

Importance Sampling for Jackson Networks*

Paul Dupuis[†] and Hui Wang[‡]
Lefschetz Center for Dynamical Systems
Brown University
Providence, R.I. 02912, U.S.A.

12 March, 2008

Abstract

Rare event simulation in the context of queueing networks has been an active area of research for more than two decades. A commonly used technique to increase the efficiency of Monte Carlo simulation is importance sampling. However, there are few rigorous results on the design of efficient or asymptotically optimal importance sampling schemes for queueing networks. Using a recently developed game/subsolution approach, we construct simple and efficient state-dependent importance sampling schemes for simulating buffer overflows in stable open Jackson networks. The sampling distributions do not depend on the particular event of interest, and hence overflow probabilities for different events can be estimated simultaneously. A by-product of the analysis is the identification of the minimizing trajectory for the calculus of variation problem that is associated with the sample-path large deviation rate function.

1 Introduction

Rare event simulation in the context of queueing networks has been an active area of research for more than two decades. One of the most commonly used techniques to increase the efficiency of Monte Carlo simulation is importance

*We should be generous with references.

[†]Research of this author supported in part by the National Science Foundation (NSF-DMS-0404806 and NSF-DMS-0706003) and the Army Research Office (W911NF-05-1-0289).

[‡]Research of this author supported in part by the National Science Foundation (NSF-DMS-0404806 and NSF-DMS-0706003).

sampling, which amounts to simulating the system under a different probability distribution, and then correcting for any induced bias by multiplying by the likelihood ratio. The literature on importance sampling for queueing networks is large; for a quick introduction and survey see [13, 15].

Despite the considerable effort put into this problem, little is known regarding the design of efficient or asymptotically optimal importance sampling schemes for networks of queues. Most of the literature bases the construction on heuristics, and the resulting importance sampling schemes are supported only by limited numerical evidence. Such an approach is perilous in rare event simulation, since changes of measure suggested by seemingly reasonable heuristics can be very inefficient, sometimes even more so than standard Monte Carlo. For instance, for total population overflow in a simple two-node tandem Jackson network, [19] suggested a change of measure based on large deviation analysis that amounted to exchanging the arrival rate and the smallest service rate. This importance sampling scheme was later on shown to be inefficient in general [11], and in some cases will lead to an estimator with infinite variance [4]. A more subtle and serious danger of heuristic approaches is that a badly designed importance sampling scheme can be mistakenly identified as a good one after moderate (but not extensive) numerical experimentation. In particular, it may yield estimates that are far off from the true values, yet with very small *empirical* variance [12, 8].

The few exceptions in the literature that offer rigorous treatments of state-independent change of measures are only applicable to special classes of queueing networks and/or buffer overflow probabilities. For example, [2, 20] study efficient importance sampling algorithms for simulating large buildups of a single/multiple server queue; [11] shows asymptotic optimality of the change of measure proposed by [19] for total population overflow in tandem Jackson networks, but only under special assumptions on the arrival and service rates; [14] considers the buffer overflow of a single queue within a Jackson network where arrival/service rates and the routing matrix satisfy certain conditions.

It is by now clear that state-independent changes of measure will not, in general, be asymptotically optimal for networks of queues [4], and one must seek among the class of dynamic (i.e., state-dependent) importance sampling schemes in order to attain optimality. The main difficulty in the design problem is the discontinuity of the state dynamics on the boundaries where one or more queues are empty. However, under the game/subsolution framework for importance sampling [8, 9], the construction of an efficient dynamic scheme reduces to the construction of a classical *subsolution* to the

associated Isaacs equation, a nonlinear partial differential equation (PDE) naturally connected to the network via large deviation asymptotics. The subsolution is required to satisfy appropriate boundary conditions that reflect the discontinuities of the state dynamics on the boundaries. The approach has already been successfully applied to tandem Jackson networks [7] and two-node tandem networks with server slowdown [6].

In this paper, we consider the construction of efficient dynamic importance sampling schemes or subsolutions for general buffer overflows in stable open Jackson networks of arbitrary dimension. The key question in the construction is to identify the appropriate gradients for the subsolutions on various boundaries of the state space. Surprisingly, it turns out that all the gradients, and hence the subsolution, can be determined by simply solving $2^d - 1$ systems of *linear* equations, where d is the dimension of the network. This greatly facilitates the implementation of the algorithm. Moreover, the resulting change of measure is *intrinsic* in the sense that it only depends on the system parameters, and is *independent* of the type of buffer overflow problem under consideration (e.g., individual buffer overflow versus total population overflow). Thus the change of measure for a particular network provides a *universal* importance sampling scheme, in that it can be used to estimate the probability of different sets, and in fact probabilities of families of events can be estimated simultaneously using a common set of simulated trajectories. One other feature of note is that while our primary interest is in obtaining an optimal decay rate for the second moment of the estimator, the schemes also have optimal decay rates for all moments of order strictly greater than one. This is discussed in Remark 4.1, and is of particular use in the analysis of the sample variance.

Given any particular event a variety of subsolutions, all of which lead to asymptotically optimal schemes, can be constructed using different approaches and motivations. The existence of universal schemes as described above is related to the fact that the negative of the well-known *quasipotential* function from large deviation theory formally defines a subsolution. We say “formally” because the negative of the quasipotential fails as a *classical sense* subsolution, in that certain boundary conditions hold only in a weak sense and must be suitably mollified to produce a classical sense subsolution. For Jackson networks the quasipotential is simply an affine function (reflecting the product form nature of the invariant distribution), and thus the main work is to identify the boundaries that require mollification and identify appropriate gradients to use along these boundaries. This task is simplified by the explicit identification of the minimizing trajectories in the definition of the quasipotential (i.e., minimizers for the sample path large

deviations rate function that connect the origin to a given point). While our study of the minimizing trajectories heavily exploits the related PDE, optimal trajectories have also been explicitly identified in [17] using a very different time-reversal argument.

While the existence of a universal asymptotically optimal scheme is certainly attractive, it is not a panacea for problems of importance sampling for Jackson networks. For example, for the estimation of any particular event the practical performance of a scheme designed specifically for that event may well be superior, since the particular parts of the boundary that require mollification will depend on the configuration of the escape region. The universal scheme will likely be more complex than is necessary for this single event, since it must account for all escape regions. On the other hand, we now have useful and provably stable schemes where there were none before, and the constructive methods developed in this paper may also find use in more specialized analyses.

Although in this paper we restrict to Jackson networks, we expect that the techniques are applicable to analogous models for stochastic networks that possess a product form stationary distribution, such as skew-symmetric and related forms of reflecting Brownian motion and Kelly networks. One may go further, and observe that since it is only the form of the large deviation asymptotics that determine the form of the associated Isaacs equation, it may be possible to prove analogous results for networks that are only “asymptotically product form” in an appropriate sense.

The paper is organized as follows. Section 2 introduces the setup. The buffer overflow probabilities of interest and their large deviations asymptotics are given in Section 3. Section 4 gives a brief discussion on the asymptotic optimality of importance sampling estimators. The mathematical model for system dynamics is established in Section 5 and the dynamic importance sampling scheme is defined in Section 6. The various Hamiltonians and their roots, which are essential to the construction of the subsolution, are analyzed in Section 7. In Section 8 we identify the related Isaacs equation. Subsolutions and the corresponding importance sampling schemes are constructed in Sections 9 and 10. In Section 11 we explicitly solve the calculus of variation problem associated with the sample path large deviation rate function. Numerical results are presented in Section 12. We note that many technical proofs are postponed to the appendix in order to ease exposition.

2 Model setup

Consider a classical d -node open Jackson network with arrival rates $\lambda = (\lambda_1, \dots, \lambda_d)'$ and service rates $\mu = (\mu_1, \dots, \mu_d)'$. Departures from node i join node j with probability P_{ij} and leave the system with probability

$$P_{i0} \doteq 1 - \sum_{j=1}^d P_{ij}. \quad (2.1)$$

Define $P = [P_{ij}]_{1 \leq i, j \leq d}$, which is a substochastic matrix. The state process is denoted by $\{Q(t) = (Q_1(t), \dots, Q_d(t))' : t \geq 0\}$, where $Q_i(t)$ is the queue length at node i and at time t . The process Q is *constrained* since the queue lengths have to be non-negative.

Assumptions throughout the paper.

- [a]. For each i , either $\lambda_i > 0$ or $\lambda_{j_1} P_{j_1 j_2} \cdots P_{j_k i} > 0$ for some j_1, \dots, j_k .
- [b]. For each i , either $P_{i0} > 0$ or $P_{i j_1} P_{j_1 j_2} \cdots P_{j_k 0} > 0$ for some j_1, \dots, j_k .
- [c]. The network is stable.

Assumptions [a,b,c] amount to that each node in the network will receive external input (possibly through other nodes), and each job will eventually leave the system. Under these assumptions, $P = [P_{ij}]_{1 \leq i, j \leq d}$ has spectral radius less than one [3], and the *utilization parameters* $\rho = (\rho_1, \dots, \rho_d)'$ given by

$$\rho_i = \frac{[\lambda'(I - P)^{-1}]_i}{\mu_i} \quad (2.2)$$

satisfy $\rho_i \in (0, 1)$ for all $i = 1, \dots, d$.

3 Buffer overflow probabilities

Consider a set $\Gamma \subset \mathbb{R}_+^d$ such that $0 \notin \bar{\Gamma}$, where $\bar{\Gamma}$ denotes the closure of Γ . The probability of interest is

$$p_n \doteq \mathbb{P}\{\text{the state process } Q \text{ reaches } n\Gamma \text{ before coming} \\ \text{back to } 0, \text{ starting from } 0\} \quad (3.1)$$

for large n . Under the stability assumption on the network, p_n is the probability of a rare event. The goal of the current paper is to design a simple and asymptotically optimal importance sampling algorithm for the Monte Carlo estimation of p_n .

The following result, which characterizes the exponential decay rate of $\{p_n\}$, will be useful. Its proof is analogous to that of [11, Theorem 2.3]. For completeness, we include the proof in the appendix. Let Γ° denote the interior of Γ .

Proposition 3.1. *We have*

$$-\inf_{x \in \Gamma^\circ} \langle -r^*, x \rangle \leq \liminf_n \frac{1}{n} \log p_n \leq \limsup_n \frac{1}{n} \log p_n \leq -\inf_{x \in \Gamma} \langle -r^*, x \rangle,$$

where

$$r^* \doteq (\log \rho_1, \dots, \log \rho_d)'. \quad (3.2)$$

Assumptions throughout the paper.

[d].

$$\gamma \doteq \inf_{x \in \Gamma^\circ} \langle -r^*, x \rangle = \inf_{x \in \Gamma} \langle -r^*, x \rangle. \quad (3.3)$$

It follows from Proposition 3.1 that γ is the exponential decay rate of p_n , or

$$-\lim_n \frac{1}{n} \log p_n = \gamma.$$

4 Basics of importance sampling

In this section, we give a brief discussion on the asymptotic optimality of importance sampling algorithms. Consider a family of events $\{A_n\}$ with

$$\lim_n -\frac{1}{n} \log \mathbb{P}(A_n) = \gamma.$$

In order to estimate $\mathbb{P}(A_n)$, importance sampling generates samples under a different probability distribution \mathbb{Q}_n (i.e., change of measure) such that $\mathbb{P} \ll \mathbb{Q}_n$, and forms an estimator by averaging independent replications of

$$\hat{p}_n \doteq 1_{A_n} \frac{d\mathbb{P}}{d\mathbb{Q}_n},$$

where $d\mathbb{P}/d\mathbb{Q}_n$ is the Radon-Nikodým derivative or likelihood ratio. It is easy to check that \hat{p}_n is unbiased.

The rate of convergence of the importance sampling estimator is determined by the variance of \hat{p}_n . Since the estimate is unbiased, minimizing the

variance is equivalent to minimizing the second moment, and the smaller the second moment, the faster the convergence. However, by Jensen's inequality

$$\limsup_n -\frac{1}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n^2] \leq \limsup_n -\frac{2}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n] = 2\gamma.$$

We say the importance sampling estimator \hat{p}_n or the change of measure \mathbb{Q}_n is *asymptotically optimal* if the upper bound is achieved, i.e., if

$$\liminf_n -\frac{1}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n^2] \geq 2\gamma.$$

Sometimes 2γ is referred to simply as the “optimal decay rate.”

Remark 4.1. It is sometimes of interest to study other moments of \hat{p}_n . One reason is that when an estimate is obtained, its standard error will be estimated from the empirical data. It would be nice to know if the empirical standard error is close to the theoretical standard error. This clearly connects with the higher order moments of \hat{p}_n . In general, when the b -th moment of \hat{p}_n is of interest with $b > 1$, we have the analogous inequality

$$\limsup_n -\frac{1}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n^b] \leq \limsup_n -\frac{b}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n] = b\gamma,$$

and say \hat{p}_n is *asymptotically optimal in the b -th moment* if

$$\liminf_n -\frac{1}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n^b] \geq b\gamma.$$

The classical definition of asymptotical optimality is just the special case with $b = 2$. As we will see, the importance sampling estimators we construct will be asymptotically optimal in the b -th moment for all $b > 1$.

5 The system dynamics

To ease exposition, we consider the embedded discrete time Markov chain $Z = \{Z(k) : k \geq 0\}$, defined on some probability space $(\Omega, \mathcal{F}, \mathbb{P})$, that represents the queue lengths at the transition epochs of the Jackson network, and assume without loss of generality

$$\sum_{i=1}^d [\lambda_i + \mu_i] = 1. \tag{5.1}$$

We introduce the following notation. Denote by e_i the unit vector with the i -th component one and zero elsewhere, and let

$$\mathbb{V} \doteq \{e_i, -e_i + e_j, -e_i : i, j = 1, 2, \dots, d\}.$$

Note that e_i represents an arrival at node i , $-e_i + e_j$ a departure from node i that joins node j , and $-e_i$ a departure from node i that leaves the system. As a convention that will help with notation later on, we distinguish between $-e_i + e_i$ (a service at station i that immediately returns to i), $-e_j + e_j$, and 0, which by our convention is not an element of \mathbb{V} .

Since the state process Z is constrained in the non-negative orthant by nature, its dynamics are discontinuous on the boundaries where one or more queues are empty. These boundaries are indexed by subsets $E \subset \{1, 2, \dots, d\}$, where

$$\partial_E \doteq \{x \in \mathbb{R}_+^d : x_i = 0 \text{ for } i \in E \text{ and } x_i > 0 \text{ for } i \notin E\},$$

with the convention

$$\partial_\emptyset \doteq \{x = (x_1, \dots, x_d) \in \mathbb{R}_+^d : x_i > 0 \text{ for all } i\}$$

for the interior of the state space.

For each $i = 1, \dots, d$, define the subset of \mathbb{V}

$$\mathbb{V}[-e_i] \doteq \{v : v = -e_i \text{ or } v = -e_i + e_j \text{ for some } j = 1, \dots, d\}. \quad (5.2)$$

The non-negativity constraint mandates that on ∂_E the state process Z cannot make physical jumps of size v if v belongs to the set $\cup_{i \in E} \mathbb{V}[-e_i]$. This motivates the definition of the constraining mapping $\pi : \mathbb{R}_+^d \times \mathbb{V} \mapsto \mathbb{V} \cup \{0\}$ by

$$\pi[x, v] \doteq \begin{cases} 0 & \text{if } x \in \partial_E, v \in \cup_{i \in E} \mathbb{V}[-e_i], \\ v & \text{otherwise.} \end{cases} \quad (5.3)$$

The evolution of the Markov chain Z can then be modeled by the equation

$$Z(k+1) = Z(k) + \pi[Z(k), Y(k+1)], \quad (5.4)$$

where $Y = \{Y(k) : k \geq 1\}$ is a sequence of random variables taking values in \mathbb{V} . In other words, on ∂_E , the process Z is still allowed to make “fictitious” jumps of sizes that are in $\cup_{i \in E} \mathbb{V}[-e_i]$, but it will be pushed back so that it stays in the positive orthant.

Under the original measure \mathbb{P} , $Y = \{Y(k) : k \geq 1\}$ are independent and identically distributed (iid) with common distribution

$$\Theta[v] = \begin{cases} \lambda_i & \text{if } v = e_i, \\ \mu_i P_{ij} & \text{if } v = -e_i + e_j, \\ \mu_i P_{i0} & \text{if } v = -e_i. \end{cases} \quad (5.5)$$

In importance sampling algorithms, the distribution of Y will be altered, i.e., change of measure. We should restrict our attention to those probability distributions that are equivalent to Θ , namely,

$$\mathcal{P}^+(\mathbb{V}) \doteq \left\{ \theta : 0 \leq \theta[v] \leq 1, \sum_{v \in \mathbb{V}} \theta[v] = 1, \theta[v] = 0 \text{ iff } \Theta[v] = 0 \right\}. \quad (5.6)$$

Remark 5.1. Note that since with our convention $-e_i + e_i \neq -e_j + e_j$, $\mathbb{V}[-e_i] \cap \mathbb{V}[-e_j] = \emptyset$ whenever $i \neq j$.

