

# STATS 217: Introduction to Stochastic Processes I

## Lecture 25

## Last time: Stationary distributions

- Let  $(X_t)_{t \geq 0}$  be a CTMC on  $\Omega$ . A probability distribution  $\pi$  on  $\Omega$  is said to be a stationary distribution if

$$\pi P^t = \pi \quad \forall t \geq 0$$

- This is equivalent to the condition that

$$\pi Q = 0.$$

- In terms of the matrix  $Q$ , the detailed balance conditions are given by

$$\pi_i q_{ij} = \pi_j q_{ji} \quad \forall i, j \in \Omega$$

# Convergence theorem

Let  $(X_t)_{t \geq 0}$  be an irreducible CTMC on a finite state space  $\Omega$ . Then, there exists a unique stationary distribution  $\pi$ , and

$$\max_{x \in \Omega} \text{TV}(P^t(x, \cdot), \pi) \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

# Convergence theorem

Let  $(X_t)_{t \geq 0}$  be an irreducible CTMC on a finite state space  $\Omega$ . Then, there exists a unique stationary distribution  $\pi$ , and

$$\max_{x \in \Omega} \text{TV}(P^t(x, \cdot), \pi) \rightarrow 0 \quad \text{as } t \rightarrow \infty.$$

We have already done the work to prove this theorem.

# Existence of stationary distribution

- The first point is the existence of the stationary distribution.

## Existence of stationary distribution

\* assume that the chain  $(X_t)_{t \geq 0}$   
is described using  $Q$ .

- The first point is the existence of the stationary distribution.
- Recall the notation  $\lambda_i = \sum_{j \neq i} q_{ij}$ ,  $\Lambda = \max_{i \in \Omega} \lambda_i$ .
- Since  $\Omega$  is finite,  $\Lambda < \infty$ . In the case when  $\Lambda < \infty$ ,

we had a simpler way of  
simulating  $(X_t)_{t \geq 0}$  given  
 $Q$ .

in this case, we were able to show  
that  $X_t = Y_{N(t)}$  where  $N(t)$  is a PPP  
w/ rate  $\Lambda$  and  $Y_n$  is a DTMC.

## Existence of stationary distribution

- The first point is the existence of the stationary distribution.
- Recall the notation  $\lambda_i = \sum_{j \neq i} q_j$ ,  $\Lambda = \max_{i \in \Omega} \lambda_i$ .
- Since  $\Omega$  is finite,  $\Lambda < \infty$ . In this case, recall that we have the representation  $X_t = Y_{N(t)}$ , where  $N(t)$  is a PPP with rate  $\lambda$  and  $Y_n$  is a DTMC with the transition matrix

$$U_{ij} = \frac{q_{ij}}{\Lambda} \quad \forall i \neq j \quad U_{ii} = 1 - \frac{\lambda_i}{\Lambda}$$

## Existence of stationary distribution

- The first point is the existence of the stationary distribution.
- Recall the notation  $\lambda_i = \sum_{j \neq i} q_j$ ,  $\Lambda = \max_{i \in \Omega} \lambda_i$ .
- Since  $\Omega$  is finite,  $\Lambda < \infty$ . In this case, recall that we have the representation  $X_t = Y_{N(t)}$ , where  $N(t)$  is a PPP with rate  $\lambda$  and  $Y_n$  is a DTMC with the transition matrix

$$U_{ij} = \frac{q_{ij}}{\Lambda} \quad \forall i \neq j \quad U_{ii} = 1 - \frac{\lambda_i}{\Lambda}$$

- Since  $(X_t)_{t \geq 0}$  is irreducible, so is  $U$ , and hence, it has a unique stationary distribution  $\pi$ .

o goal: show that  $\pi P_t = \pi \quad \forall t \geq 0$   
( $\Rightarrow \pi Q = 0$ )



## Existence of stationary distribution

- We can check that  $\pi Q = 0$ .

