

# SFWR ENG 4E03

Kemal Ahmed

Fall 2015

Instructor: Vincent Maccio

Note: material covered in [Stats 3Y03 Summary](#) will not be covered in this summary. To find a unit CTRL-F “[<unit>]”, e.g. for Number of jobs in system, CTRL-F “[N]”

## Contents

|                                 |    |
|---------------------------------|----|
| Statistics .....                | 2  |
| Variance .....                  | 2  |
| Exponential .....               | 2  |
| Uniform .....                   | 2  |
| Binomial .....                  | 2  |
| Operations Analysis .....       | 3  |
| Visitation Trick .....          | 4  |
| Summation Equations .....       | 4  |
| DTMC .....                      | 5  |
| Balance Equations .....         | 5  |
| Matrices .....                  | 6  |
| CTMC .....                      | 6  |
| Poisson Process .....           | 6  |
| Kendall notation .....          | 6  |
| Variations .....                | 6  |
| M/M/1 .....                     | 7  |
| M/M/1/N .....                   | 7  |
| M/M/C .....                     | 7  |
| M/M/ $\infty$ .....             | 7  |
| Queuing .....                   | 7  |
| Square Root Staffing Rule ..... | 9  |
| e.g.) .....                     | 9  |
| Jackson Networks .....          | 9  |
| Traffic Equations .....         | 9  |
| Questions .....                 | 10 |

## Statistics

**Poisson parameter**  $[\lambda]$ : rate

**Service rate**  $[\mu]$ :

**Continuous Random Variable (CRV):**

$$\binom{n}{k} = \frac{n!}{k!(n-k)!}$$

*Think chemistry, i.e. cancelling units*

## Variance

$$\text{var}(X) = E[X^2] - (E[X])^2$$

- Don't change probability, but square X for calculation only

$$\begin{aligned}\text{var}(x) &= E[(X - \mu)^2] \\ &= \sigma^2 \\ &= \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx \\ &= \int x^2 f(x) dx - \mu^2\end{aligned}$$

**The higher your variance, the worse your system will perform.**

## Exponential

- **Mean**  $[E[X]]$ :  $1/\lambda$ 
  - a.k.a. Expected value
- **Variance**:  $1/\lambda^2$
- **Probability Distribution Function (PDF)**  $[P(X=x)]$ :  $\lambda e^{-\lambda x}/x!$
- **Cumulative Distribution Function (CDF)**  $[f(x)]$ :  $\text{CDF} = \int \text{PDF}$ , i.e.  $1 - e^{-\lambda x}$
- Memoryless
- not always for time

## Uniform

- **Variance**:  $(b-a)^2/12$
- **Mean**:  $(a+b)/2$
- **PDF**:  $1 / (b-a)$ ,  $a \leq x \leq b$
- **CDF**: 1
- **Uniform Distribution**: no memoryless property

## Binomial

- **Mean**  $[E[X]]$ :  $n \times \text{probability}$
- **Variance**:  $n \times p \times (1 - p)$
- **Probability Distribution Function (PDF)**  $[P(X)]$ :  $(n \text{ c } x)p^x(1-p)^{n-x}$

- **Cumulative Distribution Function (CDF) [f(x)]:**  $\sum_{i=0}^{\lfloor x \rfloor} \binom{n}{i} p_i (1-p)^{n-i}$

## Operations Analysis

**Device [i]:** units that are in terms of  $i$  are specific to an individual device or node within a system

**Total devices [k]:**

**Service Time [S]:** time per specific job

$1/\mu$

**Visitation [V]:** given or projected visits/jobs (closed system); cannot be calculated; basically a probability

$E(V)$ : calculated visit/job ratio

$P(\text{visit}) \cdot \text{total visits in previous node}$

**Demand [D]:** total service time for all jobs

$$D_i = E[S_i] \cdot V_i$$

$$D = \sum_{i=0}^k D_i$$

**Bottleneck [D<sub>max</sub>]:** device with largest demand, utilization

**Time in system [T]:** time the job is in the system

$$E[T] = \frac{N}{X}$$

$$E[T] \geq \max(D, ND_{\max} - E[Z])$$

If  $E[Z] = 0$ ,  $T = R$

**Response Time [R]:** time the job is *being processed* in the system

If  $E[Z] = 0$ ,  $R = T$

M/M/1:  $E[R] = 1/(\mu - \lambda)$

M/M/1/N:  $E[R] = E[N]/\lambda$

M/M/C:  $E[R] = E[R_q] + E[S]$

**Users [M]:**

**Optimal users [M\*]:**

$$M^* = \frac{D + E[Z]}{D_{\text{bottleneck}}}$$

**Total Jobs [N]:**  $N=M$  in a closed system

- *Little's Law:*  $E[N] = \lambda E[T]$ ,  $\lambda = X$
- $E[N] = \lambda E[R]$ ,  $\lambda = X$
- M/M/1:
  - $E[N] = \lambda/(\mu - \lambda) = \rho/(1 - \rho)$ , if you have overall system  $\lambda$