6 Dynamic importance sampling schemes

Given the large deviation parameter n , it is convenient to study the scaled state process $X^n = \{X^n(k) \doteq Z(k)/n : k \geq 0\}$. Note that, by (5.3) and (5.4), X^n satisfies

$$X^n(k+1) = X^n(k) + \frac{1}{n} \pi[X^n(k), Y(k+1)]. \quad (6.1)$$

With this notation, the probability of interest p_n equals the probability of the scaled process X^n reaching Γ before coming back to the origin, after starting from the origin [i.e., $X^n(0) = 0$].

Dynamic, or state-dependent, importance sampling schemes can be characterized by alternative stochastic kernels $\bar{\Theta}^n[\cdot|\cdot]$ on \mathbb{V} given \mathbb{R}_+^d that satisfy $\bar{\Theta}^n[\cdot|x] \in \mathcal{P}^+(\mathbb{V})$ for every $x \in \mathbb{R}_+^d$. More precisely, the conditional probability of $Y(k+1) = v$ given the history $\{Y(j) : j = 1, \dots, k\}$, is just $\bar{\Theta}^n[v|X^n(k)]$ for every $v \in \mathbb{V}$ and every k . Define two hitting times

$$T_n \doteq \inf\{k \geq 1 : X^n(k) \in \Gamma\}, \quad T_0 \doteq \inf\{k \geq 1 : X^n(k) = 0\}. \quad (6.2)$$

Then the corresponding unbiased importance sampling estimator is

$$\hat{p}_n \doteq 1_{\{T_n < T_0\}} \prod_{k=0}^{T_n-1} \frac{\Theta[Y(k+1)]}{\bar{\Theta}^n[Y(k+1)|X^n(k)]}. \quad (6.3)$$

7 The Hamiltonians and their roots

The construction of efficient importance sampling schemes is intimately connected with differential games, Isaacs equations and subsolutions [9]. We start by considering various Hamiltonians for the state dynamics and important roots that are associated with these Hamiltonians. Such knowledge is essential for the construction that follows.

7.1 The various Hamiltonians

Given any $x \in \mathbb{R}_+^d$ and $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{R}^d$, the Hamiltonian at x is defined by

$$H(x, \alpha) \doteq \log \left(\sum_{v \in \mathbb{V}} \Theta[v] \cdot e^{\langle \alpha, \pi[x, v] \rangle} \right). \quad (7.1)$$

Intuitively speaking, $H(x, \cdot)$ is just the log moment generating function of the increments of the state process at state x . It is worth noting that for $x \in \partial_E$ with $E \subset \{1, \dots, d\}$, $H(x, \alpha) = H_E(\alpha)$, where

$$H_E(\alpha) \doteq \log \left(\sum_{v \notin \cup_{i \in E} \mathbb{V}[-e_i]} \Theta[v] \cdot e^{\langle \alpha, v \rangle} + \sum_{v \in \cup_{i \in E} \mathbb{V}[-e_i]} \Theta[v] \right). \quad (7.2)$$

In the literature, H_\emptyset is often referred to as the *interior* Hamiltonian and H_E as a *boundary* Hamiltonian when E is non-empty. A standard calculation using Hölder's inequality shows that $H(x, \cdot)$ and $H_E(\cdot)$ are convex.

Lemma 7.1. *For all subsets $E, E_1, E_2 \subset \{1, \dots, d\}$ such that $E_1 \cap E_2 = \emptyset$ and $E_1 \cup E_2 = E$,*

$$e^{H_\emptyset} + e^{H_E} = e^{H_{E_1}} + e^{H_{E_2}}.$$

The proof is straightforward from the definition (7.2) and Remark 5.1, and therefore omitted.

Remark 7.2. To ease exposition, we will simply use H_i to denote H_E when $E = \{i\}$.

7.2 The roots of Hamiltonians

As we will see, subsolutions for Jackson networks can be constructed from a small number of key roots to the various Hamiltonians. These roots only depend on the arrival/service rates and the matrix P . They are *independent*

of the set Γ . A detailed discussion on the motivation and intuition behind these roots is presented in Section 9.3.

Later in this section we will need the notation

$$\mathbb{V}[e_i] \doteq \{v \in \mathbb{V} : v = e_i, \text{ or } v = -e_j + e_i \text{ for some } j = 1, \dots, d\}.$$

This should be distinguished from the previous notation $\mathbb{V}[-e_i]$ as defined by (5.2). We will also need two well known identities—see, e.g., [3, Chapter IV.2]. The first identity says that the arrival rate equals the output rate at node j in equilibrium, while the second identity amounts to that the total arrival rate to the network equals the total output rate of the network.

Lemma 7.3. *For each $j = 1, \dots, d$, we have*

$$\lambda_j + \sum_{i=1}^d \mu_i \rho_i P_{ij} = \mu_j \rho_j. \quad (7.3)$$

Furthermore,

$$\sum_{j=1}^d \lambda_j = \sum_{j=1}^d \mu_j \rho_j P_{j0}. \quad (7.4)$$

Recall the definition of r^* as in (3.2). For any non-empty subset $E \subset \{1, 2, \dots, d\}$, consider the system of linear equations of $\{z_i : i \in E\}$

$$\sum_{\substack{v \in \mathbb{V}[-e_i]: \\ v \neq -e_i + e_j \\ \text{any } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i]: \\ v = -e_i + e_j \\ \text{some } j \in \bar{E}}} \Theta[v] e^{-\langle r^*, v \rangle} z_j = \mu_i z_i. \quad (7.5)$$

Later on we will define a vector r^E in terms of the solution to this system. The form of (7.5) is dictated by the requirement $H_i(-r^E) = 0$ for $i \in E$. See also the discussion at the end of Section 9.3.

Lemma 7.4. *For every non-empty subset $E \subset \{1, \dots, d\}$, the system of linear equations (7.5) has a unique solution $z^E = \{z_i^E : i \in E\}$ with strictly positive components.*

Proof. For $i \in E$, let $u_i \doteq \mu_i z_i$. Then the system of equations (7.5) can be rewritten as

$$c + Bu = u, \quad (7.6)$$

where

$$u = (u_i)_{i \in E}', \quad c = (c_i)_{i \in E}', \quad B = [B_{ij}]_{i, j \in E},$$

$$c_i \doteq \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle}, \quad B_{ij} \doteq \Theta[v] e^{-\langle r^*, v \rangle} / \mu_j \Big|_{v = -e_i + e_j}.$$

We first argue that B' is a substochastic matrix. Recalling that $r_i^* = \log \rho_i$ and the definition of $\Theta[v]$, we observe that the identity (7.3) is equivalent to the identity

$$-\mu_i + \sum_{v \in \mathbb{V}[e_i]} \Theta[v] e^{-\langle r^*, v \rangle} = 0, \quad (7.7)$$

and therefore for every $j \in E$

$$\sum_{i \in E} B_{ij} = \frac{1}{\mu_j} \sum_{\substack{v = -e_i + e_j \\ \text{some } i \in E}} \Theta[v] e^{-\langle r^*, v \rangle} \leq \frac{1}{\mu_j} \sum_{v \in \mathbb{V}[e_j]} \Theta[v] e^{-\langle r^*, v \rangle} = 1. \quad (7.8)$$

Now we argue that B' has spectral radius less than one. Define

$$E_1 \doteq \left\{ j \in E : \sum_{i \in E} B_{ij} = 1 \right\}, \quad E_2 \doteq E \setminus E_1.$$

For every $j \in E_1$, inequality (7.8) is indeed equality and therefore $\Theta[e_j] = \lambda_j = 0$ and $\Theta[-e_i + e_j] = \mu_i P_{ij} = 0$ or $P_{ij} = 0$ for all $i \notin E$. Now fix an arbitrary $j \in E_1$. By Assumption [a] and since $\lambda_j = 0$, there must exist some j_1, \dots, j_k such that $\lambda_{j_1} P_{j_1 j_2} \cdots P_{j_k j} > 0$. Clearly $j_1 \notin E_1$, and therefore we can define a $k^* \in \{1, \dots, k\}$ such that $\{j_{k^*+1}, \dots, j_k, j\} \subset E_1$, while $j_{k^*} \notin E_1$. Since $P_{ij_{k^*+1}} = 0$ [with $j_{k^*+1} \doteq j$ if $k^* = k$] for all $i \notin E$, we have $j_{k^*} \in E$. It follows that $j_{k^*} \in E_2$. Recalling the definition of B_{ij} ,

$$[B']_{jj_k} [B']_{j_k j_{k-1}} \cdots [B']_{j_{k^*+1} j_{k^*}} = B_{j_{k^*} j_{k^*+1}} \cdots B_{j_{k-1} j_k} B_{j_k j} > 0.$$

Therefore the spectral radius of B' is less than one [3, Proposition I.6.3].

It follows that $(I - B)$ is invertible and $(I - B)^{-1} = I + B + B^2 + \cdots$ has non-negative components. Therefore, $u^E \doteq (I - B)^{-1} c$ is the unique solution to the system of equations (7.6). Clearly $u^E \geq 0$. We now argue $u^E > 0$ by contradiction. Suppose this is not true, or the set $F \doteq \{i \in E : u_i^E = 0\}$ is non-empty. Then for every $i \in F$, equation (7.6) implies that $c_i = 0$ and $B_{ij} = 0$ for all those $j \in E \setminus F$. But $c_i = 0$ amounts to $P_{ij} = 0$ for all $j \notin E$ and $j = 0$, and $B_{ij} = 0$ amounts to $P_{ij} = 0$. Therefore $P_{ij} = 0$ for all $i \in F$ and $j \notin F$ or $j = 0$. Thus for $i \in F$ the Assumption [b] is violated, a contradiction. Letting $z_i^E \doteq u_i^E / \mu_i$ we complete the proof. \blacksquare

Corollary 7.5. *Suppose that a vector $\{z_i : i \in E\}$ satisfies (7.5) with “=” replaced by “ \geq ” [resp., “ \leq ”]. Then $z_i^E \geq z_i$ [resp., “ \leq ”] for every $i \in E$.*

Proof. We only show for the case “ \geq ”, the other direction being analogous. Let $u^E \doteq (\mu_i z_i^E)_{i \in E}$ and $u \doteq (\mu_i z_i)_{i \in E}$. Using the notation in the proof of Lemma 7.4, we have $c + Bu^E = u^E$ and $c + Bu \geq u$. Letting $w = u^E - u$, we have $w = Bw + (c + Bu - u)$. It follows readily that $w = (I - B)^{-1}(c + Bu - u) \geq 0$ since $(I - B)^{-1}$ has non-negative components and $c + Bu - u \geq 0$. ■

For every non-empty subset E , define a vector r^E by

$$r^E \doteq r^* - \sum_{i \in E} \log[z_i^E] \cdot e_i. \quad (7.9)$$

For completeness, we also define $r^\emptyset = r^*$.

Remark 7.6. Observe that for any i

$$\sum_{v \in \mathbb{V}[-e_i]} \Theta[v] = \sum_{j=0}^d \mu_i P_{ij} = \mu_i.$$

This implies $z_i^{\{1, \dots, d\}} = \rho_i$ for every i , which yields $r^{\{1, \dots, d\}} = 0$.

Proposition 7.7. *For every subset $E \subset \{1, 2, \dots, d\}$, the vector r^E satisfies*

$$H_F(-r^E) = 0$$

for all $F \subset E$.

This is in fact the key result for the construction of good subsolutions. When constructing a subsolution W , we will want the gradient to satisfy $H_E(DW(x)) \geq 0$ when $x \in \partial_E$. Suppose that W is smooth, and let DW denote the gradient. Proposition 7.7 tells us that if $DW(x) = -r^E$ at a point $x \in \partial_E$, then at least up to a small error $H_F(DW(y)) \geq 0$ for all y in a neighborhood of x , where F is the set of indices such that $y \in \partial_F$.

Not every root r^E is needed in the construction of subsolutions. Define

$$\mathcal{F} \doteq \{F \subset \{1, \dots, d\} : z_i^F < 1 \text{ for all } i \in F\}. \quad (7.10)$$

It will turn out that only the roots in $\{r^F : F \in \mathcal{F}\}$ are needed to construct subsolutions for importance sampling. This issue is discussed further in Section 9.3, where we clarify the role that the optimal trajectories play in determining those roots that are needed.

Proof of Proposition 7.7. We first consider the case of $E = \emptyset$, where $r^\emptyset = r^*$. By the definitions of H_\emptyset and r^* , Lemma 7.3, and condition (5.1),

$$\begin{aligned}
e^{H_\emptyset(-r^*)} &= \sum_{j=1}^d \Theta[e_j] e^{-\langle r^*, e_j \rangle} + \sum_{j=1}^d \Theta[-e_j] e^{-\langle r^*, -e_j \rangle} \\
&\quad + \sum_{i=1}^d \sum_{j=1}^d \Theta[-e_i + e_j] e^{-\langle r^*, -e_i + e_j \rangle} \\
&= \sum_{j=1}^d \frac{1}{\rho_j} \left(\lambda_j + \sum_{i=1}^d \mu_i P_{ij} \rho_i \right) + \sum_{j=1}^d \mu_j P_{j0} \rho_j \\
&= \sum_{j=1}^d [\mu_j + \lambda_j] \\
&= 1,
\end{aligned}$$

which yields $H_\emptyset(-r^*) = 0$.

Assume from now on that E is non-empty. It suffices to show that $H_\emptyset(-r^E) = H_i(-r^E) = 0$ for every $i \in E$, since Lemma 7.1 can then be inductively invoked to argue that $H_F(-r^E) = 0$ for all $F \subset E$. Fix an arbitrary $i \in E$. By the definitions of H_\emptyset and H_i ,

$$e^{H_\emptyset(-r^E)} - e^{H_i(-r^E)} = \sum_{v \in \mathbb{V}[-e_i]} \Theta[v] \left(e^{-\langle r^E, v \rangle} - 1 \right).$$

As observed in Remark 7.6

$$\sum_{v \in \mathbb{V}[-e_i]} \Theta[v] = \mu_i, \tag{7.11}$$

whereas it follows from the definition (7.9) that for any $v \in \mathbb{V}[-e_i]$

$$\begin{aligned}
\langle r^E, v \rangle &= \langle r^*, v \rangle - \sum_{k \in E} \log[z_k^E] \cdot \langle e_k, v \rangle \\
&= \langle r^*, v \rangle - \begin{cases} -\log z_i^E + \log z_j^E & \text{if } v = -e_i + e_j, j \in E, \\ -\log z_i^E & \text{otherwise.} \end{cases}
\end{aligned}$$

Since $z^E = \{z_i^E : i \in E\}$ is a solution to the linear equations (7.5), it is now straightforward calculation to show that

$$e^{H_\emptyset(-r^E)} - e^{H_i(-r^E)} = 0,$$

or $H_\emptyset(-r^E) = H_i(-r^E)$.

