Queueing Theory I

Summary

- Little's Law
- Queueing System Notation
- Stationary Analysis of Elementary Queueing Systems
 - □ M/M/1
 - \square M/M/m
 - □ M/M/1/K
 - □...





Generality of Little's Law

 $E[N] = \lambda E[T]$

- Little's Law is a pretty general result
- It does not depend on the arrival process distribution
- It does not depend on the service process distribution
- It does not depend on the number of servers and buffers in the system.



Specification of Queueing Systems

- Customer arrival and service stochastic models
- Structural Parameters
 - □ Number of servers
 - □ Storage capacity
- Operating policies
 - Customer class differentiation (are all customers treated the same or do some have priority over others?)
 - Scheduling/Queueing policies (which customer is served next)
 - Admission policies (which/when customers are admitted)

Queueing System Notation





Recall the Birth-Death Chain Example



M/M/1 Example

Meaning: Poisson Arrivals, exponentially distributed service times, one server and infinite capacity buffer.

• Using the birth-death result
$$\lambda_j = \lambda$$
 and $\mu_j = \mu$, we obtain

$$\pi_j = \left(\frac{\lambda}{\mu}\right)^j \pi_0, \quad j = 0, 1, 2, ...$$
• Therefore

$$\pi_0 \left(1 + \sum_{j=1}^{\infty} \left(\frac{\lambda}{\mu}\right)^j\right) = 1 \quad \text{for } \lambda/\mu = \rho < 1 \qquad \frac{\pi_0 = 1 - \rho}{\pi_j = (1 - \rho)\rho^j, \quad j = 1, 2, ...}$$

M/M/1 Performance Metrics

• Server Utilization $E[U] = \sum_{j=1}^{\infty} \pi_j = 1 - \pi_0 = 1 - (1 - \rho) = \rho$ • Throughput $E[R] = \mu \sum_{j=1}^{\infty} \pi_j = \mu (1 - \pi_0) = \mu \rho = \lambda$ • Expected Queue Length $E[X] = \sum_{j=0}^{\infty} j\pi_j = (1 - \rho) \sum_{j=0}^{\infty} j\rho^j = \rho (1 - \rho) \sum_{j=0}^{\infty} \frac{d\{\rho^j\}}{d\rho} =$ $= \rho (1 - \rho) \frac{d}{d\rho} \left\{ \sum_{j=0}^{\infty} \rho^j \right\} = \rho (1 - \rho) \frac{d}{d\rho} \left\{ \frac{1}{(1 - \rho)} \right\} = \frac{\rho}{(1 - \rho)}$

M/M/1 Performance Metrics

Average System Time $E[X] = \lambda E[S] \Rightarrow E[S] = \frac{1}{\lambda} E[X]$ $E[S] = \frac{1}{\lambda} \frac{\rho}{(1-\rho)} = \frac{1}{\mu(1-\rho)}$

Average waiting time in queue

$$E[S] = E[W] + E[Z] \Longrightarrow E[W] = E[S] - E[Z]$$
$$E[W] = \frac{1}{\mu(1-\rho)} - \frac{1}{\mu} = \frac{\rho}{\mu(1-\rho)}$$

M/M/1 Performance Metrics Examples



PASTA Property

- PASTA: Poisson Arrivals See Time Averages
- Let $\pi_i(t) = \Pr\{$ System state $X(t) = j \}$
- Let $a_i(t) = \Pr\{$ Arriving customer at t finds $X(t) = j \}$

• In general $\pi_i(t) \neq a_i(t)!$

Suppose a D/D/1 system with interarrival times equal to 1 and service times equal to 0.5



• Thus $\pi_0(t) = 0.5$ and $\pi_1(t) = 0.5$ while $a_0(t) = 1$ and $a_1(t) = 0!$

Theorem

For a queueing system, when the arrival process is Poisson and independent of the service process then, the probability that an arriving customer finds *j* customers in the system is equal to the probability that the system is at state *j*. In other words, $a_j(t) = \pi_j(t) \equiv \Pr\{X(t) = j\}, \quad j = 0, 1, ...$

Proof:

$$a_{j}(t) \equiv \lim_{\Delta t \to 0} \Pr \left\{ X(t) = j \mid a(t, t + \Delta t) \right\}$$

$$= \lim_{\Delta t \to 0} \frac{\Pr \left\{ X(t) = j, a(t, t + \Delta t) \right\}}{\Pr \left\{ a(t, t + \Delta t) \right\}}$$

$$= \lim_{\Delta t \to 0} \frac{\Pr \left\{ X(t) = j \right\} \Pr \left\{ a(t, t + \Delta t) \right\}}{\Pr \left\{ a(t, t + \Delta t) \right\}} = \Pr \left\{ X(t) = j \right\} = \pi_{j}(t)$$

M/M/m Queueing System

 Meaning: Poisson Arrivals, exponentially distributed service times, *m* identical servers and infinite capacity buffer.



