

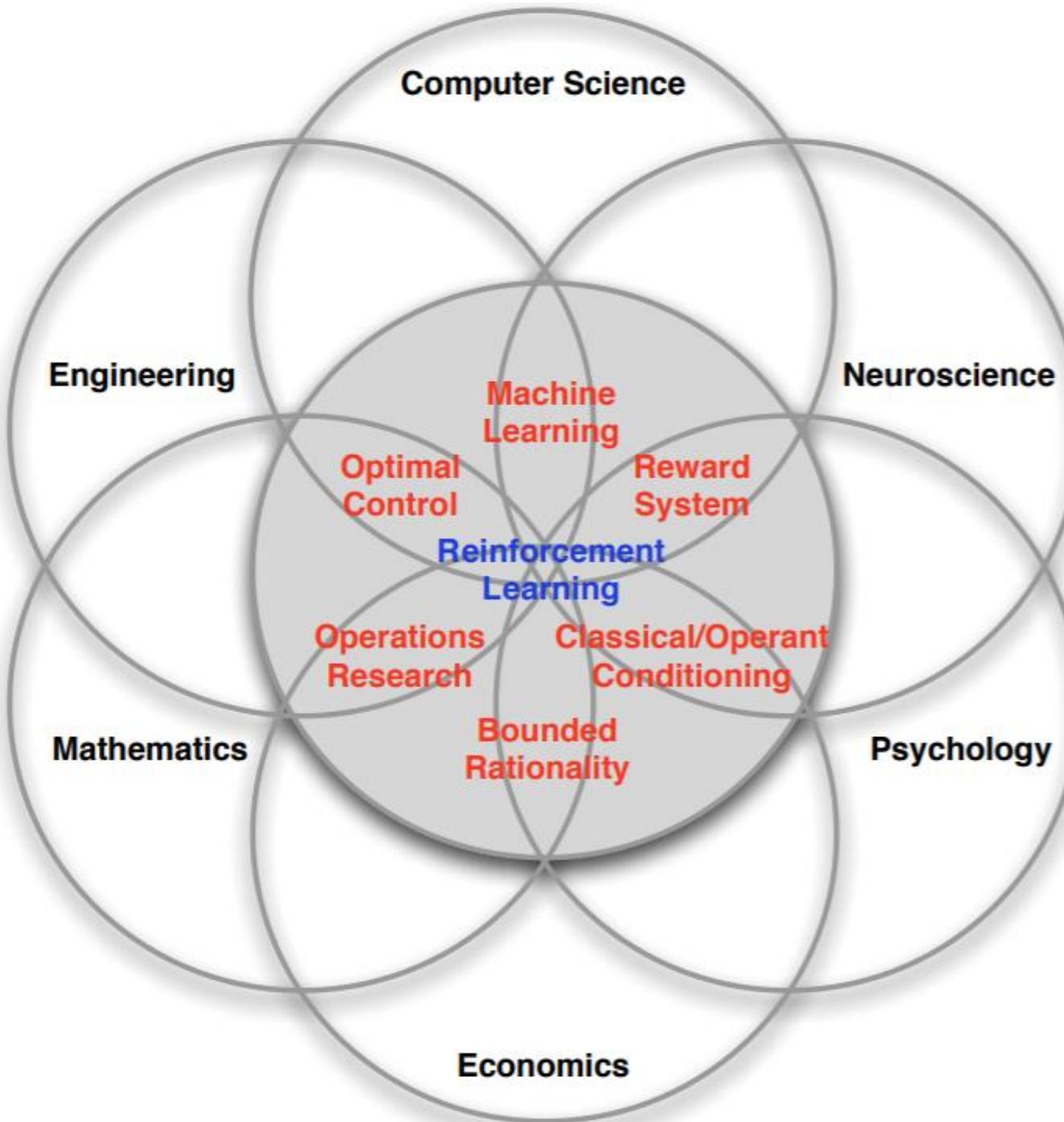
# Basics of Reinforcement Learning\*

by Alina Vereshchaka

\*Slides are based on David Silver's course "Reinforcement Learning" and the book by Sutton, Richard S., and Andrew G. Barto. "Reinforcement Learning: An introduction (2 edition)" (2018)

# Outline

- Markov Decision Process Definition
- RL Problem Definition
- Dynamic Programming
- Monte-Carlo Learning
- Temporal-Difference Learning
- SARSA
- Q-Learning
- n-Step Prediction