## Existence of stationary distribution

- We can check that  $\pi Q = 0$ . Indeed,

$$(\pi Q)_j = \sum_{i \in \Omega} \pi_i q_{ij} = \pi_j q_{jj} + \sum_{j \neq i} \pi_i q_{ij}$$

## Existence of stationary distribution

- We can check that  $\pi Q = 0$ . Indeed,

$$\begin{aligned}\sum_{i \in \Omega} \pi_i q_{ij} &= \pi_j q_{jj} + \sum_{j \neq i} \pi_i q_{ij} \\ &= -\pi_j \lambda_j + \sum_{i \neq j} \pi_i \underbrace{U_{ij} \Lambda}\end{aligned}$$

$$U_{ij} = \frac{q_{ij}}{\Lambda} \quad \forall i \neq j$$

$$\sum_j U_{jj} = 1 - \frac{\lambda_j}{\Lambda}$$

$$\lambda_j = -q_{jj}$$

$$\lambda_j = \sum_{k \neq j} q_{jk}$$

# Existence of stationary distribution

- We can check that  $\pi Q = 0$ . Indeed,

$$\pi \mathcal{U} = \pi$$

$$\begin{aligned} \sum_{i \in \Omega} \pi_i q_{ij} &= \pi_j q_{jj} + \sum_{j \neq i} \pi_i q_{ij} \\ &= -\pi_j \lambda_j + \sum_{i \neq j} \pi_i \underline{U_{ij} \Lambda} \\ &= -\pi_j \lambda_j + \Lambda \sum_{i \in \Omega} \pi_i \underline{U_{ij}} - \underline{\Lambda \pi_j U_{jj}} \\ &= -\pi_j \lambda_j + \underline{\Lambda \pi_j} - \pi_j (\Lambda - \lambda_j) \\ &= 0. \end{aligned}$$

o wts : if  $\pi Q = 0 = \mu Q \Rightarrow \pi = \mu$ .  
 $\pi P = \pi \iff \mu P = \mu$ .  $\implies$  since  $P$  is irred.  $\pi = \mu$ .

# Convergence

$$\text{wts: } \max_x TV(P^t(x, \cdot), \pi) \xrightarrow{t \rightarrow \infty} 0$$

- Note that

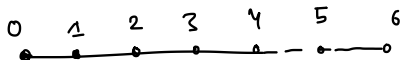
we know:  $\tilde{P}$  is an irr + aperiodic transition matrix on finite  $\Omega$ , then  $\max TV(\tilde{P}^n(x, \cdot), \pi)$

$$TV(P^{t+s}(x, \cdot), \pi) = TV(\delta_x P^t P^s, \pi P^s) \xrightarrow{n \rightarrow \infty} 0$$

$Y_n$  w/ transition matrix  $\mathbb{P} (= P^1)$

we know that  $P^1$  is irr + aperiodic (Levy's dichotomy)

$$\Rightarrow \max_x TV(\mathbb{P}^n(x, \cdot), \pi) \xrightarrow{n \rightarrow \infty} 0$$



Recall: that for DTMC:  
we found  $\exists$  some  $\epsilon_0 > 0$  s.t.

# Convergence

every entry of  $\mathbb{P}^{r_0} > 0$ .  
then we showed

$$\max_K \text{TV}(\mathbb{P}^{r_0 K}(x, \cdot), \pi) \xrightarrow{K \rightarrow \infty} 0$$

- Note that

$$\mathbb{P}^{t+s} = \mathbb{P}^t \mathbb{P}^s$$

$$\begin{aligned} \text{TV}(\underbrace{\mathbb{P}^{t+s}(x, \cdot)}_{\text{prob. distribution w/}}, \pi) &= \text{TV}(\underbrace{\delta_x}_{\text{prob}[y]} \mathbb{P}^t \mathbb{P}^s, \underbrace{\pi}_{\text{prob}[y]}) \\ &\leq \text{TV}(\delta_x \mathbb{P}^t, \pi). \end{aligned}$$

