# 10. Reinforcement learning

Based on dynamic programming.

Transition function: $x' = \delta(x, a)$

Reward function: $r' = \varsigma(x, a)$

we are searching for an optimal action sequence

$$a_0, a_1, a_2, \ldots = \{a_n\}_{k=0}^{\infty} \quad \text{and}$$

and control policy $a_n = \mu(x_n)$ leading to an optimal infinite-horizon discounted value function

$$J^*(x_0) = \max_{\{a_n\}_{k=0}^{\infty}} \sum_{k=0}^{\infty} \gamma^k \varsigma(x_k, a_k)$$

where $\gamma < 1$ is a discounting factor which determines how far ahead rewards should influence a control policy.

$\gamma < 1 \Rightarrow J^*$ is finite

$$J^*(x_0) = \max_{a_0}\left[\varsigma(x_0, a_0) + \gamma \max_{\{a_{k+1}\}_{k=0}^{\infty}} \sum_{k=0}^{\infty} \gamma^k \varsigma(x_{k+1}, a_{k+1})\right]$$

$$= \max_{a_0} \left[ \beta(x_0, a_0) + \gamma \jmath^*(x_1) \right] =$$

$$= \max_{a_0} \left[ \beta(x_0, a_0) + \gamma \jmath^*(\delta(x_0, a_0)) \right]$$

This expression can be applied to any state $x$ $\Rightarrow$

$$\boxed{\jmath^*(x) = \max_{a \in \Sigma(x)} \left[ \beta(x, a) + \gamma \jmath^*(\delta(x, a)) \right]}$$

This famous equation is called Bellman's equation which is the core of dynamic programming.

EX



$$\jmath^*(1) = \max\{\beta(1, b), \beta(1, c)\} + \gamma \jmath^*(2)$$
$$= \max\{1, 2\} + \gamma \jmath^*(2) =$$
$$= 2 + \gamma \jmath^*(2)$$

$$\jmath^*(2) = \max\{\underbrace{\beta(2, d)}_{4}, \underbrace{\beta(2, e)}_{3}\} +$$
$$+ \gamma \jmath^*(1) = 4 + \gamma \jmath^*(1)$$

Fixed point
$$\jmath^*(1) = 2 + \gamma \jmath^*(2) = 2 + \gamma(4 + \gamma \jmath^*(1))$$
$$\Rightarrow \jmath^*(1) = \frac{2 + 4\gamma}{1 - \gamma^2} \quad \jmath^*(2) = \frac{4 + 2\gamma}{1 - \gamma^2}$$

## Q-function

$$Q(x,a) = g(x,a) + \gamma J^*(\delta(x,a))$$

state action pair: $(x,a)$

$$J^*(x) = \max_{a \in \Sigma(x)} Q(x,a)$$

Bellman's equation is now expressed in terms of the Q-function.

$$J^*(\underbrace{\delta(x,a)}_{x'}) = \max_{b \in \Sigma(\delta(x,a))} Q(\delta(x,a), b)$$

$$\boxed{Q(x,a) = g(x,a) + \gamma \max_{b \in \Sigma(\delta(x,a))} Q(\delta(x,a), b)}$$

For all possible actions we are looking for the optimal one, which determines the optimal control action

$$\mu(x) = \arg \max_{a \in \Sigma(x)} Q(x,a) = a^*$$

EX

$$Q(x,a) = g(x,a) + \gamma J^*(\delta(x,a))$$

$$Q(1,b) = 1 + \gamma J^*(2)$$

$$Q(1,c) = 2 + \gamma J^*(2)$$

$$Q(2,d) = 4 + \gamma J^*(1)$$

$$Q(2,e) = 3 + \gamma J^*(1)$$

# Optimal policy

$$\mu(1) = \max_{a \in \{b,c\}} Q(1,a) = c$$

$$\mu(2) = \max_{a \in \{d,e\}} Q(2,a) = d$$

## Model-free Q-function iteration

Replace $\delta(x,a)$ with $x'$

$\quad\quad\quad\quad \delta(x,a)$ with $r'$

Send an action to the plant
an wait for the next $x'$
and the transition reward $r'$

The Q-function is then updated without any model but instead by feedback from the plant.

$$Q(x,a) = r' + \gamma \max_{b \in \Sigma(x')} Q(x',b)$$

Problem: We don't know the Q-function.