It remains to show that $H_\emptyset(-r^E) = 0$. To this end, observe that the definition of H_\emptyset and $H_\emptyset(-r^*) = 0$ imply

$$e^{H_\emptyset(-r^E)} - 1 = e^{H_\emptyset(-r^E)} - e^{H_\emptyset(-r^*)} = \sum_{v \in \mathbb{V}} \Theta[v] e^{-\langle r^*, v \rangle} \left(e^{\langle r^* - r^E, v \rangle} - 1 \right).$$

For each $i \in E$, define subsets of \mathbb{V} by

$$\begin{aligned} \mathbb{V}_1^{[i]} &\doteq \{v \in \mathbb{V}[-e_i] : v \neq -e_i + e_j, \text{ any } j \in E\}, \\ \mathbb{V}_2^{[i]} &\doteq \{v \in \mathbb{V}[-e_i] : v = -e_i + e_j, \text{ some } j \in E\}, \\ \mathbb{V}_3^{[i]} &\doteq \{v \in \mathbb{V} : v = -e_j + e_i, j \notin E\}, \\ \mathbb{V}_4^{[i]} &\doteq \{v \in \mathbb{V} : v = e_i\}. \end{aligned}$$

It is not difficult to see that $\{\mathbb{V}_k^{[i]} : i \in E, k = 1, \dots, 4\}$ are disjoint, and

$$e^{\langle r^* - r^E, v \rangle} - 1 = \begin{cases} 1/z_i^E - 1 & \text{if } v \in \mathbb{V}_1^{[i]}, \\ z_j^E/z_i^E - 1 & \text{if } v \in \mathbb{V}_2^{[i]} \text{ with } v = -e_i + e_j \text{ for some } j \in E, \\ z_i^E - 1 & \text{if } v \in \mathbb{V}_3^{[i]} \text{ or } v \in \mathbb{V}_4^{[i]}, \\ 0 & \text{otherwise.} \end{cases}$$

Therefore,

$$e^{H_\emptyset(-r^E)} - 1 = \sum_{i \in E} \left(d_1^{[i]} + d_2^{[i]} + d_3^{[i]} + d_4^{[i]} \right)$$

with

$$d_k^{[i]} \doteq \sum_{v \in \mathbb{V}_k^{[i]}} \Theta[v] e^{-\langle r^*, v \rangle} \left(e^{\langle r^* - r^E, v \rangle} - 1 \right)$$

for every $k = 1, \dots, 4$. Note that for every $i \in E$,

$$d_1^{[i]} + d_2^{[i]} = \sum_{v \in \mathbb{V}_1^{[i]}} \Theta[v] e^{-\langle r^*, v \rangle} (1/z_i^E - 1) + \sum_{\substack{v = -e_i + e_j \\ \text{some } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} (z_j^E/z_i^E - 1).$$

From (7.5) we obtain an alternative expression for the first sum, which gives

$$d_1^{[i]} + d_2^{[i]} = \mu_i (1 - z_i^E) + \sum_{\substack{v = -e_i + e_j \\ \text{some } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} (z_j^E - 1).$$

Therefore,

$$\sum_{i \in E} \left(d_1^{[i]} + d_2^{[i]} \right) = \sum_{i \in E} \mu_i (1 - z_i^E) + \sum_{i \in E} \sum_{\substack{v = -e_i + e_j \\ j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} (z_j^E - 1).$$

For the double-sum in the last display, one can use symmetry to interchange the roles of i and j and rewrite it as

$$\sum_{i \in E} \sum_{\substack{v = -e_j + e_i \\ j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} (z_i^E - 1).$$

Also note that

$$\sum_{i \in E} (d_3^{[i]} + d_4^{[i]}) = \sum_{i \in E} \sum_{\substack{v = -e_j + e_i \\ j \notin E}} \Theta[v] e^{-\langle r^*, v \rangle} (z_i^E - 1) + \sum_{\substack{i \in E \\ v = e_i}} \Theta[v] e^{-\langle r^*, v \rangle} (z_i^E - 1).$$

Therefore,

$$\sum_{i \in E} (d_1^{[i]} + d_2^{[i]} + d_3^{[i]} + d_4^{[i]}) = \sum_{i \in E} (z_i^E - 1) \left[-\mu_i + \sum_{v \in \mathbb{V}[e_i]} \Theta[v] e^{-\langle r^*, v \rangle} \right],$$

Invoking equality (7.7), we complete the proof. \blacksquare

Recall the definition of \mathcal{F} in (7.10). The proof of the following result is deferred to the appendix.

Lemma 7.8. *The following properties hold.*

1. $\emptyset \in \mathcal{F}$, $\{1, \dots, d\} \in \mathcal{F}$.
2. If $E_1, E_2 \in \mathcal{F}$, then $E_1 \cup E_2 \in \mathcal{F}$.
3. If $E, F \in \mathcal{F}$ and $E \subset F$, then $r^E \leq r^F$.

Here we adopt the notation that $u \leq v$ (or $\geq, <, >$) when u and v are vectors if the inequality holds component-wise. Taking $F = \{1, \dots, d\}$ in Part 3 of this lemma, we have the following result.

Corollary 7.9. *If $E \in \mathcal{F}$, then $z_i^E \in [\rho_i, 1)$ for all $i \in E$.*

8 The Isaacs equation

In this section we introduce the Isaacs equation that is relevant to importance sampling for Jackson networks. A heuristic motivation, which in any case would not be used in the analysis, can be found in [7] for two dimensional networks and in [8, 9] for other problem formulations. Subsolutions

to the Isaacs equation, together with a verification argument (one of the oldest techniques in optimal control), will be used to obtain bounds on the performance of related IS schemes.

With the convention $0 \log 0/0 = 0$, for each $\alpha \in \mathbb{R}^d$ we define the Hamiltonian

$$\mathbb{H}(x, \alpha) = \sup_{\bar{\Theta} \in \mathcal{P}^+(\mathbb{V})} \inf_{\theta \in \mathcal{P}^+(\mathbb{V})} \left[\langle \alpha, \mathbb{F}(x, \theta) \rangle + \sum_{v \in \mathbb{V}} \theta[v] \log \frac{\bar{\Theta}[v]}{\Theta[v]} + R(\theta \| \Theta) \right],$$

where

$$\mathbb{F}(x, \theta) \doteq \sum_{v \in \mathbb{V}} \theta[v] \cdot \pi[x, v],$$

and the relative entropy

$$R(\theta \| \Theta) \doteq \sum_{v \in \mathbb{V}} \theta[v] \log \frac{\theta[v]}{\Theta[v]}.$$

The Hamiltonian \mathbb{H} is associated with a differential game where the $\bar{\Theta}$ -player represents the choice of change of measure of an importance sampling scheme, whereas the θ -player is introduced from the relative entropy representation formula [5, Proposition 1.4.2] and characterizes certain large deviation properties.

For a function $W : \mathbb{R}_+^d \rightarrow \mathbb{R}$, denote by $DW(x)$ the gradient of W at x , if it exists. The Isaacs equation is then

$$\mathbb{H}(x, DW(x)) = 0, \quad x \in \mathbb{R}_+^d,$$

with the boundary condition

$$W(x) = 0, \quad x \in \Gamma.$$

Note that \mathbb{H} is discontinuous in x . The sense in which this equation should hold (as least with respect to subsolutions) will be discussed momentarily.

Consider a smooth subsolution to the Isaacs equation and a state nx . Then one importance sampling scheme corresponding to this subsolution is that which uses at nx the change of measure determined by the $\bar{\Theta}$ -component of the *saddle point* of the Hamiltonian $\mathbb{H}(x, \alpha)$, with α replaced by the gradient of the subsolution at x . As we will see, it is also possible to use a *mixture* of saddle point distributions when the gradient equals a convex combination of some fixed collection of vectors. This will be the case for the subsolutions we construct, and the use of mixtures will lead to schemes that are in some sense easier to implement.

The proof of the following result can be found in the appendix.

Lemma 8.1. Consider any $E \subset \{1, \dots, d\}$ and $x \in \partial_E$. For every $\alpha \in \mathbb{R}^d$,

$$\mathbb{H}(x, \alpha) = -2H(x, -\alpha/2) = -2H_E(-\alpha/2),$$

where H and H_E are defined in (7.1) and (7.2), respectively. The min/max condition holds in the definition of \mathbb{H} , and the saddle point is $(\bar{\Theta}^*(x, \alpha), \theta^*(x, \alpha)) = (\Theta_E^*(\alpha), \Theta_E^*(\alpha)) \in \mathcal{P}^+(\mathbb{V})$, where

$$\begin{aligned} \Theta_E^*(\alpha)[v] &\doteq e^{-H_E(-\alpha/2)} \cdot \Theta[v] \cdot e^{-\langle \alpha, \pi[x, v] \rangle / 2} \\ &= e^{-H_E(-\alpha/2)} \cdot \begin{cases} \Theta[v] \cdot e^{-\langle \alpha, v \rangle / 2} & \text{if } v \notin \cup_{i \in E} \mathbb{V}[-e_i], \\ \Theta[v] & \text{if } v \in \cup_{i \in E} \mathbb{V}[-e_i]. \end{cases} \end{aligned}$$

In particular, for every x , $\mathbb{H}(x, \cdot)$ is concave.

Remark 8.2. Due to the discontinuity of the state dynamics on the boundaries, a very natural interpretation of the Isaacs equation $\mathbb{H}(x, DW(x)) = 0$ is

$$\lim_{\epsilon \downarrow 0} \inf_{\{y \in \mathbb{R}_+^d : \|y-x\| \leq \epsilon\}} \mathbb{H}(y, DW(x)) = 0, \quad (8.1)$$

which is equivalent to $-2 \max\{H_F(-DW(x)/2) : F \subset E\} = 0$ if $x \in \partial_E$. The intuition is that the behavior of the controlled state process on the boundary ∂_E in the *limit* is an asymptotic characterization of the behavior of the scaled controlled process on all the boundaries ∂_F with $F \subset E$ in the *prelimit*. Owing to the infimum operation in (8.1), the subsolution property at x will imply that it holds, at least approximately, for all nearby points in the prelimit, and because of this one will be able to properly bound the performance of the corresponding importance sampling scheme. We will not explicitly refer to (8.1) in the rest of the paper, since the proof of asymptotic optimality [Theorem 10.1] does not need this formulation. It is worth noting though that the piecewise affine subsolution that we will construct in Section 9 will be shown to satisfy (8.1) on all those x where the subsolution is continuously differentiable [see the proof of Lemma 9.2]. The identification of the optimal trajectory in Section 11 calls upon (8.1) in the proof of Theorem 11.2.

9 Construction of classical subsolutions

A *classical subsolution* is a continuously differentiable function $W : \mathbb{R}_+^d \rightarrow \mathbb{R}$ such that $\mathbb{H}(x, DW(x)) \geq 0$ for all $x \in \mathbb{R}_+^d$ and $W(x) \leq 0$ for all $x \in \Gamma$. As was remarked at the beginning of the last section, these functions can be used

to bound the performance of related IS schemes. $W(0)$ is the key measure of performance, with large values indicating greater variance reduction. Hence methods for the systematic construction of subsolutions with optimal or nearly optimal $W(0)$ are very useful. Since the state dynamics are piecewise homogeneous, it is natural to start with the construction of piecewise affine subsolutions and then use an exponential weighting mollification to obtain a classical subsolution [9, 7, 6].

9.1 Piecewise affine subsolution

A concave, piecewise affine function $\bar{W} : \mathbb{R}_+^d \rightarrow \mathbb{R}$ is said to be a *subsolution* to the Isaacs equation if $\mathbb{H}(x, D\bar{W}(x)) \geq 0$ for all x where $D\bar{W}(x)$ is well defined and $\bar{W}(x) \leq 0$ for all $x \in \Gamma$. It is often identified as the minimum of a finite collection of affine functions. The idea is to start with $-\langle r^*, x \rangle$, which is a solution on all of the interior of \mathbb{R}_+^d , and then perturb this function slightly near the boundaries.

Let $\{C_F : F \in \mathcal{F}\}$ be an arbitrary collection of non-negative constants that satisfy

Condition 9.1. $C_\emptyset = 0$ and for any nonempty $F \in \mathcal{F}$,

$$C_F > \max \left\{ \frac{\log(1/z_i^F)}{\log(1/z_i^G)} \cdot C_G : i \in G \subset F, G \neq F, G \in \mathcal{F} \right\}$$

with the convention that $\max\{\emptyset\} = 0$.

There are infinitely many such collections of constants. Indeed, for all those $F \in \mathcal{F}$ that are singletons, C_F can be defined as any positive number. Once those C_F are fixed, one can recursively define C_F for those F with $|F| = 2$, and so on so forth. Moreover, it is not difficult to see that $C_F > 0$ for all nonempty $F \in \mathcal{F}$ and by Lemma 7.8 $C_G < C_F$ if G is a strict subset of F .

Fixing an arbitrary small positive number δ , we define a piecewise affine function \bar{W}^δ by

$$\bar{W}^\delta(x) \doteq \min_{F \in \mathcal{F}} \bar{W}_F^\delta(x), \quad \bar{W}_F^\delta(x) \doteq 2\gamma + \langle 2r^F, x \rangle - C_F \delta. \quad (9.1)$$

Note that for all $x \in \Gamma$, by assumption (3.3)

$$\bar{W}^\delta(x) \leq \bar{W}_\emptyset^\delta(x) = 2\gamma + \langle 2r^\emptyset, x \rangle = 2\gamma + \langle 2r^*, x \rangle \leq 0. \quad (9.2)$$

Lemma 9.2. *Fix any collection of constants $\{C_F : F \in \mathcal{F}\}$ that satisfy Condition 9.1. Then there exists a constant $a > 0$ such that for any $x \in \mathbb{R}_+^d$ and any $\delta > 0$, if $G \in \mathcal{F}$ satisfies*

$$\bar{W}_G^\delta(x) - \bar{W}^\delta(x) < a\delta,$$

then $\mathbb{H}(x, 2r^G) \geq 0$.

The proof of the lemma is deferred to the appendix. Applying this lemma when G is a minimizer in the definition of $\bar{W}^\delta(x)$, we see that the function \bar{W}^δ is indeed a piecewise affine subsolution. From now on, we will assume throughout that $\{C_F : F \in \mathcal{F}\}$ is a collection of constants that satisfy Condition 9.1.

9.2 Mollification

To obtain a smooth approximation of the piecewise affine subsolution \bar{W}^δ , we use the *exponential weighting* mollification as follows. Fix an arbitrary $\varepsilon > 0$. Define

$$W^{\varepsilon, \delta}(x) \doteq -\varepsilon \log \sum_{F \in \mathcal{F}} \exp \left\{ -\frac{1}{\varepsilon} \bar{W}_F^\delta(x) \right\}. \quad (9.3)$$

The following lemma is standard, and we omit the proof [9].

Lemma 9.3. *For any $\varepsilon, \delta > 0$, the function $W^{\varepsilon, \delta}$ is continuously differentiable. Furthermore, for every $x \in \mathbb{R}_+^d$,*

$$-\log |\mathcal{F}| \cdot \varepsilon \leq W^{\varepsilon, \delta}(x) - \bar{W}^\delta(x) \leq 0,$$

$$DW^{\varepsilon, \delta}(x) = \sum_{F \in \mathcal{F}} \rho_F^{\varepsilon, \delta}(x) \cdot 2r^F, \quad \rho_F^{\varepsilon, \delta}(x) \doteq \frac{\exp\{-\bar{W}_F^\delta(x)/\varepsilon\}}{\sum_{E \in \mathcal{F}} \exp\{-\bar{W}_E^\delta(x)/\varepsilon\}}.$$

In particular, for every x , $\{\rho_F^{\varepsilon, \delta}(x) : F \in \mathcal{F}\}$ defines a probability vector on \mathcal{F} . Moreover,

$$\rho_F^{\varepsilon, \delta}(x) \doteq \frac{\exp\{[C_F \delta - \langle 2r^F, x \rangle]/\varepsilon\}}{\sum_{E \in \mathcal{F}} \exp\{[C_E \delta - \langle 2r^E, x \rangle]/\varepsilon\}}$$

is independent of γ and hence Γ .

Proposition 9.4. *The function $W^{\varepsilon,\delta}$ satisfies $W^{\varepsilon,\delta}(x) \leq 0$ for $x \in \Gamma$. Furthermore, for all $x \in \mathbb{R}_+^d$,*

$$\mathbb{H}(x, DW^{\varepsilon,\delta}(x)) \geq \sum_{F \in \mathcal{F}} \rho_F^{\varepsilon,\delta}(x) \cdot \mathbb{H}(x, 2r^F) \geq -M \exp\{-a\delta/\varepsilon\},$$

where M and a (the constant given by Lemma 9.2) are two positive constants both independent of ε and δ .

This result establishes that $\bar{W}^{\varepsilon,\delta}$ is “approximately” a classical subsolution as long as ε/δ is small, which is sufficient for the subsequent analysis. See [7] for more discussion on exponential weighting mollification.

Proof of Proposition 9.4. $W^{\varepsilon,\delta}(x) \leq 0$ for $x \in \Gamma$ is trivial from Lemma 9.3 and (9.2). Consider any $x \in \mathbb{R}_+^d$ and assume $x \in \partial_E$. We will argue that, for all $F \subset E$,

$$2H_F(-DW^{\varepsilon,\delta}(x)/2) \leq \sum_{G \in \mathcal{F}} \rho_G^{\varepsilon,\delta}(x) \cdot 2H_F(-r^G) \leq M \exp\{-a\delta/\varepsilon\}. \quad (9.4)$$

The second part of the proposition is a special case of (9.4) with $F \doteq E$, thanks to Lemma 8.1.

The first inequality in (9.4) is trivial from the convexity of $H_F(\cdot)$. As for the second inequality, note that by the proof of Lemma 9.2, more precisely, inequality (A.5), $H_F(-r^G) \leq 0$ for all those G such that $\bar{W}_G^\delta(x) - \bar{W}^\delta(x) < a\delta$. For all those G such that $\bar{W}_G^\delta(x) - \bar{W}^\delta(x) \geq a\delta$, the formula for $\rho_G^{\varepsilon,\delta}(x)$ in Lemma 9.3 together with the fact that $\bar{W}^\delta(x) = \bar{W}_E^\delta(x)$ for some $E \in \mathcal{F}$ yield that

$$\rho_G^{\varepsilon,\delta}(x) \leq \exp\left\{\frac{-\bar{W}_G^\delta(x) + \bar{W}^\delta(x)}{\varepsilon}\right\} \leq \exp\{-a\delta/\varepsilon\}.$$

Inequality (9.4) follows readily for some M independent of ε and δ . ■

9.3 Remarks on the construction of subsolutions

The purpose of this section is to discuss the motivation and intuition of the roots r^F with $F \in \mathcal{F}$, which are clearly the key quantities in the construction of subsolutions. This section can be safely skipped by readers interested only in algorithmic aspects of the construction. The discussion here is informal, and assumes familiarity with the theory of large deviations.

