

↖ 5 straight vowels :) (1)

# Hsu Chapter 9: Queueing Theory

To do queueing theory we need two preparatory sections - one about exponential distributions and one about so called birth and death processes.

## Exponential distributions

**Def** A random variable  $T \geq 0$  is exponential distributed with parameter  $\lambda > 0$  (= with expected value  $1/\lambda > 0$ ) if  $f_T(t) = \lambda e^{-\lambda t}$  for  $t \geq 0$ . Notation  $T$  is exp( $\lambda$ ).

**Thm 1** If  $T_1, \dots, T_n$  are independent  $\text{exp}(\lambda_1), \dots, \text{exp}(\lambda_n)$  then  $\min(T_1, \dots, T_n)$  is  $\text{exp}(\lambda_1 + \dots + \lambda_n)$ .

**Proof**  $P(\min(T_1, \dots, T_n) > t) = P(\bigcap_{i=1}^n \{T_i > t\}) = \prod_{i=1}^n P(T_i > t)$   
 $= \prod_{i=1}^n e^{-\lambda_i t} = e^{-(\lambda_1 + \dots + \lambda_n)t} = P(\text{exp}(\lambda_1 + \dots + \lambda_n) > t) \#$

**Thm 2** If  $T_1, \dots, T_n$  are independent  $\text{exp}(\lambda_1), \dots, \text{exp}(\lambda_n)$  then  $P(\min(T_1, \dots, T_n) = T_i) = \lambda_i / (\lambda_1 + \dots + \lambda_n)$ .

**Proof** By Thm 1 it is enough to do proof for  $n=2$ :

$$P(\min(T_1, T_2) = T_1) = P(T_2 \geq T_1) = \int_{x=0}^{+\infty} \int_{y=x}^{+\infty} \lambda_1 e^{-\lambda_1 x} \lambda_2 e^{-\lambda_2 y} dy dx$$
$$= \int_{x=0}^{+\infty} \lambda_1 e^{-\lambda_1 x} e^{-\lambda_2 x} dx = \frac{\lambda_1}{\lambda_1 + \lambda_2} \#$$

Note For  $\lambda$  mean  $\text{exp}(\lambda)$

**Thm 3** (Lack of Memory) For  $T \sim \text{exp}(\lambda)$  we have  $P(T > t+s | T > s) = P(T > t)$  for  $s, t \geq 0$

Proof  $P(T > t+s | T > s) = \frac{P(T > t+s, T > s)}{P(T > s)} = \frac{e^{-\lambda(t+s)}}{e^{-\lambda s}} = e^{-\lambda t} = P(T > t)$  #

Note It is also the other way around: a continuously distributed r.v.  $T \geq 0$  with lack of memory is  $\text{exp}(\lambda)$ .

### Birth and death processes

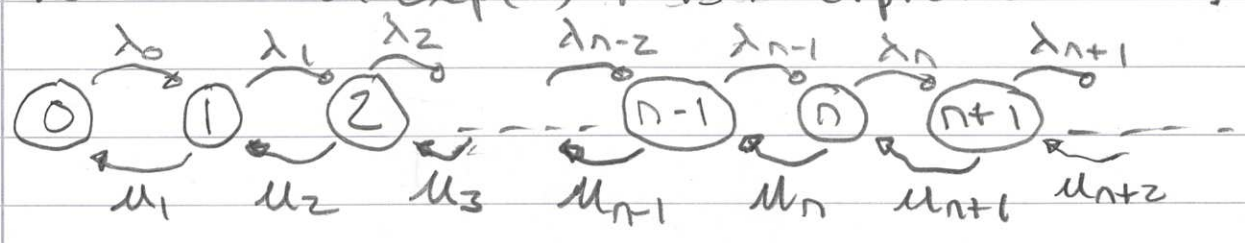
A birth and death process is an  $\mathbb{N}$ -valued random process  $\{X(t)\}_{t \geq 0}$  such that