prob. distribution  
w/

$$\text{prob}[y] = \mathbb{P}^{t+s}(x, y)$$

$$\mathbb{P}^{t+s}(x, \cdot) = \underbrace{\delta_x}_{\text{constant distribution at } x} \mathbb{P}^{t+s} = \delta_x \mathbb{P}^t \mathbb{P}^s$$

i.e. 1D - dimensional row  
vector which is 1 at  $x$   
0 elsewhere

# Convergence

- Note that

$$\begin{aligned} \text{TV}(P^{t+s}(x, \cdot), \pi) &= \text{TV}(\delta_x P^t P^s, \pi P^s) \\ &\leq \text{TV}(\delta_x P^t, \pi). \end{aligned}$$

- Hence,  $\text{TV}(P^t(x, \cdot), \pi)$  is non-increasing in  $t$ , so it suffices to show that it converges to 0 along (say) the natural numbers.

# Convergence

① first gen. joint sample from

$$\left( \hat{X}, \hat{Y} \right) = \left( \delta_x P^t, \pi \right)$$

the opt. coupling of

②  $\left( \hat{X} P^s, \hat{Y} P^s \right) \sim \left( \delta_x P^t P^s, \pi P^s \right)$

$$\left. \begin{aligned} \text{TV}(P^{t+s}(x, \cdot), \pi) &= \text{TV}(\delta_x P^t P^s, \pi P^s) \\ &\leq \text{TV}(\delta_x P^t, \pi) \end{aligned} \right\} \textcircled{3} \mathbb{P}[\hat{X} P^s \neq \hat{Y} P^s] \leq \mathbb{P}[\hat{X} \neq \hat{Y}]$$

- Hence,  $\text{TV}(P^t(x, \cdot), \pi)$  is non-increasing in  $t$ , so it suffices to show that it converges to 0 along (say) the natural numbers.
- But  $P^1$  is an irreducible and aperiodic transition matrix with unique stationary distribution  $\pi$ , so that by looking at the corresponding DTMC, we have

$$\text{TV}(P^n(x, \cdot), \pi) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$



## Example: $M/M/1$ queues

*# of servers*

- This is a popular queuing model in which the arrival of customers is modelled by a Poisson point process with rate  $\lambda$ . There is a single server, and service times are independent and exponentially distributed with parameter  $\mu$ .

## Example: M/M/1 queues

- This is a popular queuing model in which the arrival of customers is modelled by a Poisson point process with rate  $\lambda$ . There is a single server, and service times are independent and exponentially distributed with parameter  $\mu$ .
- Due to the memorylessness property of the exponential distribution, this can be modelled as a continuous time birth and death chain with jump rates

think of  
 $n$  as the  
# of customers  
in the system.

$$Q_{n,n+1} = \lambda, \quad n = 0, 1, \dots$$

$$Q_{n,n-1} = \mu, \quad n = 1, 2, \dots$$

} continuous time birth & death chain.

## Example: M/M/1 queues

- This is a popular queuing model in which the arrival of customers is modelled by a Poisson point process with rate  $\lambda$ . There is a single server, and service times are independent and exponentially distributed with parameter  $\mu$ .
- Due to the memorylessness property of the exponential distribution, this can be modelled as a continuous time birth and death chain with jump rates

$$Q_{n,n+1} = \lambda, \quad n = 0, 1, \dots$$

$$Q_{n,n-1} = \mu, \quad n = 1, 2, \dots$$

- Suppose instead that there are  $s$  servers, and customers are served if there is at least one server available. This is called the  $M/M/s$  queueing model,