It should be noted first that the Isaacs equation is closely related to the PDE associated with the corresponding large deviation problem, which takes the form

$$-H(x, -DU(x)) = 0, \quad U(x) = 0 \text{ for } x \in \partial\Gamma. \quad (9.5)$$

By Lemma 8.1, it suffices to construct a subsolution, say \bar{U} , to the large deviation PDE (9.5) and $\bar{W} \doteq 2\bar{U}$ will then be a subsolution to the Isaacs equation.

The central question in the construction of \bar{U} is to determine the appropriate gradients in different regions of the state space. There are two considerations in the determination of these gradients: (a) the subsolution properties should be satisfied everywhere; (b) along the optimal path (i.e., least cost path), the subsolution should satisfy the PDE with *equality*. The latter is required so that $\bar{U}(0) = \gamma$ or $\bar{W}(0) = 2\gamma$, which is indeed necessary for the optimality of the corresponding importance sampling estimator since $\bar{W}(0)$ will characterize the lower bound on the exponential decay rate of its second moment.

Here it useful to recall that the goal in this paper is to establish importance sampling schemes that do not depend on the particular hitting set Γ . The distinction is crucial, in that it leads to very different approaches to the construction of subsolutions. When the hitting set is fixed, it is natural to think of the subsolution as solving a relaxed form of (9.5). Suppose that ϕ is the trajectory that minimizes the large deviation rate function over all paths that connect x to Γ , and that ϕ first touches Γ at y . Then one can relax the condition $U(x) = 0$ for $x \in \partial\Gamma$ to $\bar{U}(x) \leq 0$ on $\partial\Gamma$, so long as one keeps $\bar{U}(0) = U(0)$. It is often the case, especially for queueing type problems and for sets Γ with a simple shape, that there are very simple functions which satisfy the relaxed problem, even though $U(x)$ itself has no explicit representation. See the examples in [7] and Section 12.

On the other hand, a natural starting point for schemes that are not directly tied to the hitting set Γ is the *quasipotential function* $V(x)$ from large deviations theory (see Section 11). The quasipotential is defined as the minimum of the large deviation rate over all trajectories that connect 0 to x , and is known to characterize the large deviation properties of invariant distributions [10, Chapter 6]. Observe that with the quasipotential the free variable is the terminal location rather than the initial location. Hence it does not have a direct interpretation as the value function of an optimal control problem. When the quasipotential is smooth and in the setting of processes with smooth statistics, one can show that V satisfies $-H(x, DV(x)) = 0$. Under

these circumstances, for *any* set Γ a subsolution can be found by setting $a = \inf\{V(x) : x \in \Gamma\}$ and then choosing $\bar{U}(x) = -V(x) + a$. Difficulties associated with this approach are that the quasipotential is not usually smooth, nor is it usually available in any explicit form.

For Jackson networks the invariant distribution, and hence the quasipotential, are known in closed form. In fact, one has $V(x) = \langle -r^*, x \rangle$, where r^* is defined as in Proposition 3.1. Although V is obviously smooth, because of the discontinuous statistics it does not satisfy the subsolution property on all the boundaries and must be modified there (although it is a subsolution everywhere on the interior). Hence there are two questions: how do we identify the parts of the boundary where an alternative gradient is needed, and for those parts what is a good gradient?

The issue of where alternative gradients are needed is related to how well the quasipotential approximates the invariant distribution at the *prelimit*. Consider a point $x \in \partial_E$, and consider the minimizing trajectory which connects the origin to x in the definition of the quasipotential. If the first time this trajectory reaches ∂_E after leaving the origin is when it hits x , then in some sense dynamics other than the interior dynamics do not contribute substantially to the behavior of the invariant distribution near x , and the subsolution property will hold without modification. On the other hand, if the trajectory runs along ∂_E to reach x then we must distinguish between whether the trajectory in some sense “pushes into” ∂_E or not. This distinction is quantified by the z_i^E introduced in Section 7, and in fact the trajectory pushes into ∂_E if and only if $z_i^E < 1$ for all $i \in E$. This will be shown in Section 11, where it is more convenient to first construct the time-reversed trajectories. It is only for the parts of the boundary where the condition $z_i^E < 1$ for all $i \in E$ holds that the prelimit invariant distribution reflects the different statistical behaviors of the boundary dynamics, and for which the quasipotential must be modified to produce a subsolution with respect to the boundary dynamics. This motivates the definition of \mathcal{F} in (7.10).

Having motivated the identification of the boundaries where modification is needed, the remaining question is how one should modify. Under the condition where modification is required all boundary behaviors must be accounted for by the subsolution. Thus if we consider a gradient r for use in a neighborhood of ∂_E , then we would like r to satisfy the equation $-H_F(-r) = 0$ for all $F \subset E$. By Lemma 7.1, this amounts to

$$-H_\emptyset(-r) = 0, \quad -H_i(-r) = 0 \text{ for all } i \in E. \quad (9.6)$$

Moreover, if the optimal path goes from point x to point y on ∂_E , then the

cost should satisfy

$$\langle -r^*, x - y \rangle = \langle -r, x - y \rangle.$$

In other words, we would like r to be a *projection* of r^* such that $r_i = r_i^*$ for all $i \notin E$ so that the above display holds automatically. This leads to the equation

$$r = r^* - \sum_{i \in E} u_i e_i,$$

where $\{u_i : i \in E\}$ are unknown real numbers. There are in total $|E|$ unknowns, and $|E| + 1$ equations (9.6). But Lemma 7.4 and Proposition 7.7 indeed claim that these $|E|$ unknowns are uniquely determined by these $|E| + 1$ equations, with $u_i = \log z_i^E$.

10 The importance sampling algorithm

Recall the formula for the saddle point in Lemma 8.1. One approach to the design of importance sampling schemes begins with the smooth function $W^{\varepsilon, \delta}(x)$, and then uses the stochastic kernel defined by the saddle point $\bar{\Theta}^*(x, DW^{\varepsilon, \delta}(x))[v]$. While this would yield an algorithm with very good performance, it is more difficult to implement than a closely related mixture of saddle points.

For every $F \in \mathcal{F}$, define a stochastic kernel on \mathbb{V} given \mathbb{R}_+^d by the saddle point formula in Lemma 8.1

$$\bar{\Theta}_F[v|x] \doteq \bar{\Theta}^*(x, 2r^F)[v]. \quad (10.1)$$

Note that $\bar{\Theta}_F$ is piecewise constant in that $\bar{\Theta}^*(x, 2r^F) \equiv \Theta_E^*(2r^F)$ for all $x \in \partial_E$ and all $E \subset \{1, \dots, d\}$. The change of measure we will associate with $W^{\varepsilon, \delta}$ is the mixture of $\{\bar{\Theta}_F : F \in \mathcal{F}\}$ given by

$$\bar{\Theta}^{\varepsilon, \delta}[\cdot|x] \doteq \sum_{F \in \mathcal{F}} \rho_F^{\varepsilon, \delta}(x) \bar{\Theta}_F[\cdot|x]. \quad (10.2)$$

In general, the parameters ε and δ are allowed to depend on n , denoted by ε_n and δ_n respectively. To estimate p_n , the importance sampling algorithm uses a change of measure, say \mathbb{Q}_n , under which the transition kernel is

$$\bar{\Theta}^n \doteq \bar{\Theta}^{\varepsilon_n, \delta_n}$$

and the corresponding estimator \hat{p}_n is defined as in (6.3). The next result establishes the asymptotically optimality of \hat{p}_n in any order of moments, under suitable decay conditions on the parameters $\{\varepsilon_n, \delta_n\}$.

Theorem 10.1. *Suppose that $\delta_n \rightarrow 0$, $\varepsilon_n/\delta_n \rightarrow 0$, and $n\varepsilon_n \rightarrow \infty$. Then for any real number $b > 1$, the importance sampling estimator \hat{p}_n satisfies*

$$\lim_n \frac{1}{n} \log E^{\mathbb{Q}_n} \left[\hat{p}_n^b \right] = -b\gamma.$$

In particular, \hat{p}_n is asymptotically optimal.

The proof of this result, which uses ideas developed in [7], is deferred to the appendix. It uses a verification argument, with the main difficulty due to the fact that the time interval of interest is potentially unbounded.

Remark 10.2. As we have remarked previously, a particular feature of the schemes constructed in this paper is that they are *independent* of γ and Γ , since $\rho_F^{\varepsilon, \delta}(x)$ and $\bar{\Theta}_F[\cdot|x]$ are both so by Lemma 9.3 and Lemma 8.1. In other words, for different target sets Γ , one can use the *same* state-dependent change of measure (10.2) for simulation, and indeed probabilities for distinct sets can be estimated simultaneously. This is very appealing for its simplicity and universality. However, we reiterate that with more knowledge of Γ it is sometimes possible to construct even simpler subsolutions, i.e., fewer affine pieces.

11 The optimal trajectory

A by-product of the subsolution analysis is the solution to the calculus of variation problem associated with the sample path large deviations for Jackson networks. There are two versions, one for the discrete time process Z , and the other for the continuous time process Q . These two problems can be related by a change of the time variable, and are equivalent when the calculus of variations problem does not depend on time. We will treat them simultaneously.

Discrete time version. The calculus of variation problem is

$$V(x^*) \doteq \inf \left\{ \int_0^\tau L(\phi(t), \dot{\phi}(t)) dt : \phi(0) = 0, \phi(\tau) = x^*, \phi(t) \in \mathbb{R}_+^d \text{ for all } t \right\}.$$

Here x^* is an arbitrary point in \mathbb{R}_+^d and the local rate function L is defined, for all $x \in \partial_E$ and $\beta \in \mathbb{R}^d$, as

$$L(x, \beta) \doteq \inf \left\{ \sum_{G \subset E} \rho_G L_G(\beta_G) : \rho_G \geq 0, \sum_{G \subset E} \rho_G = 1, \sum_{G \subset E} \rho_G \beta_G = \beta \right\},$$

where L_G is the Legendre transform of the convex function H_G .

Continuous time version. For every $G \subset \{1, 2, \dots, d\}$, define the convex function \bar{H}_G by

$$\bar{H}_G(\alpha) \doteq e^{H_G(\alpha)} - 1 = \sum_{v \notin \cup_{i \in E} \mathbb{V}[-e_i]} \Theta[v] \cdot \left(e^{\langle \alpha, v \rangle} - 1 \right).$$

In other words, \bar{H}_G is the Hamiltonian of the continuous time process Q on the boundary ∂_G . Furthermore $\bar{H}_G(\alpha) = 0$ (resp., “ \geq ”, “ \leq ”) if and only if $H_G(\alpha) = 0$ (resp., “ \geq ”, “ \leq ”). The calculus of variation problem associated with sample path large deviations for the continuous time process Q is

$$\bar{V}(x^*) \doteq \inf \left\{ \int_0^\tau \bar{L}(\phi(t), \dot{\phi}(t)) dt : \phi(0) = 0, \phi(\tau) = x^*, \phi(t) \in \mathbb{R}_+^d \text{ for all } t \right\}.$$

Here the local rate function \bar{L} is defined, for all $x \in \partial_E$ and $\beta \in \mathbb{R}^d$, as

$$\bar{L}(x, \beta) \doteq \inf \left\{ \sum_{G \subset E} \rho_G \bar{L}_G(\beta_G) : \rho_G \geq 0, \sum_{G \subset E} \rho_G = 1, \sum_{G \subset E} \rho_G \beta_G = \beta \right\},$$

with \bar{L}_G the Legendre transform of \bar{H}_G .

It will turn out that $V \equiv \bar{V}$ and the optimal trajectories take the same form for both problems. We introduce the notation

$$\bar{\mathcal{F}} \doteq \{F \subset \{1, \dots, d\} : z_i^F \leq 1 \text{ for all } i \in F\}$$

Recall that \mathcal{F} has the analogous definition but with \leq replaced by $<$. Clearly $\mathcal{F} \subset \bar{\mathcal{F}}$. In particular, $\emptyset \in \bar{\mathcal{F}}$ and $\{1, \dots, d\} \in \bar{\mathcal{F}}$. Almost verbatim to the proof of Lemma 7.8 [we omit the details], one can argue that $F_1 \cup F_2 \in \bar{\mathcal{F}}$ if $F_1, F_2 \in \bar{\mathcal{F}}$. Therefore, for every $E \subset \{1, \dots, d\}$, we can define the *maximal* subset of E that belong to $\bar{\mathcal{F}}$, that is,

$$\Pi[E] \doteq \cup_{F \subset E, F \in \bar{\mathcal{F}}} F.$$

Moreover, for any $F \in \bar{\mathcal{F}}$, define the vector

$$\beta_F^* \doteq \sum_{v \in \mathbb{V}} \Theta[v] e^{-\langle r^F, v \rangle} v - \sum_{i \in F} (1 - z_i^F) \sum_{v \in \mathbb{V}[-e_i]} \Theta[v] e^{-\langle r^F, v \rangle} v. \quad (11.1)$$

The structural properties and the cost of β_F^* is given by the following lemma, whose proof is deferred to the appendix.

Lemma 11.1. Fix any strict subset $E \subset \{1, \dots, d\}$ and let $F \doteq \Pi[E]$.

1. $\langle \beta_F^*, e_i \rangle = 0$ for all $i \in F$, and $\langle -\beta_F^*, e_i \rangle > 0$ for all $i \in E \setminus F$.
2. The set $G \doteq \{j : \langle -\beta_F^*, e_j \rangle < 0\}$ is non-empty. Furthermore, for all $G_1 \subset G$, we have $F \cup G_1 \in \bar{\mathcal{F}}$.
3. For any $x \in \partial_F$, $L(x, \beta_F^*) = \bar{L}(x, \beta_F^*) = \langle -r^*, \beta_F^* \rangle$.

We are now ready now to describe the optimal trajectory, which is done most conveniently in a time-reversed fashion. Assume $x^* \in \partial_E$ for some strict subset $E \subset \{1, \dots, d\}$. Denote $F_0 = \Pi[E]$. The trajectory starts with constant velocity $-\beta_{F_0}^*$. Thanks to Lemma 11.1 Part 1, the trajectory immediately leaves ∂_E (unless $E = F_0$) and travels on ∂_{F_0} for a period of time until it leaves ∂_{F_0} and hits a new boundary, say ∂_{F_1} . By Lemma 11.1 Part 2, the trajectory will hit ∂_{F_1} in finite amount of time and F_1 contains F_0 as a *strict* subset with $F_1 \in \bar{\mathcal{F}}$. The trajectory then switches to velocity $-\beta_{F_1}^*$, and will travel on ∂_{F_1} for a finite amount of time until it leaves ∂_{F_1} and hits, say ∂_{F_2} . Again, $F_2 \in \bar{\mathcal{F}}$ and contains F_1 as a strict subset. The trajectory then switches to velocity $-\beta_{F_2}^*$, and so on, until it hits the origin $\partial_{\{1, \dots, d\}} = \{0\}$. Denote this path by φ in the time reversal setting with τ^* the hitting time to the origin. Switching to the forward-time setup, we define the trajectory

$$\phi^*(t) \doteq \varphi(\tau^* - t), \quad 0 \leq t \leq \tau^*.$$

Clearly $\phi^*(0) = 0$, $\phi^*(\tau^*) = x^*$, and $\phi^*(t) \in \mathbb{R}_+^d$. Moreover, ϕ^* is piecewise affine with at most d pieces.

Theorem 11.2. $V(x^*) = \bar{V}(x^*) = \langle -r^*, x^* \rangle$ and ϕ^* is an optimal trajectory for both calculus of variation problems.

Proof. We will give the proof for the discrete time version V . The continuous time version \bar{V} is almost verbatim and thus omitted.

Let $F^*(t)$ be such that $\phi^*(t) \in \partial_{F^*(t)}$. By the definition of ϕ^* , it is easy to see that $F^*(t) \in \bar{\mathcal{F}}$ and $\dot{\phi}^*(t) = \beta_{F^*(t)}^*$ for all $t \in [0, \tau^*) \setminus S$, where S is the finite set of times when the velocity changes. Therefore, by Part 3 of Lemma 11.1,

$$\begin{aligned} \int_0^{\tau^*} L(\phi^*(t), \dot{\phi}^*(t)) dt &= \int_0^{\tau^*} \langle -r^*, \dot{\phi}^*(t) \rangle dt \\ &= \langle -r^*, \phi^*(\tau^*) - \phi^*(0) \rangle \\ &= \langle -r^*, x^* \rangle. \end{aligned}$$

It remains to show that $V(x^*) \geq \langle -r^*, x^* \rangle$. This can be done using a verification argument.

