# REVIEW: MARKOV JUMP-PROCESS (MJP)

**MJP**   $X = \{X_t, \ t \geq 0\}$ on $\mathcal{S} = \{i, j, \ldots\}$ countable.
Markov property: $P_r\{X_t = j | X_r, \ r < s; \ X_s = i\} = P_{ij}(s,t), \quad \forall \, s < t, \ \forall i, j \in \mathcal{S}$.
Time homogeneity: $P_r\{X_{s+t} = j | X_s = i\} = P_{ij}(t), \quad \forall \, s, t, \ i, j,$ transition probabilities.

Characterization: $\pi^0 =$ initial distribution and $P(t) = [P_{ij}(t)], \ t \geq 0$, stochastic.
Finite-dimensional distributions:
$P_r\{X_0 = i_0, \ X_{t_1} = i_1, \ldots, X_{t_n} = i_n\} = \pi^0(i_0) P_{i_0, i_1}(t_1) \ldots P_{i_{n-1}, i_n}(t_n - t_{n-1})$.

$P(t)$ : stochastic ; $P(s+t) = P(s)P(t), \quad \forall \, s, t$ (Chapman Kolmogorov);
$\quad \exists \, P(0) = I \ ; \ \exists \, \dot{P}(0) = Q = [q_{ij}]$, infinitesimal generator $\quad \left( \sum_{j \in \mathcal{S}} q_{ij} = 0 \right)$.

 Micro to Macro : $\dot{P}(t) = P(t)Q \ (= QP(t))$ and $P(0) = I$
   Forward (Backward) equations.

   Solution : $P(t) = \exp[tQ] = \sum_{n=0}^{\infty} \frac{t^n}{n!} \, Q^n , \ t \geq 0$.

Animation: $\quad i \xrightarrow{q_{ij}} j; \quad \forall \, i, j \in \mathcal{S} \ \exists$ exponential clock at rate $q_{ij}$, call it $(i, j)$.
Given $i$, consider clocks $(i, j), \ j \in \mathcal{S}$; move to the "winner" when rings.
Thus: stay at $i \sim \exp(q_i = \sum_{j \neq i} q_{ij})$ and switch to $j$ with probability $P_{ij} = q_{ij}/q_i$
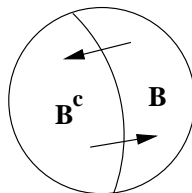$(q_{ij} = q_i P_{ij}, i \neq j; q_{ii} = -q_i)$.

 Transient analysis    vs. long-run/limit    stability/steady-state
     $\exists \, \lim_{t \uparrow \infty} P_{ij}(t) = \pi_j, \ \forall \, i; \quad \pi = \pi P(t), \ \forall \, t$.

Calculation via **steady-state equations**: $\dot{P}(\infty) = P(\infty)Q \Rightarrow \left\{ \begin{array}{l} 0 = \pi Q \\ \sum_i \pi_i = 1, \ \pi_i \geq 0 \end{array} \right\}$

or balance equations: $\sum_{i \neq j} \pi_i q_{ij} = -\pi_j q_{jj} = \sum_{i \neq j} \pi_j q_{ji}, \ \forall \, j$.

Transition rates: $\pi_i q_{ij} =$ long-run average number of switches from $i$ to $j$.

Cuts: $\quad \sum_{i \in B} \sum_{j \in B^c} \pi_i q_{ij} = \sum_{i \in B^c} \sum_{j \in B} \pi_i q_{ij}, \ \forall \, B \subset \mathcal{S}$.
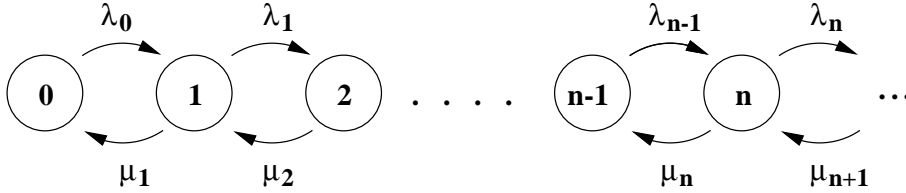
**Ergodic Theorem:** Let $X$ be *irreducible* $(i \leftrightarrow j)$. Assume that there exists a solution $\pi$ to its steady-state equations. Then, $X$ must be "unexplosive" and $\pi$ must be its stationary distribution, its limit distribution and

**SLLN** $\bullet \lim\limits_{T \uparrow \infty} \frac{1}{T} \int_0^T f(X_t)dt = \sum\limits_i \pi_i f(i)$ $(\text{``=''} Ef(X_\infty))$ ; eg. $f(x) = 1_B(x)$.