## Q-learning

Based on $(x',r')$ an estimate $\hat{Q}_k(x,a)$ is updated

$$\hat{Q}_{k+1}(x,a) = (1-\alpha_k)\hat{Q}_k(x,a) +$$
$$\alpha_k \left( r' + \gamma \max_{b \in \Sigma(x')} \hat{Q}_k(x',b) \right)$$

$\alpha_k$ = learning factor that is reduced when time $k$ is increasing.

In most cases the action is selected to maximize $\hat{Q}_k(x,a)$. But to _explore_ the whole state space

(*) This strategy is called $\varepsilon$-greedy.

it is also necessary to take actions which do not maximize $\hat{Q}$.

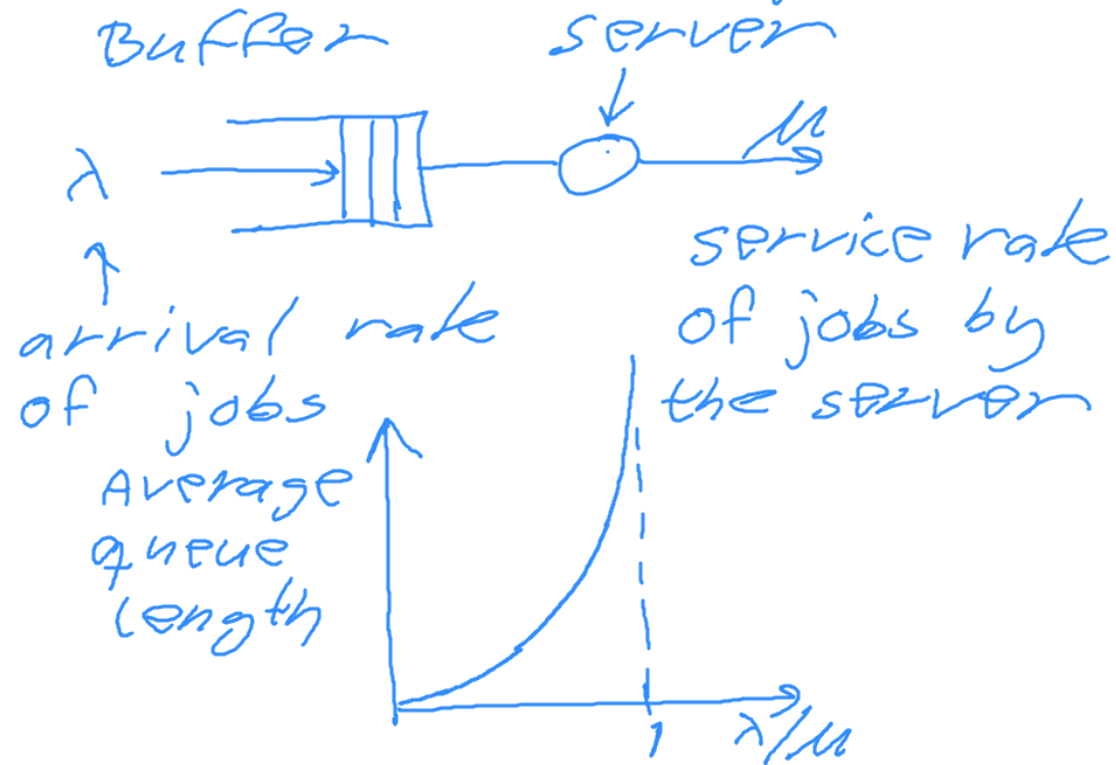An alternative arbitrary action $a \in \Sigma(x)$ is then taken with equal (uniform) but reduced probability
$$p_k(x) = \frac{1}{k|\Sigma(x)|} \text{ as time } k \text{ increases.}$$
This done for all $a \in \Sigma(x)$ and for all states $x$. The optimal (greedy) solution is then taken with prob. $1 - p_k(x) \to 1 \quad k \to \infty$ (*)