Fix arbitrarily $\varepsilon, \delta > 0$, and any absolutely continuous function ϕ such that $\phi(0) = 0$, $\phi(\tau) = x^*$, and $\phi(t) \in \mathbb{R}_+^d$ for all $t \in [0, \tau]$. Abusing the notation, define $\bar{W}^\delta, W^{\varepsilon, \delta}$ as in (9.1) and (9.3), except that γ is replaced by $\langle -r^*, x^* \rangle$. Note that

$$W^{\varepsilon, \delta}(x^*) - W^{\varepsilon, \delta}(0) = \int_0^\tau \langle DW^{\varepsilon, \delta}(\phi(t)), \dot{\phi}(t) \rangle dt. \quad (11.2)$$

However, for any $x \in \partial_E$, $L(x, \cdot)$ is the inf-convolution of $\{L_G : G \subset E\}$, and therefore it is the Legendre transform of the convex function $\max\{H_G : G \subset E\}$ [5, Appendix D]. In particular, for any $\alpha, \beta \in \mathbb{R}^d$,

$$L(x, \beta) \geq \langle \alpha, \beta \rangle - \max_{G \subset E} H_G(\alpha).$$

For any t , denote by $E(t)$ the subset of $\{1, \dots, d\}$ such that $\phi(t) \in \partial_{E(t)}$. The previous inequality implies that for $t \notin S$

$$L(\phi(t), \dot{\phi}(t)) \geq \langle -DW^{\varepsilon, \delta}(\phi(t))/2, \dot{\phi}(t) \rangle - \max_{G \subset E(t)} H_G(-DW^{\varepsilon, \delta}(\phi(t))/2).$$

It follows from the proof of Proposition 9.4, more precisely, inequality (9.4), that

$$\max_{G \subset E(t)} H_G(-DW^{\varepsilon, \delta}(\phi(t))/2) \leq M \exp\{-a\delta/\varepsilon\}/2.$$

Therefore

$$L(\phi(t), \dot{\phi}(t)) \geq \langle -DW^{\varepsilon, \delta}(\phi(t))/2, \dot{\phi}(t) \rangle - M \exp\{-a\delta/\varepsilon\}/2.$$

Recalling (11.2), we arrive at

$$W^{\varepsilon, \delta}(x^*) - W^{\varepsilon, \delta}(0) \geq -2 \int_0^\tau L(\phi(t), \dot{\phi}(t)) dt - \tau M \exp\{-a\delta/\varepsilon\}.$$

But by Lemma 9.3

$$\begin{aligned} W^{\varepsilon, \delta}(0) &\geq \bar{W}^\delta(0) - \varepsilon \cdot \log |\mathcal{F}| \\ &= \langle -2r^*, x^* \rangle - C_{\{1, \dots, d\}} \delta - \varepsilon \cdot \log |\mathcal{F}| \end{aligned}$$

and

$$W^{\varepsilon, \delta}(x^*) \leq \bar{W}^\delta(x^*) \leq \bar{W}_\emptyset^\delta(x^*) = 0.$$

Therefore,

$$\langle 2r^*, x^* \rangle + C_{\{1, \dots, d\}} \delta + \varepsilon \cdot \log |\mathcal{F}| \geq -2 \int_0^\tau L(\phi(t), \dot{\phi}(t)) dt - \tau M \exp\{-a\delta/\varepsilon\}.$$

Letting $\delta \rightarrow 0$, $\varepsilon/\delta \rightarrow 0$, we have

$$\int_0^\tau L(\phi(t), \dot{\phi}(t)) dt \geq \langle -r^*, x^* \rangle.$$

Since ϕ is arbitrary, we complete the proof. ■

12 Numerical Examples

In this section, we will give a number of numerical examples for illustration. For all the examples, three types of buffer overflow probabilities, that is,

1. total population overflow: $\Gamma = \{x \in \mathbb{R}_+^d : x_1 + \dots + x_d \geq 1\}$,
2. individual buffer overflow: $\Gamma = \{x \in \mathbb{R}_+^d : x_1 \geq 1 \text{ or } \dots \text{ or } x_d \geq 1\}$,
3. overflow of the first buffer: $\Gamma = \{x \in \mathbb{R}_+^d : x_1 \geq 1\}$,

will be estimated using the *same* change of measure; see Remark 10.2. The parameters in the simulation are set up as follows. Let $C_\emptyset = 0$ and for each $\{i\} \in \mathcal{F}$

$$C_{\{i\}} \doteq \log(1/z_i^{\{i\}}),$$

For all $F \in \mathcal{F}$ such that $|F| \geq 2$, we recursively define

$$C_F \doteq (1+k) \max \left\{ \frac{\log(1/z_i^F)}{\log(1/z_i^G)} \cdot C_G : i \in G \subset F, G \neq F, G \in \mathcal{F} \right\}$$

for some parameter $k > 0$. Also for any δ we let $\varepsilon \doteq -\delta/\log \delta$. Each estimate is based on 100,000 samples and the theoretical values are obtained by recursively solving a linear equation from first step analysis. Of course the theoretical value cannot be found in general, and we restrict the examples to those for which its computation is feasible.

12.1 A two-node feedback network

Two-node feedback networks are of great value for the illustration of the importance sampling schemes and optimal trajectories, since everything can be explicitly computed. Consider the following Jackson network with

$$\lambda = (\lambda_1, \lambda_2), \quad \mu = (\mu_1, \mu_2), \quad P = \begin{bmatrix} 0 & p \\ \bar{p} & 0 \end{bmatrix}$$

where $p\bar{p} < 1$. The utilization parameters are, by formula (2.2),

$$\rho_1 = \frac{1}{1 - p\bar{p}} \cdot \frac{\lambda_1 + \bar{p}\lambda_2}{\mu_1}, \quad \rho_2 = \frac{1}{1 - p\bar{p}} \cdot \frac{\lambda_2 + p\lambda_1}{\mu_2}.$$

We assume the system is stable or $\rho_1, \rho_2 \in (0, 1)$. See Figure 1.

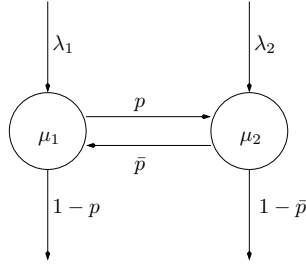


Figure 1: Two-node feedback networks

Consider all possible subsets $E \subset \{1, 2\}$.

1. For $E = \emptyset$, $r^\emptyset = r^* = (\log \rho_1, \log \rho_2)$.
2. For $E = \{1, 2\}$, $r^{\{1,2\}} = (0, 0)$ by Remark 7.6.
3. For $E = \{1\}$, equation (7.5) yields $z_1^{\{1\}} = (1 - p)\rho_1 + p\rho_1/\rho_2$, and $r^{\{1\}} = (\log \rho_1 - \log z_1^{\{1\}}, \log \rho_2)$.
4. For $E = \{2\}$, equation (7.5) yields $z_2^{\{2\}} = (1 - \bar{p})\rho_2 + \bar{p}\rho_2/\rho_1$, and $r^{\{2\}} = (\log \rho_1, \log \rho_2 - \log z_2^{\{2\}})$.

It is not difficult to show that $z_1^{\{1\}} \geq 1$ and $z_2^{\{2\}} \geq 1$ cannot hold simultaneously, and Theorem 11.2 yields the behavior of the optimal trajectories, as depicted in Figure 2.

In these figures and Figure 3 we see the relationship between the optimal trajectories, the values of z_i^E , and those parts of the boundary where the

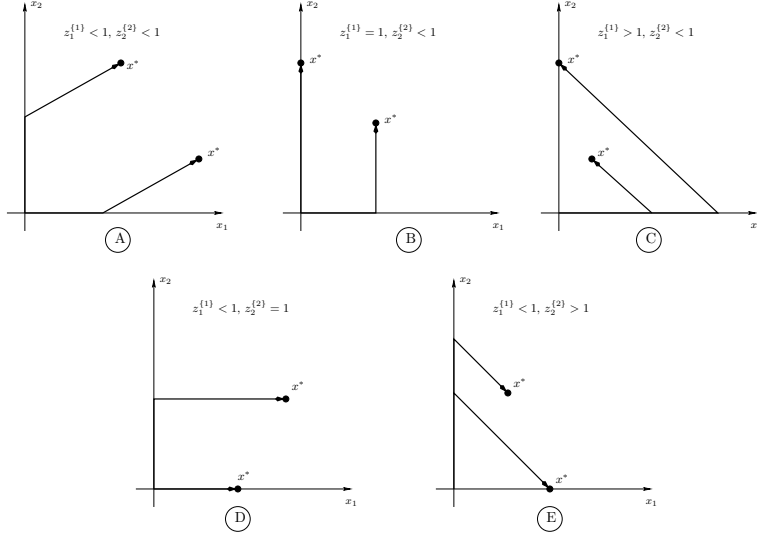


Figure 2: Optimal trajectories

affine function $-\langle r^*, x \rangle$ should be perturbed to produce a subsolution (see also the discussion in Section 9.3). In A $z_1^{\{1\}} < 1$ and $z_2^{\{2\}} < 1$, the time-reversed trajectory pushes into both boundaries, and the interior solution must be perturbed near both boundaries. In B $z_1^{\{1\}} = 1$ and $z_2^{\{2\}} < 1$, the time-reversed trajectory is parallel to $\{x : x_1 = 0\}$ and pushes into the other boundary, and the interior solution is perturbed only near this second boundary. Finally in C $z_1^{\{1\}} > 1$ and $z_2^{\{2\}} < 1$, the time-reversed trajectory moves away from $\{x : x_1 = 0\}$ and pushes into the other boundary, and again the interior solution is perturbed only near the second boundary.

As far as the importance sampling scheme is concerned, the piecewise affine subsolution \bar{W}^δ partitions the state space into several regions (this partition is independent of the value of γ), and is illustrated in Figure 3. Note that the “boundary layer” in Figure 3 on x_1 and/or x_2 axis has width δ under our choice of $\{C_F : F \in \mathcal{F}\}$.

For any target set Γ , the change of measure is the same and given by (10.1)–(10.2). Below is a collection of simulation results with

$$\lambda = (0.1, 0), \quad \mu = (0.5, 0.4), \quad p = 1, \quad \bar{p} = 0.1.$$

The simulation parameters are set as $\delta = 0.1, k = 0.5$.

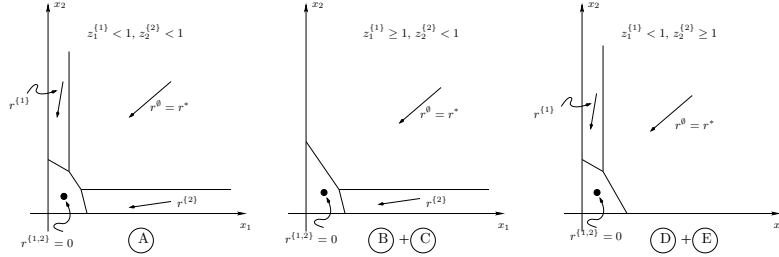


Figure 3: Gradients of the piecewise affine subsolution

	$n = 20$	$n = 40$	$n = 60$
Theoretical value	9.60×10^{-11}	7.27×10^{-22}	5.44×10^{-33}
Estimate	9.64×10^{-11}	7.42×10^{-22}	5.54×10^{-33}
Std. Err.	0.09×10^{-11}	0.11×10^{-22}	0.17×10^{-33}
95% C.I.	$[9.46, 9.82] \times 10^{-11}$	$[7.20, 7.64] \times 10^{-22}$	$[5.20, 5.88] \times 10^{-33}$

Table 1. Two-node Jackson network, total population overflow

	$n = 20$	$n = 40$	$n = 60$
Theoretical value	2.66×10^{-11}	1.96×10^{-22}	1.46×10^{-33}
Estimate	2.60×10^{-11}	1.96×10^{-22}	1.39×10^{-33}
Std. Err.	0.06×10^{-11}	0.08×10^{-22}	0.08×10^{-33}
95% C.I.	$[2.48, 2.72] \times 10^{-11}$	$[1.80, 2.12] \times 10^{-22}$	$[1.23, 1.55] \times 10^{-33}$

Table 2. Two-node Jackson network, individual buffer overflow

	$n = 20$	$n = 40$	$n = 60$
Theoretical value	4.44×10^{-13}	3.83×10^{-26}	3.31×10^{-39}
Estimate	4.32×10^{-13}	3.89×10^{-26}	3.28×10^{-39}
Std. Err.	0.09×10^{-13}	0.07×10^{-26}	0.06×10^{-39}
95% C.I.	$[4.14, 4.50] \times 10^{-13}$	$[3.75, 4.03] \times 10^{-26}$	$[3.16, 3.40] \times 10^{-39}$

Table 3. Two-node Jackson network, overflow of the first buffer

12.2 A three-node tandem Jackson networks

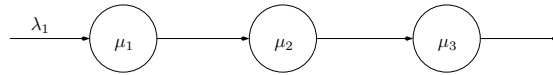


Figure 4: A three node tandem Jackson network

Consider the three-node tandem Jackson network with $\lambda_1 = 0.04$ and $\mu_1 = \mu_2 = \mu_3 = 0.32$. The simulation parameters are set again as $\delta = 0.1, k = 0.5$.

	$n = 20$	$n = 40$	$n = 60$
Theoretical value	1.26×10^{-15}	4.29×10^{-33}	8.32×10^{-51}
Estimate	1.23×10^{-15}	4.26×10^{-33}	8.60×10^{-51}
Std. Err.	0.02×10^{-15}	0.07×10^{-33}	0.30×10^{-51}
95% C.I.	$[1.19, 1.27] \times 10^{-15}$	$[4.12, 4.40] \times 10^{-33}$	$[8.00, 9.20] \times 10^{-51}$

Table 4. Three-node tandem Jackson network, total population overflow

	$n = 20$	$n = 40$	$n = 60$
Theoretical value	2.38×10^{-17}	2.06×10^{-35}	1.79×10^{-53}
Estimate	2.68×10^{-11}	2.01×10^{-22}	1.40×10^{-33}
Std. Err.	0.12×10^{-11}	0.14×10^{-22}	0.13×10^{-33}
95% C.I.	$[2.44, 2.92] \times 10^{-11}$	$[1.73, 2.29] \times 10^{-22}$	$[1.14, 1.66] \times 10^{-33}$

Table 5. Three-node tandem Jackson network, individual buffer overflow

	$n = 20$	$n = 40$	$n = 60$
Theoretical value	7.93×10^{-18}	6.88×10^{-36}	5.97×10^{-54}
Estimate	7.73×10^{-18}	6.68×10^{-36}	6.02×10^{-54}
Std. Err.	0.74×10^{-18}	0.22×10^{-36}	0.20×10^{-54}
95% C.I.	$[6.25, 9.21] \times 10^{-13}$	$[6.24, 7.12] \times 10^{-36}$	$[5.62, 6.42] \times 10^{-54}$

Table 6. Three-node tandem Jackson network, overflow of the first buffer

12.3 A four-node Jackson network

Consider the following four-node Jackson network with

$$\lambda = (0.05, 0.05, 0.05, 0), \mu = (0.15, 0.2, 0.15, 0.35), P_{12} = 0.2, P_{32} = 0.3.$$

Since we would like to have the theoretical values available as a benchmark and the recursive algorithm for theoretical values becomes exceedingly slow for big n , especially in the case of the first buffer overflow, we will present the numerical results for $n = 20, 25, 30$. The simulation parameters are set as $\delta = 0.3, k = 0.3$.

	$n = 20$	$n = 25$	$n = 30$
Theoretical value	7.08×10^{-6}	1.18×10^{-7}	1.82×10^{-9}
Estimate	7.10×10^{-6}	1.15×10^{-7}	1.81×10^{-9}
Std. Err.	0.10×10^{-6}	0.02×10^{-7}	0.05×10^{-9}
95% C.I.	$[6.90, 7.30] \times 10^{-6}$	$[1.11, 1.19] \times 10^{-7}$	$[1.71, 1.91] \times 10^{-9}$

Table 7. Four-node Jackson network, total population overflow

	$n = 20$	$n = 25$	$n = 30$
Theoretical value	2.21×10^{-7}	3.11×10^{-9}	4.44×10^{-11}
Estimate	2.34×10^{-7}	2.90×10^{-9}	4.56×10^{-11}
Std. Err.	0.17×10^{-7}	0.28×10^{-9}	0.64×10^{-11}
95% C.I.	$[2.00, 2.68] \times 10^{-7}$	$[2.34, 2.46] \times 10^{-9}$	$[3.28, 5.84] \times 10^{-11}$

Table 8. Four-node Jackson network, individual buffer overflow

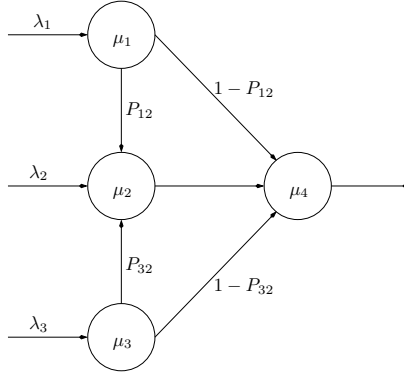


Figure 5: A four-node Jackson network

	$n = 20$	$n = 25$	$n = 30$
Theoretical value	8.03×10^{-10}	3.31×10^{-12}	1.36×10^{-14}
Estimate	7.54×10^{-10}	3.84×10^{-12}	1.49×10^{-14}
Std. Err.	0.40×10^{-10}	0.34×10^{-12}	0.17×10^{-14}
95% C.I.	$[6.74, 8.8.34] \times 10^{-10}$	$[3.20, 4.48] \times 10^{-12}$	$[1.15, 1.83] \times 10^{-14}$

Table 9. Four-node Jackson network, overflow of the first buffer

12.4 Remarks on numerical simulation

The choice of the parameters ε , δ and $\{C_F : F \in \mathcal{F}\}$ will effect the outcome of the numerical simulation, sometimes significantly, for a given n . For instance, in the example of the three-node tandem queue network, if we change $\delta = 0.1$ to $\delta = 0.15$ for $n = 20$, the numerical simulation results are better, especially for the problem of first buffer overflow.

