

BOTTLENECKS IN MARKOVIAN
QUEUEING NETWORKS

by

Phil Pollett

The University of Queensland

OUR SETTING

A closed network of queues:

- Fixed number of nodes (queues) J
- N customers circulating
- Usual Markovian assumptions in force

Examples:

- A job shop, where manufactured items are fashioned by various machines in turn.
- Provision of spare parts for a collection of machines.
- A mining operation, where coal faces are worked in turn by a number of specialized machines.

Can we identify regions of congestion (bottle-necks) from the parameters of the model?

BOTTLENECKS

Common sense:

The nodes with the smallest service effort will be the most congested.

A formal definition:

If n_j is the number of customers at node j , then this node is a *bottleneck* if, for all $m \geq 0$, $\Pr(n_j \geq m) \rightarrow 1$ as $N \rightarrow \infty$.

SIMPLE EXAMPLES

All nodes have infinitely many servers:

$\Pr(n_j = n) = \binom{N}{n} \alpha_j^n (1 - \alpha_j)^{N-n}$, $n = 0, \dots, N$, where α_j (< 1) is proportional to the arrival rate at node j divided by service rate. Clearly $\Pr(n_j = n) \rightarrow 0$ for each n as $N \rightarrow \infty$, and so *all nodes are bottlenecks*.

All nodes have a single server:

The distribution of n_j cannot be written down explicitly, but we can show that if there is a node j whose traffic intensity is *strictly greater* than the others, it is the unique bottleneck.

Moreover, for each node k in the remainder of the network, the distribution of n_k approaches a geometric distribution with parameter α_k/α_j in the limit as $N \rightarrow \infty$, and n_k , for $k \neq j$, are asymptotically *independent*.

MARKOVIAN NETWORKS

Our only assumption:

The steady-state (joint) distribution π of the numbers of customers $n = (n_1, n_2, \dots, n_J)$ at the various nodes has the product form

$$\pi(n) = B_N \prod_{j=1}^J \frac{\alpha_j^{n_j}}{\prod_{r=1}^{n_j} \phi_j(r)}, \quad n \in S,$$

where S is the finite subset of Z_+^J with $\sum_j n_j = N$ and B_N is a normalizing constant chosen so that π sums to unity over S .

Here α_j is proportional to the amount of service requirement (in items per minute) coming into node j (this will actually be *equal to* $\alpha_j B_N / B_{N-1}$). Suppose (wlog) that $\sum_j \alpha_j = 1$.

$\phi_j(n)$ is the service effort at node j (in items per minute) when there are n customers present. We shall assume that $\phi_j(0) = 0$ and $\phi_j(n) > 0$ whenever $n \geq 1$.

GENERATING FUNCTIONS

Our primary tool:

Define generating functions $\Phi_1, \Phi_2, \dots, \Phi_J$ by

$$\Phi_j(z) = 1 + \sum_{n=1}^{\infty} \frac{\alpha_j^n}{\prod_{r=1}^n \phi_j(r)} z^n.$$

It is easily shown that $B_N^{-1} = \langle \prod_{j=1}^J \Phi_j \rangle_N$, where $\langle \cdot \rangle_n$ takes the n^{th} coefficient of a power series. The marginal distribution of n_j can be evaluated as

$$\pi_j^{(N)}(n) = B_N \langle \Phi_j \rangle_n \langle \prod_{k \neq j} \Phi_k \rangle_{N-n},$$

for $n = 0, 1, \dots, N$.

SINGLE-SERVER NODES

Suppose that each node j has a single server ($\phi_j(n) = 1$ for $n \geq 1$). Then, $\langle \Phi_j \rangle_n = \alpha_j^n$ and so $\langle \Phi_j \rangle_{n+m} = \alpha_j^m \langle \Phi_j \rangle_n$. Summing

$$\pi_j^{(N)}(n) = B_N \langle \Phi_j \rangle_n \langle \prod_{k \neq j} \Phi_k \rangle_{N-n}$$

over n , and recalling that $B_N^{-1} = \langle \prod_{j=1}^J \Phi_j \rangle_N$, gives $Pr(n_j \geq m) = \alpha_j^m B_N / B_{N-m}$.

Suppose that $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_{J-1} < \alpha_J$, so that node J has maximal traffic intensity.

If we can prove that $B_{N-1}/B_N \rightarrow \alpha_J$ as $N \rightarrow \infty$, then $Pr(n_J \geq m) \rightarrow 1$ (node J is a bottleneck) and $Pr(n_j \geq m) \rightarrow (\alpha_j/\alpha_J)^m < 1$ for $j < J$ (the others are not).

WHY DOES $B_{N-1}/B_N \rightarrow \alpha_J$?

Define $\Theta_i = \Phi_1 \cdots \Phi_i$, where now $\Phi_j(z) = 1/(1 - \alpha_j z)$. Clearly Φ_j has radius of convergence (RC) $\rho_j = 1/\alpha_j$; in particular, $\Theta_1 (= \Phi_1)$ has RC $1/\alpha_1$.