# 8. Markov processes (last part of Ch8)

Probabilistic models

Motivating example

Buffer      server

$\lambda$ → [buffer] → (server) → $\mu$

$\lambda$ ↑ arrival rate of jobs

service rate of jobs by the server

Average queue length

[graph: vertical axis is average queue length, horizontal axis is $\lambda/\mu$, curve rising steeply toward asymptote at $1$]

$1$    $\lambda/\mu$

# 8. Cont. Probabilistic models

## Markov chains

stochastic process: $\{x(t) : t \in T\}$

$x(t)$ = random variable for each $t \in T$

$T$ = countable set of time instances
$$= \{t_0, t_1, t_2, \ldots\}$$

For this discrete-time set $T$ we get a discrete-time stochastic process.

If instead $t \in \mathbb{R}^+ \Rightarrow \{x(t)\}$ is a continuous stochastic process.

## Markov chains

State space of $\{x(t)\}$ is a discrete set $Q$

$x(t)$ takes values $q_i \in Q$

State probability:
$$p_i(t_k) = P(x(t_k) = q_i)$$

This stochastic process is called a Markov chain if the next state conditional probability only depends on the current state

$$P(x(t_{k+1}) = q(t_{k+1}) \mid x(t_k) = q(t_k) \wedge$$
$$x(t_{k-1}) = q(t_{k-1}) \wedge \ldots \wedge x(t_0) = q(t_0)) =$$
$$= P(x(t_{k+1}) = q(t_{k+1}) \mid x(t_k) = q(t_k))$$

Transition probability:

$$p_{ij} = P(x(t_{k+1}) = q_j \mid x(t_k) = q_i)$$

$$P\{B \mid A\} = \frac{P[A \cap B]}{P[A]} = \begin{bmatrix} \text{probability of} \\ \text{B when A has} \\ \text{already occured} \end{bmatrix}$$

$$A_i = x(t_k) = q_i \qquad B = x(t_{k+1}) = q_j$$

Total probability

$$P\{B\} = \sum_{i=1}^{n} P\{A_i \cap B\} = \sum_{i=1}^{n} P\{B \mid A_i\} P[A_i]$$

Total probability =
state probability

$$p_j(t_{k+1}) = P\{\underbrace{x(t_{k+1}) = q_j}_{B}\} =$$

$$= \sum_{i=1}^{n} P\{\underbrace{x(t_{k+1}) = q_j}_{B} \mid \underbrace{x(t_k) = q_i}_{A_i}\} P[\underbrace{x(t_k) = q_i}_{A_i}]$$

$$= \sum_{i=1}^{n} p_{ij} \cdot p_i(t_k) =$$

$$= \underbrace{[p_1(t_k) \ldots p_n(t_k)]}_{P(t_k)} \begin{bmatrix} p_{1j} \\ \vdots \\ p_{nj} \end{bmatrix}$$

For $j = 1, \ldots, n$:

$$\underbrace{[p_1(t_{k+1}) \ldots p_n(t_{k+1})]}_{p(t_{k+1})} = p(t_k) \underbrace{\begin{bmatrix} p_{11} & p_{12} \ldots & p_{1n} \\ \vdots & & \vdots \\ p_{n1} & p_{n2} \ldots & p_{nn} \end{bmatrix}}_{P}$$
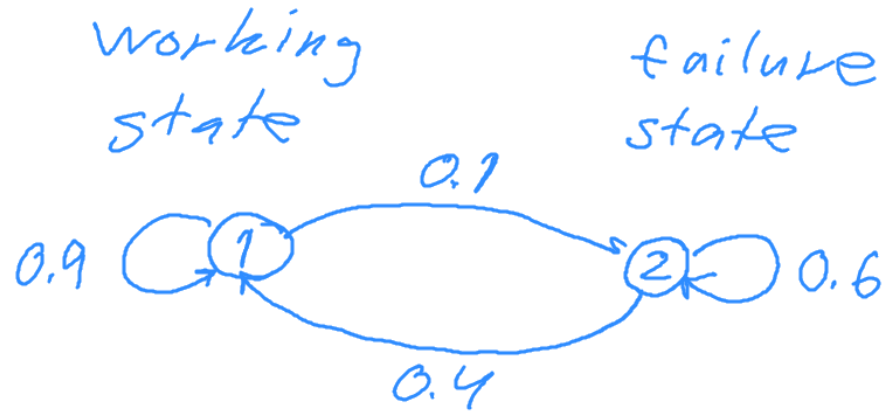
$$\boxed{p(t_{k+1}) = p(t_k)\, \mathbb{P}}$$