## Example: M/M/1 queues

- This is a popular queuing model in which the arrival of customers is modelled by a Poisson point process with rate  $\lambda$ . There is a single server, and service times are independent and exponentially distributed with parameter  $\mu$ .
- Due to the memorylessness property of the exponential distribution, this can be modelled as a continuous time birth and death chain with jump rates

$$Q_{n,n+1} = \lambda, \quad n = 0, 1, \dots$$

$$Q_{n,n-1} = \mu, \quad n = 1, 2, \dots$$

- Suppose instead that there are  $s$  servers, and customers are served if there is at least one server available. This is called the  $M/M/s$  queuing model, and the jump rates are now

$$\left\{ \begin{array}{l} Q_{n,n+1} = \lambda, \quad n = 0, 1, \dots, \\ Q_{n,n-1} = n\mu, \quad n = \overbrace{1, \dots, s}, \\ Q_{n,n-1} = s\mu, \quad n = s+1, s+2, \dots, \end{array} \right. \begin{array}{l} M/G/s \\ G/I/s \\ M/M/\infty \end{array}$$

## M/M/1 queues

arrival rate

service rate

- Suppose that  $\lambda < \mu$  i.e., the rate of arrivals is smaller than the rate of service. Otherwise, the size of the queue explodes.

## M/M/1 queues

- Suppose that  $\lambda < \mu$  i.e., the rate of arrivals is smaller than the rate of service. Otherwise, the size of the queue explodes.
- When  $\lambda < \mu$ , we can use the detailed balance conditions

$$\pi_i Q_{ij} = \pi_j Q_{ji}$$

## M/M/1 queues

- Suppose that  $\lambda < \mu$  i.e., the rate of arrivals is smaller than the rate of service. Otherwise, the size of the queue explodes.
- When  $\lambda < \mu$ , we can use the detailed balance conditions

$$\pi_i Q_{ij} = \pi_j Q_{ji}$$

to find the stationary distribution

$$\pi_n = \left(1 - \frac{\lambda}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^n, \quad n = 0, 1, \dots$$

Q 1: server occupancy  
: server occupied  $\Leftrightarrow n \geq 1$   
 $= 1 - \pi_0$

## M/M/1 queues

- Suppose that  $\lambda < \mu$  i.e., the rate of arrivals is smaller than the rate of service. Otherwise, the size of the queue explodes.
- When  $\lambda < \mu$ , we can use the detailed balance conditions

$$\pi_i Q_{ij} = \pi_j Q_{ji}$$

to find the stationary distribution

$$\pi_n = \left(1 - \frac{\lambda}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^n, \quad n = 0, 1, \dots$$

- Given this stationary distribution, one can compute many quantities of interest.



## M/M/1 queues

- Suppose that  $\lambda < \mu$  i.e., the rate of arrivals is smaller than the rate of service. Otherwise, the size of the queue explodes.
- When  $\lambda < \mu$ , we can use the detailed balance conditions

$$\pi_i Q_{ij} = \pi_j Q_{ji}$$

to find the stationary distribution

$$\pi_n = \left(1 - \frac{\lambda}{\mu}\right) \left(\frac{\lambda}{\mu}\right)^n, \quad n = 0, 1, \dots$$

- Given this stationary distribution, one can compute many quantities of interest. For instance, the long-run fraction of time that the server is busy is

$$1 - \pi_0 = \frac{\lambda}{\mu}.$$

## M/M/1 queues

(i) waiting time /  
time spent in the  
system.

- Moreover, the expected length of the queue under the equilibrium distribution is

$$L = \sum_{n=0}^{\infty} n\pi_n = \frac{\lambda}{\mu - \lambda} \quad \begin{array}{l} \lambda = 10 \\ \mu = 15 \end{array}$$

$$\frac{10}{15 - 10} = 2$$

## M/M/1 queues

- Moreover, the expected length of the queue under the equilibrium distribution is

$$L = \sum_{n=0}^{\infty} n\pi_n = \frac{\lambda}{\mu - \lambda}.$$

- Another important quantity is the total time  $T$  (waiting time + time with the server) spent by a customer in the system.