Claim: Θ_i has RC $1/\alpha_i$ for all i , so that

$$\frac{B_N}{B_{N-1}} = \frac{\langle \Theta_J \rangle_{N-1}}{\langle \Theta_J \rangle_N} \rightarrow \frac{1}{\alpha_J}, \quad \text{as } N \rightarrow \infty.$$

Proof: Suppose Θ_k has RC $1/\alpha_k$ and consider

$$\begin{aligned} \langle \Theta_{k+1} \rangle_m &= \sum_{n=0}^m \alpha_{k+1}^{m-n} \langle \Theta_k \rangle_n \\ &= \alpha_{k+1}^m \sum_{n=0}^m \rho_{k+1}^n \langle \Theta_k \rangle_n. \end{aligned}$$

Clearly $\sum_{n=0}^{\infty} \rho_{k+1}^n \langle \Theta_k \rangle_n = \Theta_k(\rho_{k+1}) < \infty$, since $\rho_{k+1} < \rho_k$, and so

$$\frac{\langle \Theta_{k+1} \rangle_m}{\langle \Theta_{k+1} \rangle_{m+1}} \rightarrow \frac{1}{\alpha_{k+1}} \quad \text{as } m \rightarrow \infty,$$

implying that Θ_{k+1} has RC $1/\alpha_{k+1}$.

THE GENERAL CASE

Message: Bottleneck behaviour depends on the relative sizes of the radii of convergence of the power series $\Phi_1, \Phi_2, \dots, \Phi_J$.

Proposition 1: Suppose Φ_j has radius of convergence ρ_j and that $\rho_J < \rho_{J-1} \leq \rho_{J-2} \leq \dots \leq \rho_1$. Suppose also that

$$\frac{\langle \Phi_1 \cdots \Phi_{J-1} \rangle_{n-1}}{\langle \Phi_1 \cdots \Phi_{J-1} \rangle_n} \quad (1)$$

has a limit as $n \rightarrow \infty$. Then, node J is a bottleneck.

Example: Suppose node j has s_j servers, so that the traffic intensity at node j is proportional to α_j/s_j . Since $\phi_j(n) = \min\{n, s_j\}$, we have $\phi_j(n) \rightarrow s_j$, and so $\langle \Phi_j \rangle_{n-1} / \langle \Phi_j \rangle_n \rightarrow s_j / \alpha_j$. Therefore ρ_j is proportional to the reciprocal of the traffic intensity at node j . It can be shown that (1) holds.

COMPOUND BOTTLENECKS

What happens when the generating functions corresponding to two or more nodes *share* the same minimal RC?

Proposition 2: In the setup of Proposition 1, suppose that $\rho_L = \rho_{L+1} = \dots = \rho_J (= \rho)$ and that $\rho < \rho_j$ for $j = 1, 2, \dots, L-1$. Then, nodes $L, L+1, \dots, J$ behave *jointly* as a bottleneck in that $\Pr(\sum_{i=L}^J n_i \geq m) \rightarrow 1$ as $N \rightarrow \infty$.

It might be conjectured that when the generating functions corresponding to two nodes share the same minimal RC, they are always bottlenecks *individually*. However, while this is true when all nodes have a single server (because $\Pr(n_j \geq m) \rightarrow (\rho/\rho_j)^m$), it is *not true* in general.

SOME EXAMPLES

Consider a network with $J = 2$ nodes and suppose that $\alpha_1 = \alpha_2 = 1/2$. In the following examples Φ_1 and Φ_2 have the same RC $\rho = 2$.

Only one node is a bottleneck: Suppose that $\phi_1(n) = (n + 1)^2/n^2$ and $\phi_2(n) = 1$ for $n \geq 1$. Then, it can be shown that $\Pr(n_1 = n) \rightarrow 6/(\pi^2(n + 1)^2)$ and $\Pr(n_2 = n) \rightarrow 0$ as $N \rightarrow \infty$.

Neither node is a bottleneck: Suppose that $\phi_1(n) = \phi_2(n) = (n + 1)^2/n^2$ for $n \geq 1$. Then, $\Pr(n_1 = n) \rightarrow 3/(\pi^2(n + 1)^2)$ as $N \rightarrow \infty$.

AND FINALLY ...

Proposition 3: Suppose that $\Phi_1, \Phi_2, \dots, \Phi_K$ have the *same strictly minimal RC* ρ , and that $\phi_j(n)$ converges monotonically for some $j \in \{2, \dots, K\}$. Then, node 1 is a bottleneck *if and only if*

$$\Pr(n_1 \geq m \mid \sum_{i=1}^K n_i = N) \rightarrow 1 \text{ as } N \rightarrow \infty.$$

A *sufficient condition* for node 1 to be a bottleneck is that Φ_1 diverges at its RC and

$$\frac{\langle \Phi_2 \cdots \Phi_K \rangle_{n-1}}{\langle \Phi_2 \cdots \Phi_K \rangle_n} \text{ converges as } n \rightarrow \infty.$$

This latter condition is not necessary: In the setup of the previous examples, suppose that $\phi_1(n) = (n+1)^2/n^2$ and $\phi_2(n) = (n+1)^3/n^3$ for $n \geq 1$. Then, Φ_1 and Φ_2 have common RC $\rho = 2$ and both *converge* at their RC. But, it can be shown that $\Pr(n_1 = n)$ is bounded above by a quantity which is $O(N^{-1})$ as $N \rightarrow \infty$, implying that node 1 is a bottleneck.