- $X(0)$  has a certain random or non-random value
- when  $X(t)$  gets a certain value  $n \in \mathbb{N}$  it stays with that value an  $\min(\text{exp}(\mu_n), \text{exp}(\lambda_n)) = \text{exp}(\mu_n + \lambda_n)$  distributed time independent and then change value to

$$\begin{cases} n-1 & \text{if } \text{exp}(\mu_n) = \min(\text{exp}(\mu_n), \text{exp}(\lambda_n)) \text{ w.p. } \frac{\mu_n}{\mu_n + \lambda_n} \\ n+1 & \text{if } \text{exp}(\lambda_n) = \min(\text{exp}(\mu_n), \text{exp}(\lambda_n)) \text{ w.p. } \frac{\lambda_n}{\mu_n + \lambda_n} \end{cases}$$

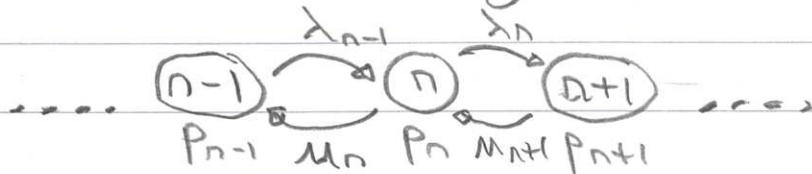
where  $\mu_0 = 0$ , and  $\lambda_1, \mu_1, \lambda_2, \mu_2, \dots \geq 0$

Note For  $\lambda = 0$  an  $\text{exp}(\lambda)$   $T$  is interpreted  $T = +\infty$ .

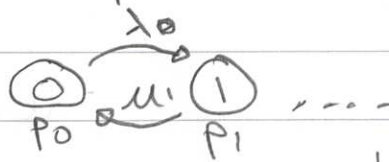


The queueing systems we consider will be special cases of birth and death processes (see below).

We will do analytical calculations on them when they are in steady state equilibrium. This means that  $p_n = P(X(t) = n)$  doesn't depend on  $t$ . This in turn means that the "flow" out from a state  $n$  is same as the flow into that state [as otherwise  $P(X(t) = n)$  will change with time], e.g.,



$$\text{outflow} = p_n(\mu_n + \lambda_n) = \text{inflow} = p_{n-1}\lambda_{n-1} + p_{n+1}\mu_{n+1} \quad \text{for } n \geq 1$$



$$\text{outflow} = p_0\lambda_0 = \text{inflow} = p_1\mu_1$$

This is a second order inhomogeneous difference equation for  $p_n$  with solution

$$p_n = p_0 \frac{\lambda_0 \cdots \lambda_{n-1}}{\mu_1 \cdots \mu_n} \quad \text{for } n \geq 1, \quad p_0 = p_0$$

$$\text{where } p_0 = \left( 1 + \sum_{n=1}^{+\infty} \frac{\lambda_0 \cdots \lambda_{n-1}}{\mu_1 \cdots \mu_n} \right)^{-1} \text{ so } \sum_{n=0}^{+\infty} p_n = 1.$$

Easy to see this  $p_n$  solves difference equations

# Queueing systems

- Our queueing system is denoted

$$M/M/s/K \text{ or } M(\lambda)/M(\mu)/s/K$$

where  $\lambda, \mu > 0$ ,  $s \in \{1, 2, \dots\}$  and  $K \in \{s, s+1, \dots, +\infty\}$  and one usually writes

$$M/M/s \text{ or } M(\lambda)/M(\mu)/s \text{ when } K = \infty.$$

- In birth and death process language this means that

$$\begin{cases} \lambda_n = \lambda & \text{for } n = 0, \dots, K-1 \text{ (=all } n \text{ if } K = \infty) \\ \lambda_n = 0 & \text{for } n \geq K \end{cases}$$

$$\begin{cases} \mu_n = n\mu & \text{for } n \in \{1, \dots, s\} \\ \mu_n = s\mu & \text{for } n \in \{s, \dots, K\} \text{ (=all } n \text{ if } K = \infty) \end{cases}$$

- Hence the birth and death process  $X(t)$  cannot pass the value  $K$  when  $K < \infty$  and is restricted to have values  $\{0, \dots, K\}$  if  $K$  is finite and values  $\mathbb{N}$  if  $K = \infty$ .

