

~~Chapter 12~~
~~Question 8~~

Tutorial 8, Question 2

To find the service demands of an HTTP request at the CPU:

The throughput of the server is

$$= \frac{10800}{3600}$$

$$= 3 \text{ HTTP requests/s}$$

By service demand law

service demand of CPU

$$= \frac{\text{utilisation of CPU}}{\text{server throughput}}$$

$$= \frac{0.3}{3}$$

$$= 0.1 \text{ s}$$

To find the throughput:

The easiest is to write a computer program
and plug the values in.

$$X(0) = 0$$

$$X_0(1) = 6.25 \text{ requests/s}$$

$$X_0(2) = 8.16 \text{ requests/s}$$

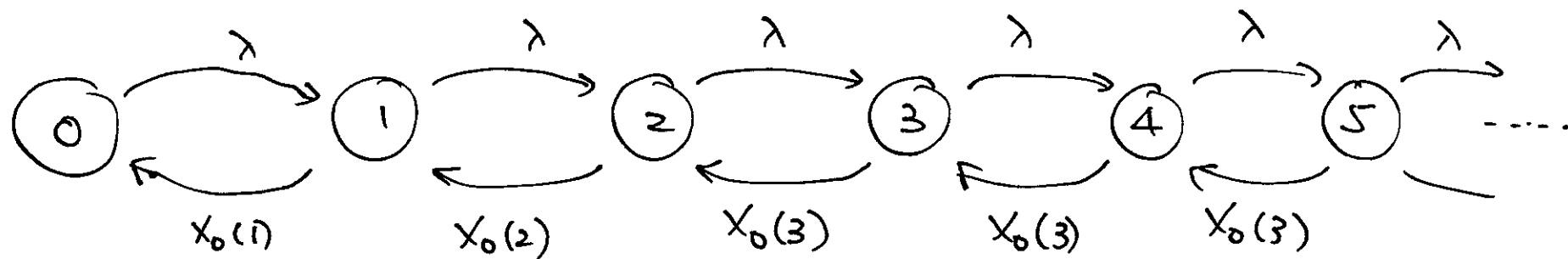
$$X_0(3) = 9.01 \text{ requests/s}$$

Part 3 of the question: To find the average response time of HTTP requests when $\lambda=5$ and the server can only process up to 3 requests at a time:

This can be modelled as a ~~generalised~~ Markov chain ~~birth-death~~ model.

Let state k ($k=0, 1, 2, \dots$) be the number of requests in the web server. Note that the number of requests in the web server includes those that are being served (up to three) and those that are in the processing queue.

The state space diagram is in the following page:



* The transition rate from state k to state $(k+1)$ (for $k=0,1,\dots$) is the arrival rate of the request.

* The transition rate from state $(k+1)$ to state k (for $k=0,\dots$) is the rate at which requests are completed.

- For state 1 to state 0, this is the same as the throughput of the web server when there is only one client. (Note that throughput is effectively the number of requests completed in an unit time.)

- For state 2 to state 1, the request completion rate is $\mu_0(2)$.

- For state 3 to state 2, the request completion rate is $X_0(3)$
- For state $(k+1)$ to state k (where $k \geq 3$), the request completion rate is always $X_0(3)$ because only 3 requests are being processed by the server. The others, _{requests} are waiting in the queue.

In order to find the response time, we need to solve the model.
Using the trick given in the notes, we know that

$$P(1) X_0(1) = \lambda P(0)$$

$$P(2) X_0(2) = \lambda P(1)$$

$$P(3) X_0(3) = \lambda P(2)$$

$$P(4) X_0(3) = \lambda P(3)$$

$$P(5) X_0(3) = \lambda P(4)$$

...

Expressing $P(1)$, $P(2)$, ... in terms of $P(0)$, we have

$$P(1) = \frac{\lambda}{X_0(1)} P(0)$$

$$P(2) = \frac{\lambda}{X_0(2)} \frac{\lambda}{X_0(1)} P(0)$$

$$P(3) = \frac{\lambda}{X_0(3)} \frac{\lambda}{X_0(2)} \frac{\lambda}{X_0(1)} P(0)$$

$$P(4) = \left(\frac{\lambda}{X_0(3)} \right)^2 \frac{\lambda}{X_0(2)} \frac{\lambda}{X_0(1)} P(0)$$

$$P(5) = \left(\frac{\lambda}{X_0(3)} \right)^3 \frac{\lambda}{X_0(2)} \frac{\lambda}{X_0(1)} P(0)$$

