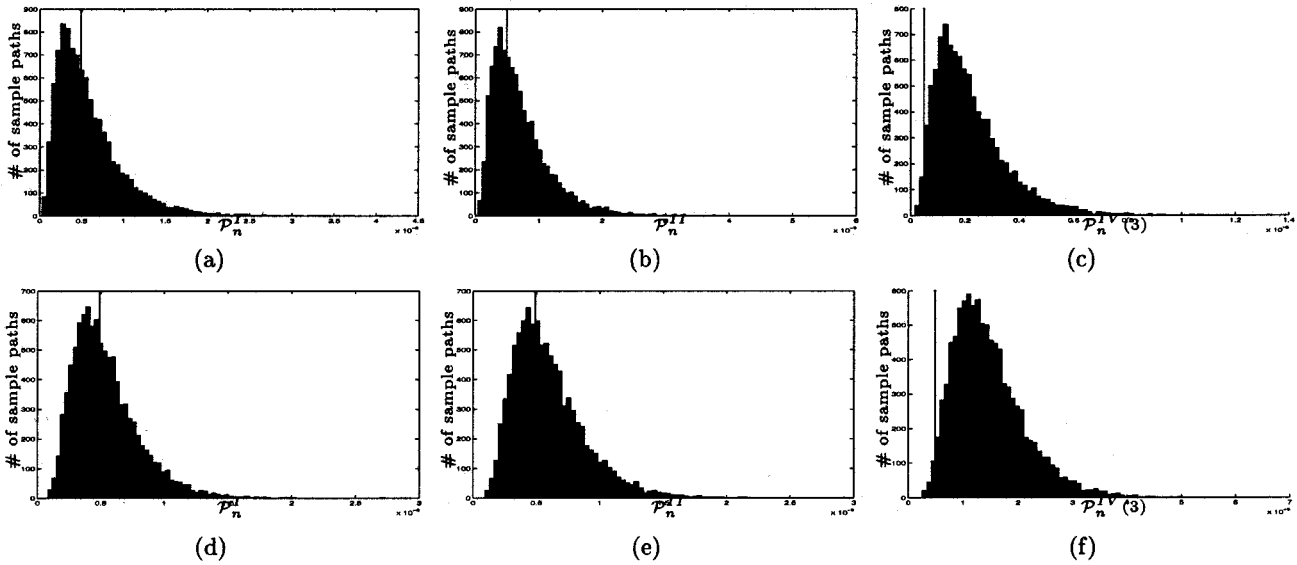
Fig. 2. Plots (a)–(c) depict histograms of the values of $\mathcal{P}_{5000}^I$, $\mathcal{P}_{5000}^{II}$, and $\mathcal{P}_{5000}^{IV}(3)$. Plots (d)–(f) depict histograms of the values of $\mathcal{P}_{10\,000}^I$, $\mathcal{P}_{10\,000}^{II}$, and $\mathcal{P}_{10\,000}^{IV}(3)$. The vertical line drawn at every one of these figures is at $\mathcal{P}^{LD} = 4.9 \times 10^{-7}$.
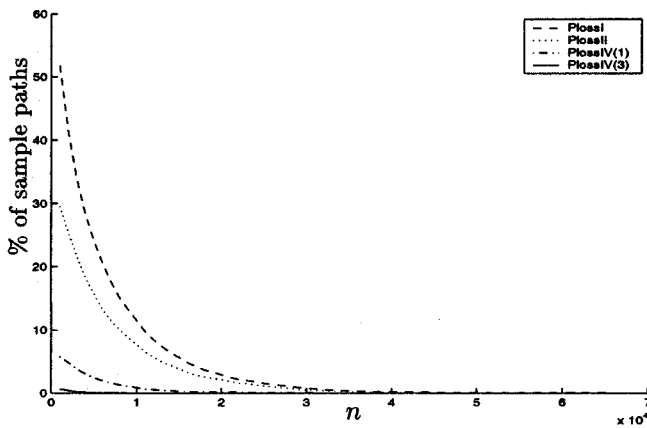


Fig. 3. Comparing the estimators $\mathcal{P}_n^I$, $\mathcal{P}_n^{II}$, and $\mathcal{P}_n^{IV}(m)$, for $m = 1, 3$. In the legend, "PlossI," "PlossII," and "PlossIV(m)" denotes $\mathcal{P}_n^I$, $\mathcal{P}_n^{II}$, and $\mathcal{P}_n^{IV}(m)$, respectively.

arrival process are not known. For any prior distribution of these transition probabilities, we characterized the large-deviations behavior of the posterior distribution. We utilized the latter to define the proposed estimators. We discussed nonlinear programming algorithms for computing these estimators and provided some illustrative numerical results. We demonstrated that the proposed estimators are preferable to the estimator based on certainty equivalence and than observing the occupancy of the queue directly.