- $E[N] = \sum_{i=0}^{\infty} i\pi_i \leftarrow \text{probability} \times \text{\#jobs, if your } \lambda \text{ or } \mu \text{ is different for each state}$
- M/M/1/N:  $E[N] = \sum_{i=0}^N i\pi_i = \pi_0 \frac{\lambda}{\mu} \left( \frac{N \left(\frac{\lambda}{\mu}\right)^{N-1} - (N+1) \left(\frac{\lambda}{\mu}\right)^N + 1}{1 - \left(\frac{\lambda}{\mu}\right)^2} \right)$
- M/M/C: go through Little's law
  - $E[N] = E[N_0] + \rho$
- M/M/ $\infty$ :
- Jackson Network:  $E[N] = \sum E[N_i] = \sum \rho \lambda / (\mu_i - \rho \lambda) = \sum (\lambda_i / (\mu_i - \lambda_i))$

**Think time [Z]:** time it takes the user to put a request in and start, it's kinda like the frequency that users put in requests (seconds / request)

$$E[Z] = E[T] - E[R]$$

**Throughput [X]:** out-rate, max jobs / hour of full system

$$X \leq \min \left( \frac{1}{D_{\max}}, \frac{N}{D + E[Z]} \right)$$

Note:  $\frac{1}{D_{\max}}$  and  $\frac{N}{D + E[Z]}$  converge at their lowest point, so equate them

$$X_i = E[V_i] X$$

**Utilization [ $\rho$ ]:** ratio that the time is busy

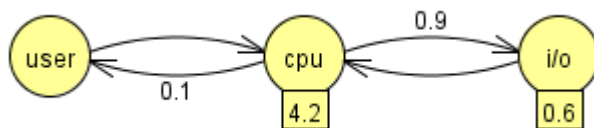
$$\rho_i = X_i E[S_i]$$

$$\rho_i = X D_i$$

$$\rho = \lambda / c_i \mu$$

Visitation Trick

If determining visitation at a node, establish a reference node from one of the incoming nodes, usually the user node, that has a returning percentage



$$V_{\text{user}} = 1 = 0.1 \cdot V_{\text{CPU}}$$

Summation Equations

**Geometric Series:**  $\sum_{i=0}^{\infty} r^i = \frac{1}{1-r}$ , where  $0 \leq r \leq 1$  (because otherwise it would be unstable)

$$\sum_{i=0}^{\infty} r^{i+1} = \sum_{i=1}^{\infty} r^i = \frac{r}{1-r}$$

$$\sum_{i=1}^{\infty} r^i = \sum_{i=0}^{\infty} r^i - 1$$

**Geometric Sequence:**  $S_n = \sum_{i=0}^n r^i = \frac{1-r^{n+1}}{1-r}$

$$S_n = \sum_{i=1}^n r^i = \frac{r(1-r^n)}{1-r}$$

Removing the annoying factors:

$$\sum_{i=0}^{\infty} i(i+1)\rho^i$$

Take out a value so the integral takes out the i and i+1

$$= \rho \sum_{i=0}^{\infty} i(i+1)\rho^{i-1}$$

$$= \rho \frac{d\rho}{di} \left( \sum_{i=0}^{\infty} (i+1)\rho^i \right)$$

$$= \rho \frac{d\rho^2}{d^2i} \left( \sum_{i=0}^{\infty} \rho^{i+1} \right)$$

$$= \rho \frac{d\rho^2}{d^2i} \left( \sum_{i=1}^{\infty} \rho^i \right)$$

$$= \rho \frac{d\rho^2}{d^2i} \left( \sum_{i=0}^{\infty} \rho^i - \rho^0 \right)$$

$$= \rho \frac{d\rho^2}{d^2i} \left( \frac{1}{1-\rho} - 1 \right)$$

## DTMC

**Discrete Time Markov Chains (DTMC):**

[n]: number of tasks in queue / system

Steady state:  $n \rightarrow \infty$

For discrete: use the sum of the X's, so  $E[X] = \sum (P(X=i) \cdot X_i)$  and  $E[X^2] = \sum (P(X=i) \cdot X_i^2)$