Overflow type	Total population	Individual buffer	1st buffer
Theoretical value	1.26×10^{-15}	2.38×10^{-17}	7.93×10^{-18}
Estimate	1.26×10^{-15}	2.40×10^{-17}	7.94×10^{-18}
Std. Err.	0.01×10^{-15}	0.08×10^{-17}	0.14×10^{-18}
95% C.I.	$[1.24, 1.28] \times 10^{-15}$	$[2.24, 2.56] \times 10^{-17}$	$[7.66, 8.22] \times 10^{-18}$

Table 10. Three-node tandem Jackson network, $n = 20$, $\delta_n = 0.15$

Therefore, it is advisable to fine tune the parameters for fixed n to obtain better numerical performance.

A problem which becomes apparent in the numerical experimentation is that sometimes the constants $\{C_F : F \in \mathcal{F}\}$ may grow too fast. This phenomenon is of course more acute when the dimension is high, and for a given n has some impact on the choice of δ . If δ is too small, the subsolution will fail to catch the important “boundary layer” information present in the

prelimit. On the other hand, if δ is too big the reduction of the value of the subsolution at the origin will impact the optimality. Sometimes one may be able to tweak the choice of $\{C_F\}$ to alleviate this situation. Another possibility, as noted in the Introduction, is to forego the universal scheme in favor of a subsolution specifically designed for the problem at hand.

A Appendix. Proofs of Certain Results

A.1 Proof of Proposition 3.1

The stability condition implies that origin is positive recurrent for the state process Q . Let τ be the first return time to the origin and define $c \doteq E_0[\tau]$, the expected value of τ starting at the origin. Denote by σ the stationary distribution of the state process Q . Then for each $y = (y_1, \dots, y_d)' \in \mathbb{Z}_+^d$, it follows that [16]

$$\sigma(y) = \prod_{i=1}^d (1 - \rho_i) \rho_i^{y_i}.$$

Fix an arbitrary $x \in \Gamma^\circ$. There exists a sequence $\{x^n\} \subset \Gamma$ such that $x^n \rightarrow x$ and $nx^n \in \mathbb{Z}_+^d$. Let $B_n \doteq \{nx^n\}$. Clearly, $B_n \subset n\Gamma$. Define u_n analogously to p_n , with $n\Gamma$ replaced by B_n . It follows trivially that $u_n \leq p_n$. By a standard result on regenerative processes [3], we have

$$\sigma(B_n) = c^{-1} E_0[I_n], \quad I_n \doteq \int_0^\tau 1_{\{Q(t) \in B_n\}} dt.$$

Observe that $E_0[I_n] = u_n E_0[I_n | I_n > 0]$. Therefore

$$u_n = \frac{c}{E_0[I_n | I_n > 0]} \sigma(B_n).$$

However, given that B_n is reached, I_n is bounded from above by the time to reach the origin from B_n . It follows from [1] that the expected value of the time to reach origin from $B_n = \{nx^n\}$ is bounded from above by $\langle nx^n, 1 \rangle t^*$, where t^* is a constant that only depends on the system parameters. Thus

$E_0[I_n | I_n > 0] \leq n \langle x^n, 1 \rangle t^*$, and therefore

$$\begin{aligned}
\liminf_n \frac{1}{n} \log u_n &\geq \liminf_n \frac{1}{n} \log \frac{c}{n \langle x^n, 1 \rangle t^*} + \liminf_n \frac{1}{n} \log \sigma(B_n) \\
&= \liminf_n \frac{1}{n} \log \sigma(B_n) \\
&= \liminf_n \frac{1}{n} \log \prod_{i=1}^d (1 - \rho_i) \rho_i^{n x_i^n} \\
&= \langle r^*, x \rangle.
\end{aligned}$$

Since $p_n \geq u_n$ and $x \in \Gamma^\circ$ is arbitrary, we have

$$\liminf_n \frac{1}{n} \log p_n \geq - \inf_{x \in \Gamma^\circ} \langle -r^*, x \rangle.$$

For the other direction, we similarly define $D \doteq \{z \in \mathbb{R}_+^d : \langle -r^*, z \rangle \geq \inf_{x \in \Gamma} \langle -r^*, x \rangle\}$. Clearly $\Gamma \subset D$. Let v_n be the probability of reaching nD before returning to 0, starting from 0, and note that $p_n \leq v_n$. Recall that $-r^*$ is a vector with strictly positive components. Define $z^* \in \mathbb{Z}_+^d$ by

$$z_i^* \doteq \left\lfloor \frac{\inf_{x \in \Gamma} \langle -r^*, x \rangle}{-r_i^*} \right\rfloor + 1, \quad i = 1, \dots, d,$$

and

$$C_n \doteq (nD) \cap \{y \in \mathbb{Z}_+^d : y_i \leq n z_i^*, i = 1, \dots, d\}.$$

Even though $C_n \subset nD$, we claim that v_n equals the probability of reaching C_n before returning to 0, starting from 0. Indeed, consider a sample path that reaches nD before coming back to 0 and assume that it first hits nD at point y . If there is a component of y , say y_i , which satisfies $y_i > n z_i^*$, then there must exist a prior point on the sample path, say \bar{y} such that $\bar{y}_i = n z_i^*$. Since $-r^*$ has positive components,

$$\langle -r^*, \bar{y}/n \rangle \geq -r_i^* z_i^* \geq \inf_{x \in \Gamma} \langle -r^*, x \rangle,$$

or $\bar{y} \in nD$. This contradicts that y is the first hitting point, and therefore $y_i \leq n z_i^*$ for all i , or $y \in C_n$.

An argument analogous to one used in the first part of the proof yields

$$v_n = \frac{c}{E_0[\bar{I}_n | \bar{I}_n > 0]} \sigma(C_n), \quad \bar{I}_n \doteq \int_0^\tau 1_{\{Q(t) \in C_n\}} dt.$$

But $E_0[\bar{I}_n | \bar{I}_n > 0] \geq \min(\mu_i^{-1})$ since once the state process Q reaches C_n , at least one service time must elapse before Q comes back to the origin. It follows that

$$\begin{aligned} \limsup_n \frac{1}{n} \log v_n &\leq \limsup_n \frac{1}{n} \log \sigma(C_n) \\ &= \limsup_n \frac{1}{n} \log \sum_{y \in C_n} \prod_{i=1}^d (1 - \rho_i) \rho_i^{y_i}. \end{aligned}$$

But for every $y \in C_n$, we have $y \in nD$ and therefore

$$\prod_{i=1}^d (1 - \rho_i) \rho_i^{y_i} = \prod_{i=1}^d (1 - \rho_i) \cdot e^{\langle r^*, y \rangle} \leq \prod_{i=1}^d (1 - \rho_i) \cdot e^{-n \inf_{x \in \Gamma} \langle -r^*, x \rangle},$$

which in turn implies

$$\begin{aligned} \limsup_n \frac{1}{n} \log v_n &\leq - \inf_{x \in \Gamma} \langle -r^*, x \rangle + \limsup_n \frac{1}{n} \log \sum_{y \in C_n} 1 \\ &\leq - \inf_{x \in \Gamma} \langle -r^*, x \rangle + \limsup_n \frac{1}{n} \log \prod_{i=1}^d (nz_i^* + 1) \\ &= - \inf_{x \in \Gamma} \langle -r^*, x \rangle. \end{aligned}$$

Recalling $p_n \leq v_n$ completes the proof. ■

A.2 Proofs for Section 7

Proof of Lemma 7.8. Part 1 is trivial by the definition of \mathcal{F} and Remark 7.6. As for Part 2, let $E_1, E_2 \in \mathcal{F}$ and $E \doteq E_1 \cup E_2$. Define $\bar{z} \doteq \{\bar{z}_i : i \in E\}$ by

$$\bar{z}_i \doteq \begin{cases} z_i^{E_1} & \text{if } i \in E_1 \setminus (E_1 \cap E_2), \\ z_i^{E_1} \wedge z_i^{E_2} & \text{if } i \in E_1 \cap E_2, \\ z_i^{E_2} & \text{if } i \in E_2 \setminus (E_1 \cap E_2). \end{cases}$$

It suffices to show that $z^E \leq \bar{z}$. Thanks to Corollary 7.5, we only need to show that $\{\bar{z}_i : i \in E\}$ satisfies (7.5) with “=” replaced by “ \leq ”.

Fix an arbitrary $i \in E = E_1 \cup E_2$. Assume for now that $i \in E_1$. Since $z^{E_2} < 1$, the left-hand-side [LHS] of equation (7.5) with z replaced by \bar{z}

satisfies

$$\begin{aligned}
\text{LHS} &= \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ \text{some } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} \bar{z}_j \\
&\leq \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ \text{some } j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} \bar{z}_j \\
&\leq \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ \text{some } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} z_j^{E_1} \\
&= \mu_i z_i^{E_1}.
\end{aligned}$$

Similarly, for $i \in E_2$, $\text{LHS} \leq \mu_i z_i^{E_2}$. Thus $\text{LHS} \leq \mu_i \bar{z}_i$ for all $i \in E$, and therefore by Corollary 7.5 the proof of Part 1 is complete.

For Part 3, assume that $E, F \in \mathcal{F}$ and $E \subset F$. The proof of Part 2 implies, with $E_1 = E$ and $E_2 = F$, that $z^F \leq \bar{z}$ where

$$\bar{z}_i \doteq \begin{cases} z_i^E \wedge z_i^F & \text{if } i \in E, \\ z_i^F & \text{if } i \in F \setminus E. \end{cases}$$

In particular, for $i \in E$,

$$z_i^F \leq \bar{z}_i = z_i^E \wedge z_i^F,$$

or equivalently $z_i^F \leq z_i^E$. It follows immediately that $r^E \leq r^F$ by the definition (7.9) and because $z_i^F < 1$ for $i \in F \setminus E$. \blacksquare

A.3 Proofs for Section 8

We will prove a stronger version of Lemma 8.1, which will be useful in the proof of the main theorem. Fix an arbitrary $s \geq 1$. We define, for any $x \in \mathbb{R}_+^d$ and $\bar{\Theta}, \theta \in \mathcal{P}^+(\mathbb{V})$,

$$L_s(x, \alpha; \bar{\Theta}, \theta) \doteq \frac{s}{2} \langle \alpha, \mathbb{F}(x, \theta) \rangle + (s-1) \sum_{v \in \mathbb{V}} \theta[v] \log \frac{\bar{\Theta}[v]}{\Theta[v]} + R(\theta \| \bar{\Theta}).$$

and

$$\mathbb{H}_s(x, \alpha) = \sup_{\bar{\Theta} \in \mathcal{P}^+(\mathbb{V})} \inf_{\theta \in \mathcal{P}^+(\mathbb{V})} L_s(x, \alpha; \bar{\Theta}, \theta).$$

Note that \mathbb{H} is the special case of \mathbb{H}_s with $s = 2$. The interpretation of \mathbb{H}_s is that if one used the s -th moment instead of the second moment of \hat{p}_n as the criteria for optimality, \mathbb{H}_s instead of \mathbb{H} would be the associated Hamiltonian.

Lemma A.1. Consider any $E \subset \{1, \dots, d\}$ and $x \in \partial_E$. For every $\alpha \in \mathbb{R}_+^d$,

$$\mathbb{H}_s(x, \alpha) = -sH(x, -\alpha/2) = -sH_E(-\alpha/2).$$

The min/max condition holds in the definition of \mathbb{H}_s and the saddle point is $(\Theta_E^*(\alpha), \Theta_E^*(\alpha))$ as defined in Lemma 8.1. In particular, for every x , $\mathbb{H}_s(\cdot)$ is concave.

Proof. It follows from the definition of relative entropy that

$$\begin{aligned} L_s(x, \alpha; \bar{\Theta}, \theta) &= \frac{s}{2} \langle \alpha, \mathbb{F}(x, \theta) \rangle - (s-1) \sum_{v \in \mathbb{V}} \theta[v] \log \frac{\theta[v]}{\bar{\Theta}[v]} \\ &\quad + (s-1) \sum_{v \in \mathbb{V}} \theta[v] \log \frac{\theta[v]}{\Theta[v]} + R(\theta \| \Theta) \\ &= \frac{s}{2} \langle \alpha, \mathbb{F}(x, \theta) \rangle - (s-1) R(\theta \| \bar{\Theta}) + s R(\theta \| \Theta). \end{aligned}$$

Assume that $x \in \partial_E$ for some $E \subset \{1, \dots, d\}$. It follows from the definition of H_E (7.2) that Θ_E^* as defined in Lemma 8.1 is a probability distribution on \mathbb{V} . We now show that $(\Theta_E^*(\alpha), \Theta_E^*(\alpha))$ is a saddle point, or for any $\bar{\Theta}, \theta \in \mathcal{P}^+(\mathbb{V})$,

$$L_s(x, \alpha; \bar{\Theta}, \Theta_E^*(\alpha)) \leq L_s(x, \alpha; \Theta_E^*(\alpha), \Theta_E^*(\alpha)) \leq L_s(x, \alpha; \Theta_E^*(\alpha), \theta).$$

The first inequality is trivial due to that $s \geq 1$, the non-negativity of relative entropy, and that $R(\theta \| \bar{\theta}) = 0$ if and only if $\theta = \bar{\theta}$. As for the second inequality, straightforward calculation yields that

$$L_s(x, \alpha; \Theta_E^*(\alpha), \theta) = -(s-1)H_E(-\alpha/2) + \frac{1}{2} \sum_{v \in \mathbb{V}} \theta[v] \langle \alpha, \pi[x, v] \rangle + R(\theta \| \Theta).$$

By elementary calculus, the right-hand-side is minimized at $\theta = \Theta_E^*(\alpha)$. Therefore, the second inequality holds, and $(\Theta_E^*(\alpha), \Theta_E^*(\alpha))$ is a saddle point. In particular,

$$\mathbb{H}_s(x, \alpha) = L_s(x, \alpha; \Theta_E^*(\alpha), \Theta_E^*(\alpha)) = -sH_E(-\alpha/2).$$

This completes the proof. ■

A.4 Proofs for Section 9

Lemma A.2. *Let $E_1, E_2 \subset \{1, \dots, d\}$ be such that $E_1 \cap E_2 = \emptyset$. Let $E = E_1 \cup E_2$ and define the vector $\{\bar{z}_i : i \in E\}$ by*

$$\bar{z}_i \doteq \begin{cases} z_i^{E_1} & \text{if } i \in E_1, \\ 1 & \text{if } i \in E_2. \end{cases}$$

Suppose that $H_i(-r^{E_1}) \leq 0$ [resp., “ \geq ”] for all $i \in E_2$. Then

1. $\bar{z}_i \leq z_i^E$ [resp., “ \geq ”] for all $i \in E = E_1 \cup E_2$,
2. $1 < z_i^E$ [resp., “ $>$ ”] for all those $i \in E_2$ such that $H_i(-r^{E_1}) < 0$ [resp., “ $>$ ”].

Proof. We will argue the “ \leq ” case, the other case is analogous and thus omitted. For any $i \notin E_1$, thanks to Proposition 7.7 and the definitions of H_\emptyset and H_i , it follows that

$$e^{H_i(-r^{E_1})} - 1 = e^{H_i(-r^{E_1})} - e^{H_\emptyset(-r^{E_1})} = \sum_{v \in \mathbb{V}[-e_i]} \Theta[v] \left(1 - e^{-\langle r^{E_1}, v \rangle}\right).$$

For $i \notin E_1$ and $v \in \mathbb{V}[-e_i]$

$$e^{-\langle r^{E_1}, v \rangle} = \begin{cases} e^{-\langle r^*, v \rangle} & \text{if } v = -e_i \text{ or } v = -e_i + e_j, j \notin E_1 \\ e^{-\langle r^*, v \rangle} z_j^{E_1} & \text{if } v = -e_i + e_j, j \in E_1. \end{cases}$$

Thus by (7.11), for all $i \notin E_1$

$$e^{H_i(-r^{E_1})} - 1 = \mu_i - \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} - \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ \text{some } j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} z_j^{E_1}. \quad (\text{A.1})$$

Thus the assumption $H_i(-r^{E_1}) \leq 0$ for $i \in E_2$ amounts to

$$\sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ \text{some } j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} z_j^{E_1} \geq \mu_i. \quad (\text{A.2})$$

It follows that for any $i \in E = E_1 \cup E_2$,

$$\sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ \text{some } j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} \bar{z}_j \geq \mu_i \bar{z}_i.$$

Indeed, for $i \in E_2$, the above inequality is implied by (A.2) and for $i \in E_1$ the inequality is equality by the definition of z^{E_1} . Thanks to Corollary 7.5, we have

$$z_j^E \geq \bar{z}_j,$$

for all $j \in E$. This ends the proof of Part 1.