## APPENDIX

In this appendix we will calculate the gradient $\nabla \theta^*(\boldsymbol{p})$, in the region where $0 < \theta^*(\boldsymbol{p}) < +\infty$. Consider an $M \times M$ matrix $\boldsymbol{A} = \{a(i, j)\}_{i, j=1, \dots, M}^M$. We will denote by $\lambda_i$ its $i$th eigenvalue, by $\rho(\boldsymbol{A})$ its spectral radius, by $\boldsymbol{x}^i$ its $i$th (normalized) right eigenvector (solution to $\boldsymbol{Ax} = \lambda_i \boldsymbol{x}$), and by $\boldsymbol{y}^i$ its $i$th (normalized) left eigenvector (solution to $\boldsymbol{y}'\boldsymbol{A} = \lambda_i \boldsymbol{y}'$). Assume that $\rho(\boldsymbol{A})$ has multiplicity 1 (as in the case of $\rho(\boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A)$)

and corresponding right (respectively, left) eigenvector denoted by $\boldsymbol{v} = (v_1, \dots, v_M)$ (respectively, $\boldsymbol{u} = (u_1, \dots, u_M)$). Then (see [26, Sec. 2.5–2.8])

$$\frac{\partial \rho(\boldsymbol{A})}{\partial a(i, j)} = \frac{u_i v_j}{\boldsymbol{u}'\boldsymbol{v}}. \tag{47}$$

Recall that

$$\boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A = \left\{ \pi_{\theta, \boldsymbol{p}}^A(i, j) \right\}_{i, j=1, \dots, M}^M$$

$$\triangleq \left\{ p(i, j) e^{\theta f_A(j)} \right\}_{i, j=1, \dots, M}^M$$

and

$$\Lambda_A(\theta, \boldsymbol{p}) = \log \rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right).$$

We obtain

$$\frac{\partial \Lambda_A(\theta, \boldsymbol{p})}{\partial p(i, j)} = \frac{\partial \log \rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)}{\partial p(i, j)}$$

$$= \frac{1}{\rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)} \frac{\partial \rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)}{\partial p(i, j)} = \frac{e^{\theta f_A(j)}}{\rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)} \frac{u_i v_j}{\boldsymbol{u}'\boldsymbol{v}} \tag{48}$$

where $\rho(\boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A)$ is the Perron–Frobenius eigenvalue of $\boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A$ and $\boldsymbol{u}, \boldsymbol{v}$ the corresponding left and right eigenvectors. Similarly

$$\frac{\partial \Lambda_A(\theta, \boldsymbol{p})}{\partial \theta} = \frac{1}{\rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)} \frac{\partial \rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)}{\partial \theta}$$

$$= \frac{1}{\rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)} \sum_{i,j} \left[ \frac{\partial \rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)}{\partial \pi_{\theta, \boldsymbol{p}}^A(i, j)} \frac{\partial \pi_{\theta, \boldsymbol{p}}^A(i, j)}{\partial \theta} \right]$$

$$= \frac{1}{\rho \left( \boldsymbol{\Pi}_{\theta, \boldsymbol{p}}^A \right)} \sum_{i,j} \left[ \frac{u_i v_j}{\boldsymbol{u}'\boldsymbol{v}} f_A(j) p(i, j) e^{\theta f_A(j)} \right]. \tag{49}$$

Fig. 4. We plot $\log \mathbf{P}[e^{-U\theta^*(\mathbf{p}_n^+(\boldsymbol{\varepsilon}_{n,2}^{\mathbf{Y}}))} \geq \epsilon]$. For comparison we draw three vertical lines at $\mathcal{P}_n^I$, $\mathcal{P}_n^{II}$, and $\mathcal{P}_n^{IV}(3)$ (from left to right, respectively).