M/M/m Queueing System

• Using the general birth-death result

$$\pi_{j} = \frac{1}{j!} \left(\frac{\lambda}{\mu}\right)^{j} \pi_{0}, \text{ if } j < m \quad \pi_{j} = \frac{m^{m}}{m!} \left(\frac{\lambda}{m\mu}\right)^{j} \pi_{0}, \text{ if } j \ge m$$
• Letting $\rho = \lambda/(m\mu)$ we get

$$\pi_{j} = \begin{cases} \frac{(m\rho)^{j}}{j!} \pi_{0} & \text{if } j < m \\ \frac{m^{m}\rho^{j}}{m!} \pi_{0} & \text{if } j \ge m \end{cases}$$
• To find π_{0}

$$\left(-\frac{m-1}{m!}(m\rho)^{j} - \frac{\infty}{m}m^{m}\rho^{j}\right) = \left(-\frac{m-1}{m!}(m\rho)^{j} - (m\rho)^{m}\right)^{m}$$

$$\pi_0 \left(1 + \sum_{j=1}^{m-1} \frac{(m\rho)^j}{j!} + \sum_{j=m}^{\infty} \frac{m^m \rho^j}{m!} \right) = 1 \implies \pi_0 = \left(1 + \sum_{j=1}^{m-1} \frac{(m\rho)^j}{j!} + \frac{(m\rho)^m}{m!(1-\rho)} \right)^{-1}$$

M/M/m Performance Metrics

• Server Utilization

$$E[U] = \sum_{j=1}^{m-1} j\pi_j + m \Pr\{X \ge m\} = \pi_0 \left(\sum_{j=1}^{m-1} j \frac{(m\rho)^j}{j!} + m \sum_{j=m}^{\infty} \frac{m^m \rho^j}{m!} \right)$$

$$= \pi_0 \left((m\rho) + \sum_{j=2}^{m-1} \frac{(m\rho)^j}{(j-1)!} + m \frac{(m\rho)^m}{m!(1-\rho)} \right)$$

$$= \pi_0 m\rho \left(1 + \sum_{j=2}^{m-1} \frac{(m\rho)^{j-1}}{(j-1)!} + \frac{(m\rho)^{m-1}}{(m-1)!} - \frac{(m\rho)^{m-1}}{(m-1)!} + \frac{m(m\rho)^{m-1}}{m!(1-\rho)} \right)$$

$$= \pi_0 m\rho \left(1 + \sum_{j=1}^{m-1} \frac{(m\rho)^j}{j!} + \frac{(m\rho)^m}{m!(1-\rho)} \right)$$

$$= \pi_0 m\rho \frac{1}{\pi_0} = m\rho = \frac{\lambda}{\mu}$$

M/M/m Performance Metrics

Throughput
$$E[R] = \mu \sum_{j=1}^{m-1} j\pi_j + m\mu \sum_{j=m}^{\infty} \pi_j = \lambda$$
Expected Queue Length
$$E[X] = \sum_{j=0}^{\infty} j\pi_j = \pi_0 \left(\sum_{j=1}^{m-1} j \frac{(m\rho)^j}{j!} + \frac{m^m}{m!} \sum_{j=m}^{\infty} j\rho^j \right) = \dots$$

$$E[X] = m\rho + \frac{(m\rho)^m}{m!} \frac{\rho}{(1-\rho)^2} \pi_0$$
Using Little's Law
$$E[S] = \frac{1}{\lambda} E[X] = \frac{1}{\lambda} \left(m\rho + \frac{(m\rho)^m}{m!} \frac{\rho}{(1-\rho)^2} \pi_0 \right)$$
Average Waiting time in queue
$$E[W] = E[S] - \frac{1}{\mu}$$

M/M/m Performance Metrics

Queueing Probability

$$P_{\mathcal{Q}} = \Pr\left\{X \ge m\right\} = \sum_{j=m}^{\infty} \pi_{j} = \pi_{0} \sum_{j=m}^{\infty} \frac{m^{m} \rho^{j}}{m!} = \frac{\pi_{0} (m\rho)^{m}}{m!(1-\rho)}$$

Erlang C Formula

Example

Suppose that packets arrive according to a Poisson process with rate λ =1. You are given the following two options,

- \Box Install a single transmitter with transmission capacity $\mu_1 = 1.5$
- $\hfill\square$ Install two identical transmitters with transmission capacity $\mu_2{=}~0.75$ and $\mu_3{=}~0.75$
- □ Split the incoming traffic to two queues each with probability 0.5 and have $\mu_2 = 0.75$ and $\mu_3 = 0.75$ transmit from each queue.