Important property:

$$\sum_{j=1}^{n} p_{ij} = \sum_{j=1}^{n} P\{x(t_{k+1}) = q_j \mid x(t_k) = q_i\} = 1$$

Ex Machine with failure

working state      failure state



$$\mathbb{P} = \begin{bmatrix} 0.9 & 0.1 \\ 0.4 & 0.6 \end{bmatrix}$$

$$p(0) = [1 \quad 0] \quad \left(\begin{array}{l}\text{initial state =}\\ \text{working state}\end{array}\right)$$

$$p(t_1) = p(0)\,\mathbb{P} = [1 \ 0]\begin{bmatrix} 0.9 & 0.1 \\ 0.4 & 0.6 \end{bmatrix} =$$

$$= [0.9 \ 0.1]$$

$$p(t_2) = [0.9 \ 0.1]\begin{bmatrix} 0.9 & 0.1 \\ 0.4 & 0.6 \end{bmatrix} = [0.85 \ 0.15]$$

stationary solution

$$\bar{p} = \lim_{k \to \infty} p(t_k) \qquad \bar{p} = \bar{p}\,\mathbb{P}$$

$$\bar{p} = [p \ \ 1-p] = [p \ \ 1-p]\begin{bmatrix} 0.9 & 0.1 \\ 0.4 & 0.6 \end{bmatrix}$$
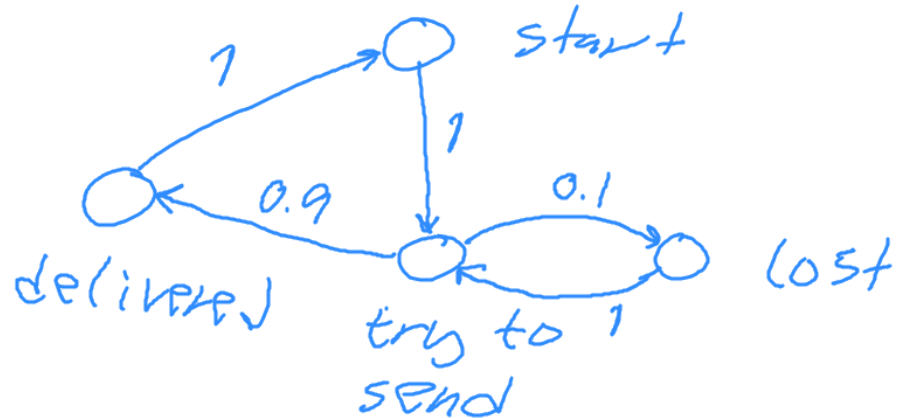
One unknown variable $p \Rightarrow$ enough to solve one equation

$$p = 0.9p + 0.4(1-p) = 0.4 + 0.5p$$

$$0.5p = 0.4 \Rightarrow p = 0.8$$

$$\bar{p} = [0.8 \quad 0.2]$$

Ex Communication protocol



# Markov processes

Continuous-time stochastic process with the Markov conditional probability property

$$t_{k+1} - t_k = \Delta t \quad \text{where} \quad \Delta t \to 0$$

Assume a given transition probability $\quad p_{ij} = a_{ij} \Delta t$

Here $a_{ij} =$ transition rate

The row sum in $\mathbb{P}$

$$\sum_{\ell=1}^{n} p_{j\ell} = \underbrace{p_{j1} + \dots + p_{jj-1} + p_{jj+1} \dots + p_{jn}}_{\sum_{\ell \neq j}^{n} p_{j\ell}} + \overset{=1}{p_{jj}}$$