As for Part 2, we assume from now on that $H_i(-r^{E_1}) < 0$ for some $i \in E_2$. Thanks to Part 1, $z_j^E \geq \bar{z}_j$ for all $j \in E$. That is, $z_j^E \geq 1$ for all $j \in E_2$ and $z_j^E \geq z_j^{E_1}$ for all $j \in E_1$. Therefore,

$$\begin{aligned} \mu_i z_i^E &= \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ j \in E}} \Theta[v] e^{-\langle r^*, v \rangle} z_j^E \\ &\geq \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} + \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ j \in E_1}} \Theta[v] e^{-\langle r^*, v \rangle} z_j^{E_1}. \end{aligned}$$

Since $H_i(-r^{E_1}) < 0$, the inequality (A.2) is satisfied with \geq replaced by $>$. It follows that $\mu_i z_i^E > \mu_i$ or $z_i^E > 1$. \blacksquare

For any set $E \subset \{1, 2, \dots, d\}$, Part 2 of Lemma 7.8 implies that

$$\eta[E] \doteq \cup_{F \subset E, F \in \mathcal{F}} F \tag{A.3}$$

is an element of \mathcal{F} and it is the *maximal* subset of E that belongs to \mathcal{F} .

Lemma A.3. *Given any subset $E \subset \{1, \dots, d\}$ and $F \subset E$,*

$$H_F(-r^{\eta[E]}) \leq 0.$$

Proof. Thanks to Lemma 7.1 and that $H_\emptyset(r^{-\eta[E]}) = 0$, it suffices to argue that for all $i \in E$ that

$$H_i(-r^{\eta[E]}) \leq 0. \tag{A.4}$$

By Proposition 7.7, (A.4) holds with equality for $i \in \eta[E]$. Suppose for now that $i \in E \setminus \eta[E]$, and that (A.4) does not hold. By Lemma A.2 [with $E_1 = \eta[E]$ and $E_2 = \{i\}$], we have $\eta[E] \cup \{i\} \in \mathcal{F}$, which contradicts the maximality of $\eta[E]$. Therefore (A.4) holds. \blacksquare

Proof of Lemma 9.2. Assume $x = (x_1, \dots, x_d)' \in \partial_E$. Since $\mathbb{H}(x, 2r^G) = -2H_E(-r^G)$ by Lemma 8.1, it suffices to show that for any $F \subset E$,

$$H_F(-r^G) \leq 0. \tag{A.5}$$

We prove this stronger form since equation (A.5) is needed later on. Define a to be a small positive constant such that

$$a \leq \min \{C_F - C_G : G, F \in \mathcal{F}, G \subset F, G \neq F\}$$

and for all nonempty $F \in \mathcal{F}$

$$C_F \geq \max \left\{ \frac{\log(1/z_i^F)}{\log(1/z_i^G)} \cdot (C_G + a) : i \in G \subset F, G \neq F, G \in \mathcal{F} \right\}.$$

Note that such a positive constant a always exists.

For any $F \subset \{1, \dots, d\}$, the definitions (7.9) and (9.1), and that $x_i = 0$ for $i \in E$ imply that

$$\bar{W}_G^\delta(x) - \bar{W}_F^\delta(x) = 2 \sum_{i \in F \setminus E} x_i \log z_i^F - 2 \sum_{i \in G \setminus E} x_i \log z_i^G + (C_F - C_G) \delta.$$

By definition, $\bar{W}_F^\delta(x) \geq \bar{W}^\delta(x)$, and therefore by the assumption of the lemma

$$\bar{W}_G^\delta(x) - \bar{W}_F^\delta(x) < a\delta.$$

Let $\eta[E]$ be the *maximal* subset of E that belongs to \mathcal{F} as defined above Lemma A.3. Plugging in $F = \eta[E]$ and $F = \eta[G \cup E] \supset G$, we arrive at the following two inequalities.

$$-2 \sum_{i \in G \setminus E} x_i \log z_i^G + (C_{\eta[E]} - C_G) \delta < a\delta, \quad (\text{A.6})$$

$$\begin{aligned} 2 \sum_{i \in \eta[G \cup E] \setminus E} x_i \log z_i^{\eta[G \cup E]} - 2 \sum_{i \in G \setminus E} x_i \log z_i^G \\ + (C_{\eta[G \cup E]} - C_G) \delta < a\delta. \end{aligned} \quad (\text{A.7})$$

If $G \subset E$, then $G \subset \eta[E]$ by the maximality of $\eta[E]$. But now inequality (A.6) reduces to $C_{\eta[E]} - C_G < a$. However, since $C_F - C_G \geq a$ if G is a strict subset of F , $G = \eta[E]$ and by Lemma A.3 the proof is completed.

Assume from now on that $G \setminus E \neq \emptyset$. Note that the definition of η implies

$$G \cup E \supset \eta[G \cup E] \supset \eta[G] \cup \eta[E] = G \cup \eta[E].$$

It follows that

$$\eta[G \cup E] \setminus E = G \setminus E.$$

Therefore (A.7) yields that

$$C_{\eta[G \cup E]} < C_G + 2\delta^{-1} \sum_{i \in G \setminus E} x_i \left[\log(1/z_i^{\eta[G \cup E]}) - \log(1/z_i^G) \right] + a.$$

Let

$$c^* \doteq \max \left\{ \frac{\log(1/z_i^{\eta[G \cup E]})}{\log(1/z_i^G)} : i \in G \setminus E = \eta[G \cup E] \setminus E \right\}.$$

Note that $c^* \geq 1$. It follows from the definition of c^* and (A.6) that

$$\begin{aligned} C_{\eta[G \cup E]} &< C_G + (c^* - 1) \cdot 2\delta^{-1} \sum_{i \in G \setminus E} x_i \log(1/z_i^G) + a \\ &\leq C_G + (c^* - 1)(C_G - C_{\eta[E]} + a) + a \\ &\leq c^*(C_G + a). \end{aligned}$$

By the definition of a , if G is a strict subset of $\eta[G \cup E]$, then $C_{\eta[G \cup E]} \geq c^*(C_G + a)$, a contradiction. It follows that $\eta[G \cup E] = G$. Invoking Lemma A.3 we complete the proof. \blacksquare

A.5 Proof of Theorem 10.1

The proof of this main result is analogous to that of [7], which is based on a verification type of argument. Even though [7] only treats the special case of $b = 2$ for tandem queueing networks, the differences in the proofs are largely notational. For this reason we only give a sketch of the proof.

An important technicality of the proof is the validity of an exponential bound on the return time to origin. We use a different but much simpler approach to establish a slightly stronger bound [18]; compare the current Lemma A.4 with [7, Proposition A.1].

Recall the definition

$$T_0 \doteq \inf\{k \geq 1 : X^n(k) = 0\} = \inf\{k \geq 1 : Z(k) = 0\},$$

and that

$$Z(k+1) = Z(k) + \pi[Z(k), Y(k+1)],$$

where Y is a sequence of iid random variables taking values in \mathbb{V} with common distribution Θ . Denote by E_z the conditional expectation given that $Z(0) = z$ for each $z \in \mathbb{Z}_+^d$.

Lemma A.4. *There exists a constant $c > 0$ such that $E_z[\exp\{cT_0\}]$ is finite for all $z \in \mathbb{Z}_+^d$.*

Proof. Let $\sigma_0 \doteq \inf\{k \geq 0 : Z(k) = 0\}$. Note that $\sigma_0 = T_0$ if $Z(0) \neq 0$ and $\sigma_0 = 0$ if $Z(0) = 0$. For each $z \in \mathbb{Z}_+^d$, define $h(z) \doteq E_z[\sigma_0]$. Since the network is positive recurrent, $h(z)$ is finite for every z . Define the process

$$S(k) \doteq \begin{cases} h(Z(k)) & \text{if } k \leq \sigma_0, \\ \sigma_0 - k & \text{if } k > \sigma_0. \end{cases}$$

There are two key properties regarding S , which can be established using an analogous argument to [7, Lemmas A.2, A.4, and A.5]. We omit the proof.

(i) Let $\{\mathcal{F}_k = \sigma(Z(0), Y(1), \dots, Y(k))\}$ be the filtration. Then

$$E_z[S(k+1) - S(k) | \mathcal{F}_k] = -1$$

for all $z \in \mathbb{Z}_+^d$ and all $k \geq 0$.

(ii) The increments of the process S are *uniformly* bounded.

Note that there exists a $\varepsilon_0 > 0$ such that, for all $|x| \leq \varepsilon_0$,

$$e^x \leq 1 + x + x^2.$$

Let B be the uniform bound of the increments of S . Define

$$\bar{\varepsilon} \doteq \min \left\{ \frac{\varepsilon_0}{B}, \frac{1}{2B^2} \right\}.$$

Then for every $k \geq 0$, thanks to property (ii) and that $\bar{\varepsilon}B \leq \varepsilon_0$,

$$\begin{aligned} e^{\bar{\varepsilon}(S(k+1)-S(k))} &\leq 1 + \bar{\varepsilon}(S(k+1) - S(k)) + \bar{\varepsilon}^2(S(k+1) - S(k))^2 \\ &\leq 1 + \bar{\varepsilon}(S(k+1) - S(k)) + \bar{\varepsilon}^2 B^2. \end{aligned}$$

But by property (i) and that $\bar{\varepsilon}^2 B^2 \leq \bar{\varepsilon}/2$, it follows that

$$E \left[e^{\bar{\varepsilon}(S(k+1)-S(k))} \middle| \mathcal{F}_k \right] \leq 1 - \bar{\varepsilon}/2.$$

This is equivalent to that the process $\{(1 - \bar{\varepsilon}/2)^{-k} e^{\bar{\varepsilon}S(k)}, \mathcal{F}_k\}$ is a supermartingale. In particular, given $Z(0) = z \neq 0$, the Optional Sampling Theorem and $S(T_0) = S(\sigma_0) = h(0) = 0$ imply that ,

$$e^{\bar{\varepsilon}h(z)} = e^{\bar{\varepsilon}S(0)} \geq E_z \left[(1 - \bar{\varepsilon}/2)^{-T_0} \right].$$

Letting $c \doteq -\log(1 - \bar{\varepsilon}/2)$, we have

$$E_z \left[e^{cT_0} \right] < \infty$$

for all $z \in \mathbb{Z}_+^d$ and $z \neq 0$. For $Z(0) = 0$, we only need to note that

$$E_0 [e^{cT_0}] = E_0 [[e^{cT_0} | Z(1)]] = e^c \sum_{i=1}^d \mathbb{P}_0\{Z(1) = e_i\} E_{e_i} [e^{cT_0}] < \infty.$$

We complete the proof. ■

Proof of Theorem 10.1. To ease exposition, we write $W^n \doteq W^{\varepsilon_n, \delta_n}$, $\rho_F^n \doteq \rho_F^{\varepsilon_n, \delta_n}$, and set $\bar{\varepsilon}_n \doteq M \exp\{-a\delta_n/\varepsilon_n\}$ where M is given by Proposition 9.4. Fix an arbitrary real number $s \geq 1$. The value of s will be later determined, depending on n and b . Recall the definitions of L_s and \mathbb{H}_s in Section A.3. By concavity of $\log x$, Lemma 9.3, and the definition of $\bar{\Theta}^n$ in (10.2), we have

$$L_s(x, DW^n(x); \bar{\Theta}^n[\cdot|x], \theta) \geq \sum_{F \in \mathcal{F}} \rho_F^n(x) L_s(x, 2r_F; \bar{\Theta}_F[\cdot|x], \theta)$$

for any $\theta \in \mathcal{P}^+(\mathbb{V})$. In particular,

$$\inf_{\theta \in \mathcal{P}^+(\mathbb{V})} L_s(x, DW^n(x); \bar{\Theta}^n[\cdot|x], \theta) \geq \sum_{F \in \mathcal{F}} \rho_F^n(x) \inf_{\theta \in \mathcal{P}^+(\mathbb{V})} L_s(x, 2r_F; \bar{\Theta}_F[\cdot|x], \theta).$$

However, by definition (10.1) of $\bar{\Theta}_F[\cdot|x]$ and Lemma A.1,

$$\inf_{\theta \in \mathcal{P}^+(\mathbb{V})} L_s(x, 2r_F; \bar{\Theta}_F[\cdot|x], \theta) = \mathbb{H}_s(x, 2r_F) = \frac{s}{2} \mathbb{H}(x, 2r_F).$$

It follows from Proposition 9.4 that

$$\begin{aligned} \inf_{\theta \in \mathcal{P}^+(\mathbb{V})} L_s(x, DW^n(x); \bar{\Theta}^n[\cdot|x], \theta) &\geq \sum_{F \in \mathcal{F}} \rho_F^n(x) \cdot \frac{s}{2} \mathbb{H}(x, 2r_F) \\ &\geq -\frac{s}{2} \bar{\varepsilon}_n. \end{aligned} \tag{A.8}$$

For any $\theta \in \mathbb{V}$ and $x \in \mathbb{R}_+^d$, Taylor's expansion yields that

$$\begin{aligned} &n \sum_{v \in \mathbb{V}} \theta[v] \cdot \left[W^n \left(x + \frac{1}{n} \pi[x, v] \right) - W^n(x) \right] - \langle DW^n(x), \mathbb{F}(x, \theta) \rangle \\ &= n \sum_{v \in \mathbb{V}} \theta[v] \cdot \left[W^n \left(x + \frac{1}{n} \pi[x, v] \right) - W^n(x) - \left\langle DW^n(x), \frac{1}{n} \pi[x, v] \right\rangle \right] \\ &= \sum_{v \in \mathbb{V}} \theta[v] \cdot \frac{1}{2n} \pi[x, v]^T D^2 W(\bar{x}_v) \pi[x, v], \end{aligned}$$

where \bar{x}_v is some point on the line connecting x and $x + \pi[x, v]/n$. By straightforward computation, it is not difficult to show that each component of the Hessian matrix $D^2W^n(x)$ is uniformly bounded by C/ε_n for some constant C that only depends on the system parameters. This, the boundedness of $\pi[x, v]$, the assumption $s \geq 1$, and inequality (A.8) imply that

$$\begin{aligned} \inf_{\theta \in \mathcal{P}^+(\mathbb{V})} \left\{ \frac{s}{2} \sum_{v \in \mathbb{V}} \theta[v] \cdot n \left[W^n \left(x + \frac{1}{n} \pi[x, v] \right) - W^n(x) \right] \right. \\ \left. + (s-1) \sum_{v \in \mathbb{V}} \theta[v] \log \frac{\bar{\Theta}^n[v|x]}{\Theta[v]} + R(\theta||\Theta) \right\} \geq -\frac{s}{2} \left[\bar{\varepsilon}_n + \frac{\bar{C}}{n\varepsilon_n} \right], \end{aligned} \quad (\text{A.9})$$

for some constant \bar{C} that only depends on the system parameters.