Example

Throughput

 It is easy to see that all three systems have the same throughput E[R_A]= E[R_B]= E[R_C]=λ

Server Utilization

$$E[U_{A}] = \frac{\lambda}{\mu_{1}} = \frac{1}{1.5} = \frac{2}{3}$$

$$E[U_{B}] = \frac{\lambda}{\mu_{2}} = \frac{1}{0.75} = \frac{4}{3}$$
 Therefore, each server is 2/3 utilized
$$E[U_{C}] = \frac{0.5\lambda}{\mu_{2}} = \frac{1}{2 \times 0.75} = \frac{2}{3}$$

Therefore, all transmitters are similarly loaded.

Example

Probability of being idle

$$\pi_{0A} = 1 - \frac{\lambda}{\mu_1} = \frac{1}{3}$$

$$\pi_{0B} = \left(1 + \sum_{j=1}^{m-1} \frac{(m\rho)^j}{j!} + \frac{(m\rho)^m}{m!(1-\rho)}\right)^{-1} = \left(1 + \frac{4}{3} + \frac{\left(\frac{4}{3}\right)^2}{2\left(1 - \frac{2}{3}\right)}\right)^{-1} = \frac{1}{5}$$

$$\pi_{_{0C}} = 1 - \frac{\lambda}{2\mu_2} = \frac{1}{3}$$
 For each transmitter

Example

• Queue length and delay

$$E[X_{A}] = \frac{\lambda}{\mu_{1} - \lambda} = \frac{1}{1.5 - 1} = 2 \qquad E[S_{A}] = \frac{1}{\lambda}E[X_{A}] = 2$$

$$E[X_{B}] = m\rho + \frac{(m\rho)^{m}}{m!} \frac{\rho}{(1 - \rho)^{2}} \pi_{0} = \frac{12}{5} \qquad E[S_{B}] = \frac{1}{\lambda}E[X_{B}] = \frac{12}{5}$$

$$E[X_{1C}] = \frac{\lambda/2}{\mu_{2} - \lambda/2} = \frac{0.5}{0.75 - 0.5} = 2 \qquad \text{For each queue!}$$

$$\Rightarrow E[X_{C}] = 2 \times E[X_{1C}] = 4 \qquad E[X_{C}] = \frac{1}{\lambda}E[X_{C}] = 4$$

M/M/∞ Queueing System



$M/M/\infty$ Performance Metrics

Expected Number in the System

$$E[X] = \sum_{j=0}^{\infty} j\pi_{j} = \sum_{j=0}^{\infty} j\frac{\rho^{j}}{j!}e^{-\rho} = \rho e^{-\rho} \sum_{j=1}^{\infty} \frac{\rho^{j-1}}{(j-1)!} = \rho$$

Using Little's Law

$$E[S] = \frac{1}{\lambda} E[X] = \frac{1}{\lambda} \frac{\lambda}{\mu} = \frac{1}{\mu}$$
 No queueing!



M/M/1/K Performance Metrics

Server Utilization

$$E[U] = 1 - \pi_0 = 1 - \frac{(1 - \rho)}{1 - \rho^{K+1}} = \frac{\rho(1 - \rho^K)}{1 - \rho^{K+1}}$$

Throughput

$$E[R] = \mu (1 - \pi_0) = \lambda \frac{1 - \rho^K}{1 - \rho^{K+1}} < \lambda$$

Blocking Probability

$$P_{B} = \pi_{K} = \frac{(1-\rho)\rho^{K}}{1-\rho^{K+1}}$$

Probability that an arriving customer finds the queue full (at state *K*)

M/M/1/K Performance Metrics



M/M/m/m – Queueing System

M/M/m/m Performance Metrics

Blocking Probability

$$P_{B} = \pi_{m} = \frac{\rho^{m} / m!}{\sum_{j=0}^{m} \frac{\rho^{j}}{j!}}$$
Erlang B Formula

Probability that an arriving customer finds all servers busy (at state *m*)

Throughput

$$E[R] = \lambda (1 - \pi_m) = \lambda \left(1 - \frac{\rho^m / m!}{\sum_{j=0}^m \frac{\rho^j}{j!}} \right) < \lambda$$

M/M/1//N – Closed Queueing System

 Meaning: Poisson Arrivals, exponentially distributed service times, one server and the number of customers are fixed to N.

• Using the birth-death result, we obtain $\pi_{j} = \frac{N!}{(N-j)!} \rho^{j} \pi_{0}, \quad j = 1, 2, ... N$ $\pi_{0} = \left[\sum_{j=0}^{N} \frac{N!}{(N-j)!} \rho^{j}\right]^{-1}$ $N\lambda \quad (N-1)\lambda \quad (N-2)\lambda \quad 2\lambda \qquad \lambda$ $0 \qquad \mu \qquad 1 \qquad \mu \qquad 2 \qquad \mu \qquad \mu \qquad N-1 \qquad \mu \qquad N$