$\bullet \lim\limits_{T \uparrow \infty} \frac{1}{T} \sum\limits_{t \leq T} g(X_{t-}, X_t) = \sum\limits_i \pi_i \sum\limits_j q_{ij} g(i, j)$, for $g(x, x) = 0$, $\forall x$; e.g. $g(x, y) = 1_C(x, y)$.

## Birth & Death Model of a Service Station (*Hall, §5.4*)



Cuts at $n \leftrightarrow n+1$ yield: $\pi_n \lambda_n = \pi_{n+1} \mu_{n+1}$, $n \geq 0$;

$$\pi_{n+1} = \frac{\lambda_n}{\mu_{n+1}} \pi_n = \frac{\lambda_n \lambda_{n-1}}{\mu_{n+1} \mu_n} \pi_{n-1} = \cdots = \frac{\lambda_0 \lambda_1 \ldots \lambda_n}{\mu_1 \mu_2 \ldots \mu_{n+1}} \pi_0 .$$

The required solution exists if and only if

$$\sum_{n=0}^{\infty} \frac{\lambda_0 \ldots \lambda_n}{\mu_1 \ldots \mu_{n+1}} < \infty .$$

The Ergodic Theorem then yields

$$\begin{cases} \pi_n &= \frac{\lambda_0 \ldots \lambda_{n-1}}{\mu_1 \ldots \mu_n} \pi_0 \ , \ n \geq 0 \\ \pi_0 &= \left[ \sum_{n \geq 0} \frac{\lambda_0 \ldots \lambda_n}{\mu_1 \ldots \mu_{n+1}} \right]^{-1} \end{cases}$$

**Measures of Performance (MOP's), in steady-state:**

$L \ = \ $ number of customers at the service station (also denoted $L_s$);
$L_q \ = \ $ number of customers in the queue;
$W \ = \ $ sojourn time of a customer at the service station (also denoted $W_s$);
$W_q \ = \ $ waiting time of a customer in the queue;

$$E(L) = \sum_{n \geq 0} n\pi_n = \lim_{T \uparrow \infty} \frac{1}{T} \int_0^T L(t)dt.$$

Let $m(n)$ = number of active servers at state $n$, $0 \leq m(n) \leq n$; the servers are **statistically identical**.

$$E(L_q) = \sum_{n \geq 0} [n - m(n)]\pi_n \quad = \text{ also long-run average, as above.}$$

Service rate per server is $\mu(n)/m(n)$, $n \geq 1$.

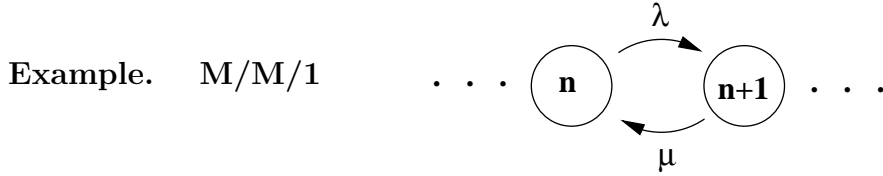**Average** (actual) **service rate:** $\qquad \sum_{n \geq 1} \frac{\mu(n)}{m(n)} \pi_n = E[\mu(L)/m(L)]$.

**Potential service rate** of each server: $\quad E[\mu(L)/m(L) | L > 0] = \frac{E[\mu(L)/m(L)]}{1 - \pi_0}$.

**Inflow rate:** $\qquad\qquad\qquad \lambda = \sum_{n \geq 0} \pi_n q_{n,n+1} = \sum_{n \geq 0} \pi_n \lambda(n) = E\lambda(L)$.

$\uparrow$ arrival

$\downarrow$ departure

**Outflow rate:** $\qquad\qquad\qquad \delta = \sum_{n \geq 0} \pi_n q_{n,n-1} = \sum_{n \geq 0} \pi_n \mu(n) = E\mu(L)$. (Assume $\mu(0) = 0$.)