Recall also that $\theta^*(\mathbf{p})$ is the largest root of the equation $\Lambda_A(\theta, \mathbf{p}) + \Lambda_B(-\theta) = 0$. Based on the discussion in Section III, $\theta^*(\mathbf{p})$ can be alternatively written as

$$\theta^*(\mathbf{p}) = \sup_{\{\theta \mid \Lambda_A(\theta, \mathbf{p}) + \Lambda_B(-\theta) < 0\}} \theta$$

$$= \sup_{\{\theta \mid \Lambda_A(\theta, \mathbf{p}) + \Lambda_B(-\theta) \leq 0\}} \theta \qquad (50)$$

$$= \inf_{a \geq 0} \sup_{\theta} [\theta - a\Lambda_A(\theta, \mathbf{p}) - a\Lambda_B(-\theta)] \qquad (51)$$

where the second equality is due to the convexity of the limiting log-moment-generating functions and the fact that $\Lambda_A(\theta, \mathbf{p}) + \Lambda_B(-\theta)$ is negative for sufficiently small $\theta > 0$. The third equality above is due to strong duality which holds since we are dealing with a convex programming problem (see [25, Ch. 5]). Using the envelope theorem [27], we obtain

$$\nabla \theta^*(\mathbf{p}) = -a^* \nabla_{\mathbf{p}} \Lambda_A(\theta^*, \mathbf{p}) \qquad (52)$$

where $a^*$, $\theta^*$ are optimal solutions of the optimization problems in (51), and the elements of $\nabla_{\mathbf{p}} \Lambda_A(\theta^*, \mathbf{p})$ are given by (48). Note that $a^*$ is a Lagrange multiplier for the original optimization problem in (50) and, therefore, satisfies the first-order optimality condition

$$1 - a^* \frac{\partial(\Lambda_A(\theta^*, \mathbf{p}) + \Lambda_B(-\theta^*))}{\partial \theta^*} = 0$$

which implies

$$a^* = \left[ \frac{\partial(\Lambda_A(\theta^*, \mathbf{p}) + \Lambda_B(-\theta^*))}{\partial \theta^*} \right]^{-1}. \qquad (53)$$

Another way of deriving this result is by using the fact that $\Lambda_A(\theta^*(\mathbf{p}), \mathbf{p}) + \Lambda_B(-\theta^*(\mathbf{p}))$ is equal to zero for all $\mathbf{p}$ satisfying $0 < \theta^*(\mathbf{p}) < +\infty$. Thus,

$$\nabla_{\mathbf{p}} \Lambda_A(\theta^*(\mathbf{p}), \mathbf{p}) + \nabla_{\mathbf{p}} \Lambda_B(-\theta^*(\mathbf{p})) = 0$$

and by using the chain rule we obtain

$$\nabla \theta^*(\mathbf{p}) = -\frac{\nabla_{\mathbf{p}} \Lambda_A(\theta^*, \mathbf{p})}{\frac{\partial}{\partial \theta^*}(\Lambda_A(\theta^*, \mathbf{p}) + \Lambda_B(-\theta^*))} \qquad (54)$$

which is in agreement with (52).

## ACKNOWLEDGMENT

The authors wish to thank A. J. Ganesh for his comments on an earlier conference version of the present paper and for bringing [19] to their attention. They are also very grateful to the anonymous reviewers for their constructive remarks and suggestions.