Observing the pattern, we've

$$P(k) = \left(\frac{\lambda}{X_0(3)} \right)^{k-2} \frac{\lambda}{X_0(2)} \frac{\lambda}{X_0(1)} P(0)$$

for $k \geq 3$

Define

$$P_1 = \frac{\lambda}{X_0(1)}$$

$$P_2 = \frac{\lambda}{X_0(2)}$$

$$P_3 = \frac{\lambda}{X_0(3)}$$

we have

$$P(1) = p_1 P(0)$$

$$P(2) = p_2 p_1 P(0)$$

$$P(k) = p_3^{k-2} p_2 p_1 P(0) \text{ for } k \geq 3$$

Since the sum of all probabilities must be 1,

$$P(0) + P(1) + \dots + \dots = 1$$

$$\cancel{P(0)} + p_1 P(0) + p_2 p_1 P(0) + p_3 p_2 p_1 P(0) + \dots = 1$$

$$\underbrace{p_3^2 p_2 p_1 P(0) + \dots}_{\text{(geometric progression)}} = 1$$

(Note: you can use the form given in the question)

$$\Leftrightarrow P(0) + p_1 P(0) + \frac{p_2 p_1}{1 - p_3} P(0) = 1$$

(Note that $p_3 < 1$, so the geometric progression converges)

$$\Rightarrow P(0) = \frac{1}{1 + p_1 + \frac{p_2 p_1}{1 - p_3}}$$

In order to calculate the response time, we need to compute the throughput and ~~queue length~~ ~~first~~ the mean # requests in the server.

Throughput

$$= X_0(1) P(1) + X_0(2) P(2) + X_0(3) (P(3) + P(4) + \dots)$$

$$= X_0(1) \cdot \rho_1 + X_0(2) \cdot \rho_2 + \rho_1 + \rho_2 + \dots$$

$$X_0(3) (\rho_3 \rho_2 \rho_1 + \rho_3^2 \rho_2 \rho_1 + \dots)$$

$$= X_0(1) \rho_1 + X_0(2) \rho_2 + \rho_1 + \rho_2 + \dots$$

$$X_0(3) \rho_3 \rho_2 \rho_1 \cdot \frac{1}{1 - \rho_3}$$

$$= \left(X_0(1) \rho_1 + X_0(2) \rho_2 + X_0(3) \frac{\rho_3 \rho_2 \rho_1}{1 - \rho_3} \right) \cdot \frac{1}{1 + \rho_1 + \frac{\rho_2 \rho_1}{1 - \rho_3}}$$

You can plug the values of $\rho_1, \rho_2, \rho_3, X_0(1), X_0(2)$ and $X_0(3)$ into the expression.

The mean # requests in the server is

$$0 P(0) + 1 \cdot P(1) + 2 P(2) + 3 P(3) + \dots$$

$$= P_1 P(0) + 2 P_2 P_1 P(0) +$$

$$3 P_3 P_2 P_1 P(0) + 4 P_3^2 P_2 P_1 P(0) + 5 P_3^3 P_2 P_1 P(0)$$

$$= P_1 P(0) +$$

$$P_2 P_1 P(0) \left[2 + 3 P_3 + 4 P_3^2 + 5 P_3^3 + \dots \right]$$

Need to recognize that this is an arithmetic-geometric progression. You can ~~use~~ ^{use the given} formula ~~but~~ or you can derive the result. Let us derive the result:

$$\text{Let } Z = 2 + 3 P_3 + 4 P_3^2 + 5 P_3^3 + \dots \quad \text{--- (1)}$$

$$\text{then } P_3 Z = 2 P_3 + 3 P_3^2 + 4 P_3^3 + \dots \quad \text{--- (2)}$$

Subtracting (2) from (1), we have

$$(1 - P_3) Z = 2 + P_3 + P_3^2 + \dots$$

$$= 2 + \frac{P_3}{1 - P_3}$$

(3)

$$\Rightarrow Z = \frac{\rho_3}{(1-\rho_3)^2} + \frac{2}{1-\rho_3}$$

Alternatively,
 substitute
 $p=2, m=\rho_3$
 $\rho=1$ in the
 given formula

Thus the mean # requests in the web server is

$$= \rho, \rho(0) + \rho_2 \rho, \rho(0) \cdot \left[\frac{\rho_3}{(1-\rho_3)^2} + \frac{2}{1-\rho_3} \right]$$

You can now plug the values of ρ_1, ρ_2, ρ_3 and

$\rho(0)$ to find the mean # requests

Finally, to obtain the mean response time,

we use Little's law

$$\text{Mean response time} = \frac{\text{mean \# requests in the server}}{\text{mean throughput}}$$