Note: in steady state, $\pi_n \lambda_n = \pi_{n+1}\mu_{n+1}$, $\forall\, n \geq 0 \Rightarrow$ inflow rate = outflow rate.

**Throughput rate:** $\qquad E\lambda(L) = E\mu(L) \qquad$ (the common quantity).

**Example.** M/M/1 $\qquad \cdots$



$\lambda_j = \lambda$, $j \geq 0$; $\quad \mu_j = \mu \cdot 1_{j \geq 1}$.

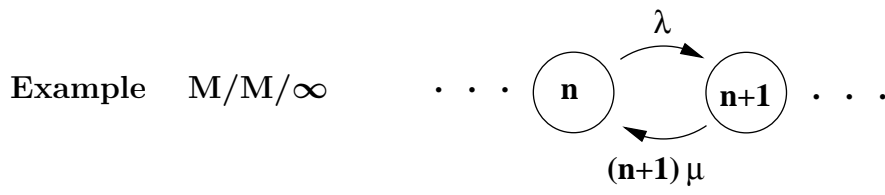$\rho = \frac{\lambda}{\mu} < 1 \quad$ assumed for steady state (traffic intensity).

$\pi_n = (1 - \rho)\rho^n$, $n \geq 0$. $\qquad$ Geometric distribution!

Actual service rate $= \sum_{n \geq 1} \pi_n \cdot \mu = \mu(1 - \pi_0) = \mu \cdot \rho = \mu \cdot \frac{\lambda}{\mu} = \lambda$, contrasted with

Potential service rate $= \frac{\lambda}{1 - \pi_0} = \frac{\lambda}{\rho} = \mu$ , as anticipated.

Additional properties: $W \sim \exp\left(\text{mean} = \frac{1}{\mu(1-\rho)} = \frac{1}{\mu}\left[1 + \frac{\rho}{1-\rho}\right]\right)$, geometric mixture of exp's.

Departure process is Poisson $(\lambda)$ (Burke's Theorem).

$$\frac{W_q}{1/\mu} \stackrel{d}{=} \begin{cases} 0 & \text{wp} \quad 1 - \rho \\[2mm] \exp\left(\text{mean} = \frac{1}{1-\rho}\right) & \text{wp} \quad \rho \end{cases}$$

3

**Example** **M/M/∞**



Always stable.

$\pi_n = e^{-\rho}\frac{\rho^n}{n!}$ , $n \geq 0$ , Poisson distribution!

$E\,(\# \text{ busy servers}) = \lambda \cdot \frac{1}{\mu} = \frac{\lambda}{\mu} = \rho.$

Very *useful*: $\infty$-server models provide upper bound (e.g., Israel Electric Company).

**Example** **M/M/S**

$$\mu_j = (j \wedge s)\mu , \quad \lambda_j \equiv \lambda ,$$

$$\rho = \frac{\lambda}{s\mu} < 1 \qquad \text{assumed, as before, to ensure stability.}$$

$$\pi_k = \frac{a^k}{k!}\,\pi_0, \qquad k \leq S,$$

$$= \frac{s^s \rho^k}{s!}\,\pi_0, \qquad k \geq S,$$

$$\pi_0 = \left[\sum_{j=0}^{s-1} \frac{a^j}{j!} + \frac{a^s}{s!(1-\rho)}\right]^{-1}, \text{ where } a = \frac{\lambda}{\mu} , \text{ offered load.}$$

**Note:** "Wait | Wait $> 0$" is exponential, having the same distribution as that in an M/M/1 queue with arrival rate $\lambda$ and service rate $S \cdot \mu$.

**Erlang-C Formula** (1917):

$$E_{2,S} = \sum_{k \geq s}\pi_k = \frac{a^S}{S!}\,\frac{1}{1-\rho} \cdot \pi_0, \qquad \text{delay probability (PASTA).}$$

**Example** **M/M/S/S**

$$\lambda_j \equiv \lambda, j = 0, \ldots, S-1, \quad \mu_j = j \cdot \mu \text{ for } j = 1, 2, \ldots, S .$$

Always reaches steady state.

$$\pi_k = \frac{a^k}{k!} \Big/ \sum_{j=0}^{S} \frac{a^j}{j!} , \qquad k = 0, 1, \ldots, S.$$

**Erlang-B Formula**:

$$E_{1,S} = \pi_S = \frac{a^s}{s!} \Big/ \sum_{j=0}^{S} \frac{a^j}{j!}, \qquad \text{loss probability (PASTA).}$$

$\lambda \pi_s$ – rate of lost customers,
$\lambda(1 - \pi_s)$ – effective throughput.