## Balance Equations

$$\pi_n = \frac{\prod_{i=0}^n \lambda_i}{\prod_{i=0}^n \mu_i} \pi_0$$

^think of it like series / parallel, where you add multiple connections out in different directions (parallel) and multiply connections stacked onto each other (series)

OR  $jobs_{in} = jobs_{out}$

## Matrices

Rows: equations for nodes going out (add up to 1)

Columns: equations for nodes coming in

## CTMC

### Poisson Process

**Counting Process:** a way of determining the time between consecutive occurrences of an event

**Poisson Process:** a *counting process*, whose time between arrivals uses Exponential Distribution

- $\lambda_{\text{total}} = \sum \lambda_i$ 
  - you can also split up  $\lambda$  into multiple  $\lambda$ s
- Not only do you see each second as time independent, each stream of probabilities is independent
- $P(x; \lambda) = e^{-\lambda} \lambda^x / x!$ 
  - $[x]$ : things will happen
  - $[\lambda]$ : rate;  $\lambda = \alpha t$
- $[\alpha]$ : expected number of events during unit interval
- $[t]$ : time interval length
- $P_x(t) = \frac{e^{-\alpha t} \cdot (\alpha t)^x}{X!}$

### Kendall notation

**Job Processing time**  $[\mu]$ : rate of jobs leaving system (jobs/sec)

$\mu = 1 / \text{processing\_time\_per\_job}$

M/M/1 Queue

$[M]$ : time between arrivals is Markovian (Memoryless)  $\sim \exp(\lambda)$

$[M]$ : job processing times are Markovian (Memoryless)  $\sim \exp(\mu)$

$[1]$ : single server

$$(\sum p_{\text{out}}) \times \pi_i = \sum p_j \pi_j, j=0..n, j \neq i$$

$\pi_0$ : percent of time that the queue is empty

Attributes:

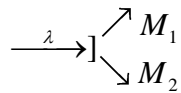
- FIFO
- Infinite buffer

### Variations

- M/M/2 Queue: same, except 2 servers
- M/M/C Queue: C servers
- M/E<sub>k</sub>/C: Erlang k, i.e. series of exponential
- H()/M/C: hyperexponential distribution
- PH/M/C: phase type, i.e. any combination of any number of exponentials with any rate
- M/G/C: General distribution
- G/G/1: has not been solved yet
- M/M/1/1: 1 server, 1 job

[c]: number of servers

Think: one queue goes to multiple servers



M/M/1

$$\pi_0 = 1 - \lambda/\mu$$

M/M/1/N

When you can only have up to N jobs in system queue.

[λ']: rate jobs enter the system, until the queue is full

$$\lambda' = \lambda(1 - \pi_N)$$

$$\pi_0 = \frac{1 - \frac{\lambda}{\mu}}{1 - \left(\frac{\lambda}{\mu}\right)^N} = \frac{1}{1 + \sum_{i=1}^N \left(\frac{\lambda}{\mu}\right)^i}$$

$$\pi_i = \left(\frac{\lambda}{\mu}\right)^i$$

**Waiting:** jobs put into the queue

**Blocked:** jobs not allowed in the queue

M/M/C

Useful if multiple jobs are sharing the same queue

Does the μ you use for equations double in M/M/2? No, but you'll see jobs coming out of a system at a rate of c·μ.

$$\pi_0 = \left[ 1 + \sum_{i=1}^{c-1} \frac{1}{i!} \left(\frac{\lambda}{\mu}\right)^i + \frac{1}{c!} \left(\frac{\lambda}{\mu}\right)^c \left(\frac{1}{1-\rho}\right) \right]^{-1}$$
$$\pi_i = \begin{cases} \frac{1}{n!} \left(\frac{\lambda}{\mu}\right)^n \pi_0, & n < c \\ \frac{1}{c! c^{n-c}} \left(\frac{\lambda}{\mu}\right)^n \pi_0, & n \geq c \end{cases}$$

M/M/∞

Same as M/M/C, except:

$$\pi_0 = e^{-\frac{\lambda}{\mu}}$$

and just find the unit

Queuing

**Blocking Probability** [P<sub>Q</sub>]: probability that a process will be blocked when entering the system and be placed in the queue