## REFERENCES

[1] J. Y. Hui, "Resource allocation for broadband networks," *IEEE J. Select. Areas Commun.*, vol. 6, pp. 1598–1608, Sept. 1988.
[2] R. J. Gibbens and P. J. Hunt, "Effective bandwidths for the multi-type UAS channel," *Queueing Syst.*, vol. 9, pp. 17–28, 1991.
[3] F. P. Kelly, "Effective bandwidths at multi-class queues," *Queueing Syst.*, vol. 9, pp. 5–16, 1991.
[4] R. Guérin, H. Ahmadi, and Naghshineh, "Equivalent capacity and its applications to bandwidth allocation in high-speed networks," *IEEE J. Select. Areas Commun.*, vol. 9, pp. 968–981, 1991.
[5] F. P. Kelly, "Notes on effective bandwidths," in *Stochastic Networks: Theory and Applications*, S. Zachary, I. B. Ziedins, and F. P. Kelly, Eds. Oxford, U.K.: Oxford Univ. Press, 1996, vol. 9, pp. 141–168.
[6] I. Ch. Paschalidis, "Class-specific quality of service guarantees in multimedia communication networks," in *Automatica, (Special Issue on Control Methods for Communication Networks)*, V. Anantharam and J. Walrand, Eds., 1999, vol. 35, pp. 1951–1968.
[7] D. Bertsimas and I. Ch. Paschalidis. (1999, Feb.) Probabilistic service level guarantees in make-to-stock manufacturing systems. Dept. Manuf. Eng., Boston Univ. Tech. Rep. [Online]. Available http://ionia.bu.edu
[8] P. Glasserman, "Bounds and asymptotics for planning critical safety stocks," *Operations Res.*, vol. 45, no. 2, pp. 244–257, 1997.
[9] I. Ch. Paschalidis and Y. Liu. (2000, April) Large deviations-based asymptotics for inventory control in supply chains. Dept. Manuf. Eng., Boston Univ. [Online]. Available http://ionia.bu.edu
[10] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd ed. New York: Springer-Verlag, 1998.
[11] D. Bertsimas, I. Ch. Paschalidis, and J. N. Tsitsiklis, "On the large deviations behavior of acyclic networks of G/G/1 queues," *Ann. Appl. Probab.*, vol. 8, no. 4, pp. 1027–1069, 1998.
[12] ——, "Asymptotic buffer overflow probabilities in multiclass multiplexers: An optimal control approach," *IEEE Trans. Automat. Contr.*, vol. 43, pp. 315–335, Mar. 1998.
[13] M. Grossglauser and D. N. C. Tse, "A framework for robust measurement-based admission control," *IEEE/ACM Trans. Networking*, vol. 7, pp. 293–309, June 1999.
[14] ——, "A time-scale decomposition approach to measurement-based admission control," in *Proc. IEEE INFOCOM*, 1999.
[15] N. G. Duffield, "A large deviation analysis of errors in measurement based admission control to buffered and bufferless resources," in *Proc. IEEE INFOCOM*, 1999.
[16] N. G. Duffield, J. T. Lewis, N. O'Connell, R. Russell, and F. Toomey, "Entropy of ATM traffic streams: A tool for estimating QoS parameters," *IEEE J. Select. Areas Commun.*, vol. 13, pp. 981–990, June 1995.
[17] A. Ganesh, P. Green, N. O'Connell, and S. Pitts, "Bayesian network management," *Queueing Syst.*, vol. 28, pp. 267–282, 1998.
[18] A. Ganesh and N. O'Connell, "An inverse of Sanov's theorem," BRIMS, HP Labs, Bristol, U.K., Tech. Rep., 1998.
[19] F. Papangelou, "Large deviations and the Bayesian estimation of higher-order Markov transition functions," *J. Appl. Probab.*, vol. 33, pp. 18–27, 1996.
[20] H. Cramér, "Sûr un nouveau théorème-limite de la théorie des probabilités," in *Actualités Scient. Industrielles (Colloque Consacré à la Théorie des Probabilités)*. Paris, France: Hermann, 1938, pp. 5–23.
[21] R. T. Rockafellar, *Convex Analysis*. Princeton, NJ: Princeton Univ. Press, 1970.
[22] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
[23] P. W. Glynn and W. Whitt, "Logarithmic asymptotics for steady-state tail probabilities in a single-server queue," *J. Appl. Probab.*, vol. 31A, pp. 131–156, 1994.
[24] P. Lancaster, *Theory of Matrices*. New York: Academic, 1969.
[25] D. P. Bertsekas, *Nonlinear Programming*. Belmont, MA: Athena Scientific, 1995.
[26] J. H. Wilkinson, *The Algebraic Eigenvalue Problem*. Oxford, U.K.: Oxford Univ. Press, 1965.
[27] H. R. Varian, *Microeconomic Analysis*, 3rd ed. New York: W. W. Norton, 1992.
[28] D. Bertsimas and I. Ch. Paschalidis, "Probabilistic service level guarantees in make-to-stock manufacturing systems," *Oper. Res.*, to be published.