- In queueing system language we have a queueing system with

- $X(t) \in \{0, \dots, K\}$  ( $= \mathbb{N}$  if  $K = \infty$ ) customers in system at time  $t$

5

- new customers arrive with independent  $\exp(\lambda)$  interarrival times
- system has  $S$  servers that need independent  $\exp(\mu)$  service times to serve a customer
- system has  $K$ -s queueing slots (infinitely many if  $K = \infty$ ) where customer wait to be served if all servers were busy when they arrived to system
- if  $K < \infty$  customers that arrive to system when it is full  $X(t) = K$  "bounce away" and disappears / never joins system
- \* when  $X(t)$  gets the value  $n$  the next value of  $X(t)$   $n-1$  or  $n+1$  will be a competition of which happens first (is the smallest) of an  $\exp(\mu_n)$  waiting time until first server is ready and  $\exp(\lambda_n)$  waiting time until next customer arrives and when that competition is over a new competition start (by lack of memory) that decides what the next value after that is.

- each competition lasts an  $\exp(\mu_n + \lambda_n)$  time with mean  $\frac{\lambda_n}{\mu_n + \lambda_n}$  and next value is  $X(t) = n+1$  w.p.  $\frac{\lambda_n}{\mu_n + \lambda_n}$  and  $X(t) = n-1$  w.p.  $\frac{\mu_n}{\mu_n + \lambda_n}$

- when we do calculations we assume  $X(t)$  has the steady state equilibrium distribution  $P(X(t) = n) = p_n$


- when we do calculations we are interested in six quantities besides  $p_n$ , namely

$L$  = mean number of customers in queueing system

$L_q$  =  queueing for service

$L_s$  =  that are served

$W$  = mean time customer spends in system

$W_q$  =  queueing

$W_s$  =  in service

The mean time between customers that tries to join queueing system is  $1/\lambda$  so on average  $\lambda$  customers tries to join per time unit. For  $K < \infty$  not all customers that try to join really joins as system can be full  $X(t) = K$  when they try to join so that they bounce away/disappears.

**Def** Traffic intensity  $\rho = \frac{\lambda}{\mu s}$  (must be  $< 1$ ) (for  $K = \infty$ )

• The real/actual average number of customers joining the system is 7

$$\lambda_e = \lambda_{\text{efficient}} = \begin{cases} \lambda & \text{for } K = \infty \\ \lambda(1 - P_K) & \text{for } K < \infty \end{cases}$$

The errata list for the book is particularly important for Chapter 9 as Hsu has used  $\lambda_e$  erroneously in some of his calculations creating contradictions.

**Thm**

$$L = \sum_{n=0}^K n p_n = L_q + L_s$$

$$W = W_q + W_s$$

$$L_q = \lambda_e W_q, L_s = \lambda_e W_s$$

$$W_s = 1/\mu$$

Proof By inspection. #

Consequence  $W_s$  gives  $L_s$  gives  $L_q = L - L_s$   
 gives  $W_q = L_q / \lambda_e$  gives  $W = W_q + W_s$  ;)

$M(\lambda) / M(\mu) / 1$   $p_n = p_0 \left(\frac{\lambda}{\mu}\right)^n = \left(1 - \frac{\lambda}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^n$

$$L = \sum_{n=0}^{+\infty} n \left(1 - \frac{\lambda}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^n = \left[ \text{mean-1 of geometric} \left(1 - \frac{\lambda}{\mu}\right) \right] = \frac{1}{1 - \frac{\lambda}{\mu}} - 1 = \frac{\lambda}{\mu - \lambda}$$

$$W_s = 1/\mu, L_s = \lambda/\mu, L_q = \frac{\lambda}{\mu - \lambda} - \frac{\lambda}{\mu} = \frac{\lambda^2}{\mu(\mu - \lambda)}$$

$$W_q = \frac{\lambda}{\mu(\mu - \lambda)}, W = \frac{\lambda}{\mu(\mu - \lambda)} + \frac{1}{\mu} = \frac{\mu}{\mu(\mu - \lambda)}$$