## M/M/1 queues

- Moreover, the expected length of the queue under the equilibrium distribution is

$$L = \sum_{n=0}^{\infty} n\pi_n = \frac{\lambda}{\mu - \lambda}.$$

- Another important quantity is the total time  $T$  (waiting time + time with the server) spent by a customer in the system.
- If there are  $n$  customers already in the system when a new customer joins the queue,

## M/M/1 queues

- Moreover, the expected length of the queue under the equilibrium distribution is

$$L = \sum_{n=0}^{\infty} n\pi_n = \frac{\lambda}{\mu - \lambda}.$$

- Another important quantity is the total time  $T$  (waiting time + time with the server) spent by a customer in the system.
- If there are  $n$  customers already in the system when a new customer joins the queue, then since service times are i.i.d. exponentials with parameter  $\mu$ , the total time spent by the customer is distributed as a sum of  $n + 1$  i.i.d. exponentials with parameter  $\mu$ .

## M/M/1 queues

- Then, using the law of total probability, we have

$$\mathbb{P}[\mathcal{T} \leq t] = \mathbb{P}[\mathcal{T} \leq t \mid n \text{ customers already in the system}] \cdot \pi_n$$

## M/M/1 queues

- Then, using the law of total probability, we have

$$\begin{aligned}\mathbb{P}[T \leq t] &= \mathbb{P}[T \leq t \mid n \text{ customers already in the system}] \cdot \pi_n \\ &= 1 - \exp(-t(\mu - \lambda)),\end{aligned}$$

# M/M/1 queues

- Then, using the law of total probability, we have

$$\begin{aligned}\mathbb{P}[T \leq t] &= \mathbb{P}[T \leq t \mid n \text{ customers already in the system}] \cdot \pi_n \\ &= 1 - \exp(-t(\mu - \lambda)),\end{aligned}$$

i.e.  $T$  has exponential distribution with mean

$$W = \frac{1}{\mu - \lambda}$$

$$\boxed{L = \lambda W}$$

$$\begin{aligned} & \mu, \lambda \\ & W, L = \frac{\lambda}{\mu - \lambda} \\ & \frac{1}{\mu - \lambda} \end{aligned}$$

$$\begin{aligned} \lambda &= 10 \\ \mu &= 15 \\ W &\sim \text{Exp}(5). \\ & \text{w. mean} \\ & \frac{1}{15 - 10} = \frac{1}{5} \end{aligned}$$



## M/M/1 queues

- Then, using the law of total probability, we have

$$\begin{aligned}\mathbb{P}[T \leq t] &= \mathbb{P}[T \leq t \mid n \text{ customers already in the system}] \cdot \pi_n \\ &= 1 - \exp(-t(\mu - \lambda)),\end{aligned}$$

i.e.  $T$  has exponential distribution with mean

$$W = \frac{1}{\mu - \lambda} = \frac{L}{\lambda}.$$

# Little's law

- The relationship

$$L = \lambda W$$

is called **Little's law**

# Little's law

- The relationship

$$L = \lambda W$$

is called **Little's law** and is true even without the specific distributional assumptions (i.e. Poisson arrivals and exponential waiting times).

# Little's law

- The relationship

$$L = \lambda W$$

is called **Little's law** and is true even without the specific distributional assumptions (i.e. Poisson arrivals and exponential waiting times). Such queues are called *GI/G/1* queues.

# Little's law

- The relationship

$$L = \lambda W$$

is called **Little's law** and is true even without the specific distributional assumptions (i.e. Poisson arrivals and exponential waiting times). Such queues are called *GI/G/1* queues.

- Here's the intuition: Suppose each customer pays \$1 for each minute of time they spend in the system. When there are  $n$  customers in the system, the establishment is earning \$ $n$  per minute, and hence, the establishment is earning an average of \$ $L$  per minute.

# Little's law

- The relationship

$$L = \lambda W$$

is called **Little's law** and is true even without the specific distributional assumptions (i.e. Poisson arrivals and exponential waiting times). Such queues are called *GI/G/1* queues.

- Here's the intuition: Suppose each customer pays \$1 for each minute of time they spend in the system. When there are  $n$  customers in the system, the establishment is earning \$ $n$  per minute, and hence, the establishment is earning an average of \$ $L$  per minute.
- On the other hand, if each customer pays for their entire duration when they arrive, then the average rate of earning is  $\lambda \times W$  per minute.