Applying the relative entropy representation [7, Remark 3.1] to the left-hand-side of (A.9) and using the notation $\beta_n = \bar{\varepsilon}_n + \bar{C}/(n\varepsilon_n)$, we arrive at

$$e^{-s\beta_n/2} \cdot \sum_{v \in \mathbb{V}} e^{-sn[W^n(x+\pi[x,v]/n) - W^n(x)]/2} \left(\frac{\Theta[v]}{\bar{\Theta}^n[v|x]} \right)^{s-1} \cdot \Theta[v] \leq 1.$$

This inequality amounts to that the process $U = \{U(k) : k \geq 0\}$ where

$$U(k) \doteq e^{-ks\beta_n/2} e^{-snW^n(X^n(k))/2} \left(\prod_{j=0}^{k-1} \frac{\Theta[Y(j+1)]}{\bar{\Theta}^n[Y(j+1)|X^n(j)]} \right)^{s-1}$$

is a supermartingale under the original probability measure \mathbb{P} . In particular, by the Optional Sampling Theorem and the non-negativity of U ,

$$E^{\mathbb{P}} U(T_n \wedge T_0) \leq E^{\mathbb{P}} U(0) = e^{-snW^n(0)/2}.$$

Since $W^n(x) \leq 0$ for $x \in \Gamma$ and $\hat{p}_n = \hat{p}_n \cdot 1_{\{T_n < T_0\}}$, it follows that

$$e^{-snW^n(0)/2} \geq E^{\mathbb{P}} [U(T_n) 1_{\{T_n < T_0\}}] \geq E^{\mathbb{P}} \left[e^{-s\beta_n T_n/2} \hat{p}_n^{s-1} \right].$$

Assume from now on that $s > b$. Then by Hölder's inequality,

$$\begin{aligned} E^{\mathbb{P}} [\hat{p}_n^{b-1}] &\leq E^{\mathbb{P}} \left[e^{-\frac{s\beta_n T_n}{2}} \hat{p}_n^{s-1} \right]^{\frac{b-1}{s-1}} \cdot E^{\mathbb{P}} \left[e^{\frac{s\beta_n(b-1)T_n}{2(s-b)}} \cdot 1_{\{T_n < T_0\}} \right]^{\frac{s-b}{s-1}} \\ &\leq e^{-\frac{sn(b-1)}{2(s-1)} W^n(0)} \cdot E^{\mathbb{P}} \left[e^{\frac{s\beta_n(b-1)T_n}{2(s-b)}} \cdot 1_{\{T_n < T_0\}} \right]^{\frac{s-b}{s-1}}. \end{aligned}$$

Let c be the constant determined by Lemma A.4 and choose s such that

$$\frac{s\beta_n(b-1)}{2(s-b)} = c \quad \text{or} \quad s = b \cdot \frac{2c}{2c - \beta_n(b-1)}.$$

It follows that

$$\begin{aligned} \frac{1}{n} \log E^{\mathbb{P}}[\hat{p}_n^{b-1}] &\leq -\frac{s(b-1)}{2(s-1)} W^n(0) + \frac{s-b}{s-1} \frac{1}{n} \log E^{\mathbb{P}}[e^{cT_0}] \\ &= -\frac{bc}{2c + \beta_n} W^n(0) + \frac{b\beta_n}{2c + \beta_n} \frac{1}{n} \log E^{\mathbb{P}}[e^{cT_0}]. \end{aligned} \quad (\text{A.10})$$

Note that by Lemma 9.3 and definition (9.1),

$$2\gamma - \log |\mathcal{F}| \cdot \varepsilon_n - C_d \delta_n \leq W^n(0) \leq 2\gamma.$$

Letting $n \rightarrow \infty$ on both sides of (A.10), and observing that $\beta_n \rightarrow 0$, we have

$$\limsup \frac{1}{n} \log E^{\mathbb{P}}[\hat{p}_n^{b-1}] \leq -b\gamma.$$

But

$$E^{\mathbb{Q}_n}[\hat{p}_n^b] = E^{\mathbb{P}}[\hat{p}_n^{b-1}]$$

and by Jensen's inequality

$$\liminf \frac{1}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n^b] \geq \liminf \frac{b}{n} \log E^{\mathbb{Q}_n}[\hat{p}_n] = \liminf \frac{b}{n} \log p_n = -b\gamma.$$

This completes the proof. \blacksquare

A.6 Proofs of Section 11

Proof of Lemma 11.1. To ease notation, we write, for each $v \in \mathbb{V}$,

$$\bar{\Theta}[v] \doteq \Theta[v] \cdot e^{-\langle r^*, v \rangle}.$$

Fix an $i \in F$. Note that for $v \in \mathbb{V}$, $\langle v, e_i \rangle \neq 0$ only if $v \in \mathbb{V}[-e_i] \cup \mathbb{V}[e_i]$, and that $\mathbb{V}[e_i] \cap \mathbb{V}[-e_i] = \{-e_i + e_i\}$ with $\langle -e_i + e_i, e_i \rangle = 0$. It follows that

$$\sum_{v \in \mathbb{V}} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_i \rangle = \sum_{v \in \mathbb{V}[-e_i]} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_i \rangle + \sum_{v \in \mathbb{V}[e_i]} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_i \rangle.$$

But by the definition (7.9) of r^F and $\bar{\Theta}[v]$,

$$\sum_{v \in \mathbb{V}[e_i]} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_i \rangle = \sum_{\substack{v \in \mathbb{V}[e_i] \\ v \neq e_i - e_j \\ \text{any } j \in F}} \bar{\Theta}[v] z_i^F + \sum_{\substack{v \in \mathbb{V}[e_i] \\ v = e_i - e_j \\ \text{some } j \in F}} \bar{\Theta}[v] z_i^F / z_j^F,$$

and thanks to equation (7.5)

$$\begin{aligned}
\sum_{v \in \mathbb{V}[-e_i]} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_i \rangle &= - \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v \neq -e_i + e_j \\ \text{any } j \in F}} \bar{\Theta}[v] / z_i^F - \sum_{\substack{v \in \mathbb{V}[-e_i] \\ v = -e_i + e_j \\ \text{some } j \in F, j \neq i}} \bar{\Theta}[v] z_j^F / z_i^F \\
&= -\mu_i + \bar{\Theta}[-e_i + e_i].
\end{aligned} \tag{A.11}$$

Moreover, for $j \in F$ with $j \neq i$,

$$\sum_{v \in \mathbb{V}[-e_j]} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_i \rangle = \bar{\Theta}[-e_j + e_i] z_i^F / z_j^F.$$

Therefore, by combining (A.11) with the second term in (11.1), we can write

$$\begin{aligned}
\langle \beta_F^*, e_i \rangle &= \sum_{\substack{v \in \mathbb{V}[e_i] \\ v \neq e_i - e_j \\ \text{any } j \in F}} \bar{\Theta}[v] z_i^F + \sum_{\substack{v \in \mathbb{V}[e_i] \\ v = e_i - e_j \\ \text{some } j \in F, j \neq i}} \bar{\Theta}[v] z_i^F / z_j^F + z_i^F (-\mu_i + \bar{\Theta}[-e_i + e_i]) \\
&\quad - \sum_{j \in F, j \neq i} (1 - z_i^F) \bar{\Theta}[-e_j + e_i] z_i^F / z_j^F \\
&= \sum_{\substack{v \in \mathbb{V}[e_i] \\ v \neq e_i - e_j \\ \text{any } j \in F}} \bar{\Theta}[v] z_i^F + \sum_{\substack{v \in \mathbb{V}[e_i] \\ v = e_i - e_j \\ \text{some } j \in F, j \neq i}} \bar{\Theta}[v] z_i^F - \mu_i z_i^F + z_i^F \bar{\Theta}[-e_i + e_i] \\
&= z_i^F \left(-\mu_i + \sum_{v \in \mathbb{V}[e_i]} \bar{\Theta}[v] \right) \\
&= 0,
\end{aligned}$$

where the last equality follows from (7.7).

Now consider an arbitrary $j \notin F$. A similar argument leads to

$$\begin{aligned}
\langle -\beta_F^*, e_j \rangle &= - \sum_{v \in \mathbb{V}[e_j]} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_j \rangle - \sum_{v \in \mathbb{V}[-e_j]} \Theta[v] e^{-\langle r^F, v \rangle} \langle v, e_j \rangle \\
&\quad + \sum_{i \in F} (1 - z_i^F) \bar{\Theta}[-e_i + e_j] e^{-\langle r^F, -e_i + e_j \rangle} \\
&= - \sum_{v \in \mathbb{V}[e_j]} \bar{\Theta}[v] + \sum_{\substack{v \in \mathbb{V}[-e_j] \\ v = -e_j + e_i \\ \text{some } i \in F}} \bar{\Theta}[v] z_i^F + \sum_{\substack{v \in \mathbb{V}[-e_j] \\ v \neq -e_j + e_i \\ \text{any } i \in F}} \bar{\Theta}[v] \\
&= -\mu_j + \sum_{\substack{v \in \mathbb{V}[-e_j] \\ v = -e_j + e_i \\ \text{some } i \in F}} \bar{\Theta}[v] z_i^F + \sum_{\substack{v \in \mathbb{V}[-e_j] \\ v \neq -e_j + e_i \\ \text{any } i \in F}} \bar{\Theta}[v].
\end{aligned}$$

Thanks to the proof of Lemma A.2, more precisely, equation (A.1),

$$\langle -\beta_F^*, e_j \rangle = 1 - e^{H_j(-r^F)}$$

for $j \notin F$. In particular, if $j \in E \setminus F$, then Lemma A.2 implies that $H_j(-r^F) < 0$ or $\langle -\beta_F^*, e_j \rangle > 0$ since otherwise $F \cup \{j\} \in \bar{\mathcal{F}}$, which contradicts the maximality of $F = \Pi[E]$.

We now argue by contradiction that G is nonempty. Suppose that $G = \emptyset$. Then $H_j(-r^F) \leq 0$ for all $j \notin F$. By Lemma A.2 with $E_1 = F$ and $E_2 = \{1, \dots, d\} \setminus F$, we arrive at

$$z_j^{\{1, \dots, d\}} \geq 1$$

for all $j \notin F$. But this is impossible since $z_j^{\{1, \dots, d\}} = \rho_j < 1$ by Remark 7.6. Therefore G is non-empty. The claim that $F \cup G_1 \in \bar{\mathcal{F}}$ for all $G_1 \subset G$ is immediate from Lemma A.2 by setting $E_1 \doteq F$ and $E_2 \doteq G_1$.

It remains to show $L(x, \beta_F^*) = \bar{L}(x, \beta_F^*) = \langle -r^*, \beta_F^* \rangle$ for all $x \in \partial_F$. Since by definition (7.9),

$$\langle r^* - r^F, \beta_F^* \rangle = \sum_{i \in F} \log[z_i^F] \cdot \langle e_i, \beta_F^* \rangle = 0,$$

it suffices to show that $L(x, \beta_F^*) = \bar{L}(x, \beta_F^*) = \langle -r^F, \beta_F^* \rangle$. We first argue the easy direction that $L(x, \beta_F^*) \geq \langle -r^F, \beta_F^* \rangle$ and $\bar{L}(x, \beta_F^*) \geq \langle -r^F, \beta_F^* \rangle$. Consider any $\{(\rho_G, \beta_G) : G \subset F\}$ such that

$$\rho_G \geq 0, \quad \sum_{G \subset F} \rho_G = 1, \quad \sum_{G \subset F} \rho_G \beta_G = \beta_F^*. \quad (\text{A.12})$$

Note that for any $G \subset F$, $H_G(-r^F) = 0$ by Proposition 7.7. It now follows from the duality of L_G and H_G that

$$\begin{aligned} \sum_{G \subset F} \rho_G L_G(\beta_G) &\geq \sum_{G \subset F} \rho_G [\langle -r^F, \beta_G \rangle - H_G(-r^F)] & (\text{A.13}) \\ &= \sum_{G \subset F} \rho_G \langle -r^F, \beta_G \rangle \\ &= \langle -r^F, \beta_F^* \rangle. \end{aligned}$$

Taking infimum on the left-hand-side, we arrive at $L(x, \beta_F^*) \geq \langle -r^F, \beta_F^* \rangle$. The other inequality $\bar{L}(x, \beta_F^*) \geq \langle -r^F, \beta_F^* \rangle$ is shown in the same fashion, observing that $\bar{H}_G(-r^F) = 0$.

In order to show the other direction “ \leq ”, for every $G \subset F$, we denote by $\hat{\beta}_G$ the conjugate to $-r^F$ through the conjugacy of L_G and H_G , that is,

$$\hat{\beta}_G \doteq DH_G(-r^F). \quad (\text{A.14})$$

Note that $\hat{\beta}_G$ is also the conjugate of $-r^F$ through the conjugacy of \bar{L}_G and \bar{H}_G , namely,

$$\hat{\beta}_G = D\bar{H}_G(-r^F).$$

It suffices to find a collection of weights $\{\hat{\rho}_G\}$ such that $\{(\hat{\rho}_G, \hat{\beta}_G) : G \subset F\}$ satisfies the constraints (A.12). Indeed, if this is the case, the inequality in (A.13) becomes equality with $(\hat{\rho}_G, \hat{\beta}_G)$ in place of (ρ_G, β_G) , and therefore

$$L(x, \beta_F^*) \leq \sum_{G \subset F} \hat{\rho}_G L_G(\hat{\beta}_G) = \langle -r^F, \beta_F^* \rangle.$$

An analogous argument works for $\bar{L}(x, \beta_F^*)$.

For any $G \subset F$, we define

$$\hat{\rho}_G \doteq \prod_{i \in G} (1 - z_i^F) \cdot \prod_{i \in F \setminus G} z_i^F.$$

Then

$$\sum_{G \subset F} \hat{\rho}_G = \prod_{i \in F} [(1 - z_i^F) + z_i^F] = 1.$$

By the definitions in (A.14) and (7.2), it is not difficult to check that

$$\hat{\beta}_G = \sum_{v \in \mathbb{V}} \Theta[v] e^{-\langle r^F, v \rangle} v - \sum_{i \in G} \sum_{v \in \mathbb{V}[-e_i]} \Theta[v] e^{-\langle r^F, v \rangle} v.$$

On the other hand,

$$\begin{aligned} \sum_{G \subset F} \hat{\rho}_G \hat{\beta}_G &= \sum_{G \subset F} \hat{\rho}_G \sum_{v \in \mathbb{V}} \Theta[v] e^{-\langle r^F, v \rangle} v - \sum_{G \subset F} \hat{\rho}_G \sum_{i \in G} \sum_{v \in \mathbb{V}[-e_i]} \Theta[v] e^{-\langle r^F, v \rangle} v \\ &= \sum_{v \in \mathbb{V}} \Theta[v] e^{-\langle r^F, v \rangle} v - \sum_{i \in F} \sum_{\{G \subset F : i \in G\}} \hat{\rho}_G \sum_{v \in \mathbb{V}[-e_i]} \Theta[v] e^{-\langle r^F, v \rangle} v. \end{aligned}$$

However, for every $i \in F$,

$$\sum_{\{G \subset F : i \in G\}} \hat{\rho}_G = (1 - z_i^F) \cdot \prod_{j \in F \setminus \{i\}} [(1 - z_j^F) + z_j^F] = 1 - z_i^F.$$

Therefore

$$\sum_{G \subset F} \hat{\rho}_G \hat{\beta}_G = \beta_F^*,$$

and the proof is completed. ■

References

- [1] V. Anantharam. The optimal buffer allocation problem. *IEEE Trans. Inf. Theor.*, 35:721–725, 1989.
- [2] S. Asmussen. Conditioned limit theorems relating a random walk to its associates, with applications to risk reverse process and GI/G/1 queue. *Adv. Appl. Prob.*, pages 143–170, 1982.
- [3] S. Asmussen. *Applied Probability and Queues*. Springer-Verlag, New York, 2003.
- [4] P.J. De Boer. Analysis of state-independent importance sampling measures for the two-node tandem queue. *Preprint*, 2007.
- [5] P. Dupuis and R. S. Ellis. *A Weak Convergence Approach to the Theory of Large Deviations*. John Wiley & Sons, New York, 1997.
- [6] P. Dupuis, K. Leder, and H. Wang. Large deviations and importance sampling for a tandem network with slow-down. *QUESTA*, 57:71–83, 2007.
- [7] P. Dupuis, A. Sezer, and H. Wang. Dynamic importance sampling for queueing networks. *Ann. Appl. Prob.*, pages 1306–1346, 2007.
- [8] P. Dupuis and H. Wang. Importance sampling, large deviations, and differential games. *Stoch. and Stoch. Reports.*, 76:481–508, 2004.
- [9] P. Dupuis and H. Wang. Subsolutions of an Isaacs equation and efficient schemes for importance sampling. *Math. Oper. Res.*, 32:1–35, 2007.
- [10] M. I. Freidlin and A. D. Wentzell. *Random Perturbations of Dynamical Systems*. Springer-Verlag, New York, 1984.
- [11] P. Glasserman and S. Kou. Analysis of an importance sampling estimator for tandem queues. *ACM Trans. Modeling Comp. Simulation*, 4:22–42, 1995.
- [12] P. Glasserman and Y. Wang. Counter examples in importance sampling for large deviations probabilities. *Ann. Appl. Prob.*, 7:731–746, 1997.
- [13] P. Heidelberger. Fast simulation of rare events in queueing and reliability models. *ACM Trans. Modeling Comp. Simulation*, 4:43–85, 1995.

- [14] S. Juneja and V. Nicola. Efficient simulation of buffer overflow probabilities in Jackson networks with feedback. *ACM Transactions on Modelling and Computer Simulation*, 15:281–315, 2005.
- [15] S. Juneja and P. Shahabuddin. Rare-event simulation techniques: an introduction and recent advances. In Shane Henderson and Barry Nelson, editors, *Handbook on Simulation*, pages 291–350. Elsevier, 2006.
- [16] F.J. Kelly, editor. *Reversibility and Stochastic Networks*. Wiley, New York, 1979.
- [17] K. Majewski and K. Ramanan. How large queues build up in jackson networks. *Preprint*, 2008.
- [18] S.P. Meyn and D. Down. Stability of generalized Jackson networks. *Ann. Appl. Prob.*, 4:124–148, 1994.
- [19] S. Parekh and J. Walrand. A quick simulation method for excessive backlogs in networks of queues. *IEEE. Trans. Autom. Control*, 34:54–66, 1989.
- [20] J.S. Sadowsky. Large deviations and efficient simulation of excessive backlogs in a $GI/G/m$ queue. *IEEE. Trans. Automat. Control*, 36:1383–1394, 1991.