**Note:** Useful relations between the Erlang-B and Erlang-C formulae are

$$E_{1,S} = \frac{(S-a)E_{2,S}}{S - aE_{2,S}} \; ; \; E_{2,S} = \frac{E_{1,S}}{(1-\rho) + \rho E_{1,S}} \; ;$$

$$E_{2,S} > E_{1,S}, \text{ as expected: why?}$$

The expression of $E_{2,S}$ in terms of $E_{1,S}$ will become especially useful later on.

**Example   M/M/S/N**          $(S \leq N)$

$$\lambda_j = \lambda, \qquad 0 \leq j \leq N-1, \qquad (\lambda_N = 0)$$
$$\mu_j = (j \wedge S)\mu, \quad 1 \leq j \leq N. \qquad (\mu_0 = 0)$$

Formulae straightforward but cumbersome (simply truncate M/M/S).

Always reaches steady state.

**Note:** Mainly M/M/S (Erlang-C) and sometimes M/M/S/S (Erlang-B) are the prevalent models used in the world of call centers. However, M/M/S/N is more appropriate, and even more so M/M/S/N + Abandonment: Erlang-A.

But the following question then arises: How to model Abandonment?

## Erlang's Formulae
(Exact Results for M/M/m = Erlang-C, and M/M/m/m = Erlang-B)

$$R = \text{offered load } \left(= \lambda/\mu = m \cdot \rho \; ; \; \rho = \tfrac{R}{m}\right)$$

Erlang B: $\qquad E_{1,m} = \dfrac{\frac{R^m}{m!}}{\sum_{k=0}^{m} \frac{R^k}{k!}}$ $\qquad\qquad\qquad$ Probability of *blocking*/loss

Erlang C: $\qquad E_{2,m} = \dfrac{\frac{R^m}{m!} \frac{1}{1-\rho}}{\sum_{k=0}^{m-1} \frac{R^k}{k!} + \frac{R^m}{m!} \frac{1}{1-\rho}}$ $\qquad$ Probability of *delay*

*Relations* (Palm, 1943?)

- Some observations on the Erlang formulae.........pg. 18

- Contributions to the Theory of Delay Systems ......pg. 37

1. $\qquad E_{2,n} = \dfrac{nE_{1,n}}{(n-R) + RE_{1,n}} = \dfrac{E_{1,n}}{(1-\rho) + \rho E_{1,n}}$ $\qquad$ for $\qquad \begin{matrix} \rho < 1 \\ \left(\rho = \tfrac{R}{n}\right) \end{matrix}$

$\qquad E_{2,n} > E_{1,n} \quad ; \quad \frac{d}{dR} E_{2,n}(n) = \frac{1}{nE_{1,n}(n)}$

2. $\qquad E_{2,n} = \dfrac{R(n-1-R)E_{2,n-1}}{(n-1)(n-R) - RE_{2,n-1}}$ $\qquad$ for $R < n-1$.

$\qquad$ (Must have $R < 1$ to start with $E_{2,1} = \rho$)

3. $\qquad E_{1,n} = \dfrac{RE_{1,n-1}}{n + RE_{1,n-1}} = \dfrac{\rho E_{1,n-1}}{1 + \rho E_{1,n-1}} \quad ; \quad E_{1,0} = 1.$

Recursions are useful for calculations.
For example, to calculate $E_{2,n}$, it is convenient to calculate recursively $E_{1,n}$ via 3. and
then calculate $E_{2,n}$ via 1.
They will also be useful for us in asymptotic analysis of systems with many servers.
For example, to analyze the behavior of $E_{2,n}$, as $n \uparrow \infty$, it is convenient to analyze first
$E_{1,n}$, and then use 1.

Recall: Erlang B/C/A formulae, and much more, are implemented in *4CallCenters* that
you have been using.