Erlang-C Equation:  $P_Q = \sum_{i=0}^{\infty} \pi_i = \frac{1}{c!} \left( \frac{\lambda}{\mu} \right)^c \left( \frac{1}{1-\rho} \right) \pi_0$

Given  $\lambda$  and  $\mu$ , what should  $c$  be so  $P_Q < \rho$

**Waiting time in queue**  $[R_Q]$ : response time of queue

$$E[R_Q] = \frac{1}{\lambda} P_Q \left( \frac{\rho}{1-\rho} \right)$$

$$M/M/1: E[R_Q] = \frac{1}{\mu - \lambda} - \frac{1}{\mu}$$

$$M/M/C: E[R_Q] = \left( \frac{(\lambda / \mu)^c \mu}{(c-1)!(c\mu - \lambda)^2} \right) \pi_0$$

$$M/M/\infty: E[R_Q] = 0$$

**Number of jobs in queue**  $[N_Q]$ :

$$M/M/1: \rho^2 / (1 - \rho)$$

$$M/M/C: E[N_Q] = \pi_0 \frac{\lambda \mu \rho^{c+1}}{(c-1)!(c\mu - \lambda)^2}$$

$$M/M/\infty: E[N_Q] = 0$$

You need to know what is in the progression of each step

*e.g.*

When you have varying

$$\pi_n = (n+1) \left( \frac{\lambda}{\mu} \right)^n \pi_0$$

$$1 = \sum_{i=0}^{\infty} (i+1) \rho^i \pi_0$$

$$= \pi_0 \sum_{i=1}^{\infty} (i+1) \rho^i$$

$$= \pi_0 \frac{d}{d\rho} \left( \sum_{i=0}^{\infty} \rho^{i+1} \right)$$

$$= \pi_0 \frac{d}{d\rho} \left( \sum_{i=1}^{\infty} \rho^i \right)$$

$$= \pi_0 \frac{d}{d\rho} \left( \sum_{i=0}^{\infty} \rho^{i+1} \right)$$



## Square Root Staffing Rule

Given an M/M/c queue with arrival rate,  $\lambda$ , server speed,  $\mu$ , and  $\rho$  is *large* (assume this means over 100, but we don't actually know what it means),  $\alpha$  is a bound on  $P_Q$ , let  $c_\alpha^*$  denote the least # of servers needed to ensure that  $P_Q < \alpha$ . Then

$c_\alpha^* \approx \rho + k\sqrt{\rho}$ , where  $k$  is the solution to

$$\frac{k\Phi(k)}{\phi(k)} = \frac{1-\alpha}{\alpha}, \text{ where } \Phi(\cdot) \text{ is the CDF of the standard normal and } \phi(\cdot) \text{ is its pdf}$$

[K]: minimum # servers to stay stable  $\lambda/\mu$  or  $\rho$

[k]: a constant...just assume 1 for now

e.g.)

| $\alpha$ | k     | $\rho + k\sqrt{\rho}$ |
|----------|-------|-----------------------|
| 0.8      | 0.178 | 10, 018               |
| 0.5      | 0.506 | 10, 051               |
| 0.2      | 1.06  | 10, 106               |
| 0.1      | 1.42  | 10, 142               |

[Q]: transition matrix

$$q_{ii} = -\sum_{j \neq i} q_{ij}$$

$$q_{ij} = \lim_{\Delta t \rightarrow 0} \frac{P\{X_{t+\Delta t} = j \mid X_t = i\} - \delta_{ij}}{\Delta t}$$

Replace  $i \leftrightarrow j$  to get  $q_{ji}$  and  $q_{ji}$ .

## Jackson Networks

$$P(N_1 = n_1) =$$

$$\pi_n = P_1^{n_1} (1 - \rho_1) \rho_2^{n_2} (1 - \rho_2) \cdots \rho_k^{n_k} (1 - \rho_k)$$

$$\pi_{n^*} = P(\text{state of system } n^*) = \prod_{i=1}^k P(n \text{ jobs at node } i)$$

$$= \prod_{i=1}^k \rho_i^{n_i} (1 - \rho_i)$$

**Poisson Arrivals See Time Averages property (PASTA):** the probability of the state (i.e.  $\pi_i$ ) as seen by an outside random observer is the same as the probability of the state seen by an arriving customer

$$\lambda_{\text{total}} = \sum \lambda_{\text{in},i}$$

## Traffic Equations

For each node, what is the number of jobs entering?

$$\lambda_x = R + \sum P_{i,\text{entering}} \cdot \lambda_{i,\text{entering}}$$

response rate + probability of each job entering

## Questions

- 275