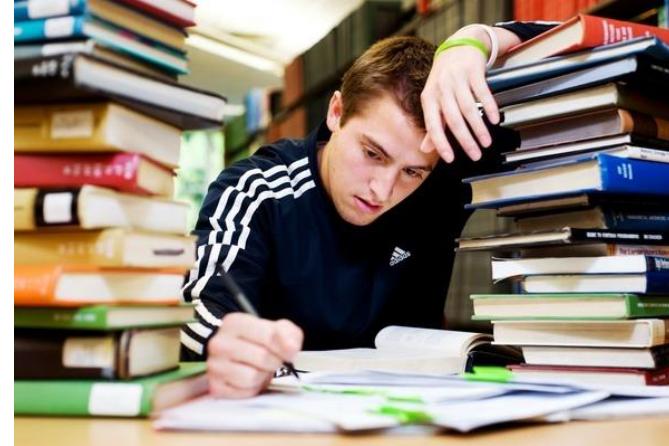
# What is Reinforcement learning?



# Reinforcement Learning in real world



**State:** tasty smell, sound  
**Actions:** give a paw  
**Reward:** food

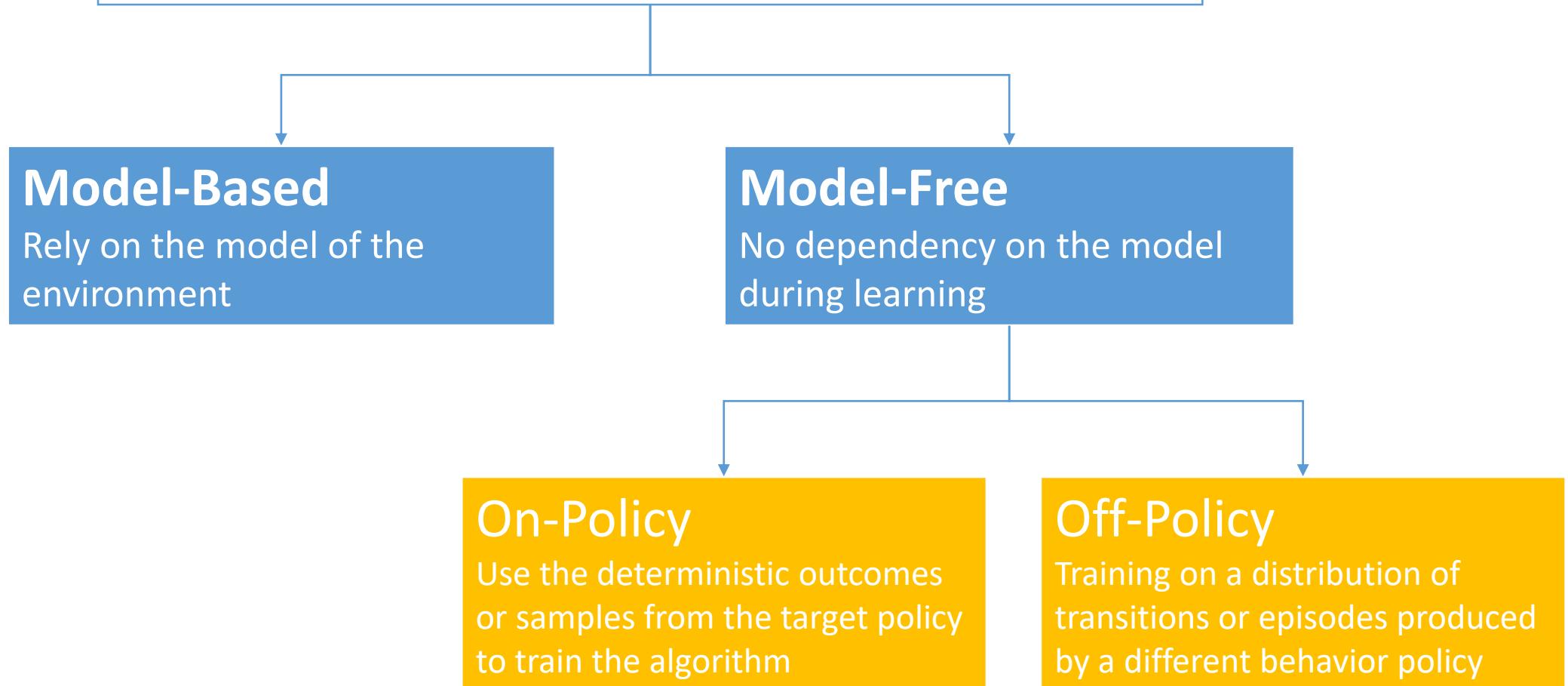


**State:** want a good grade  
**Actions:** studying  
**Reward:** good grade



**State:** opponent movement, current board state  
**Actions:** make a move  
**Reward:** points

# Reinforcement Learning Algorithms



# Markov Chain

## Definition

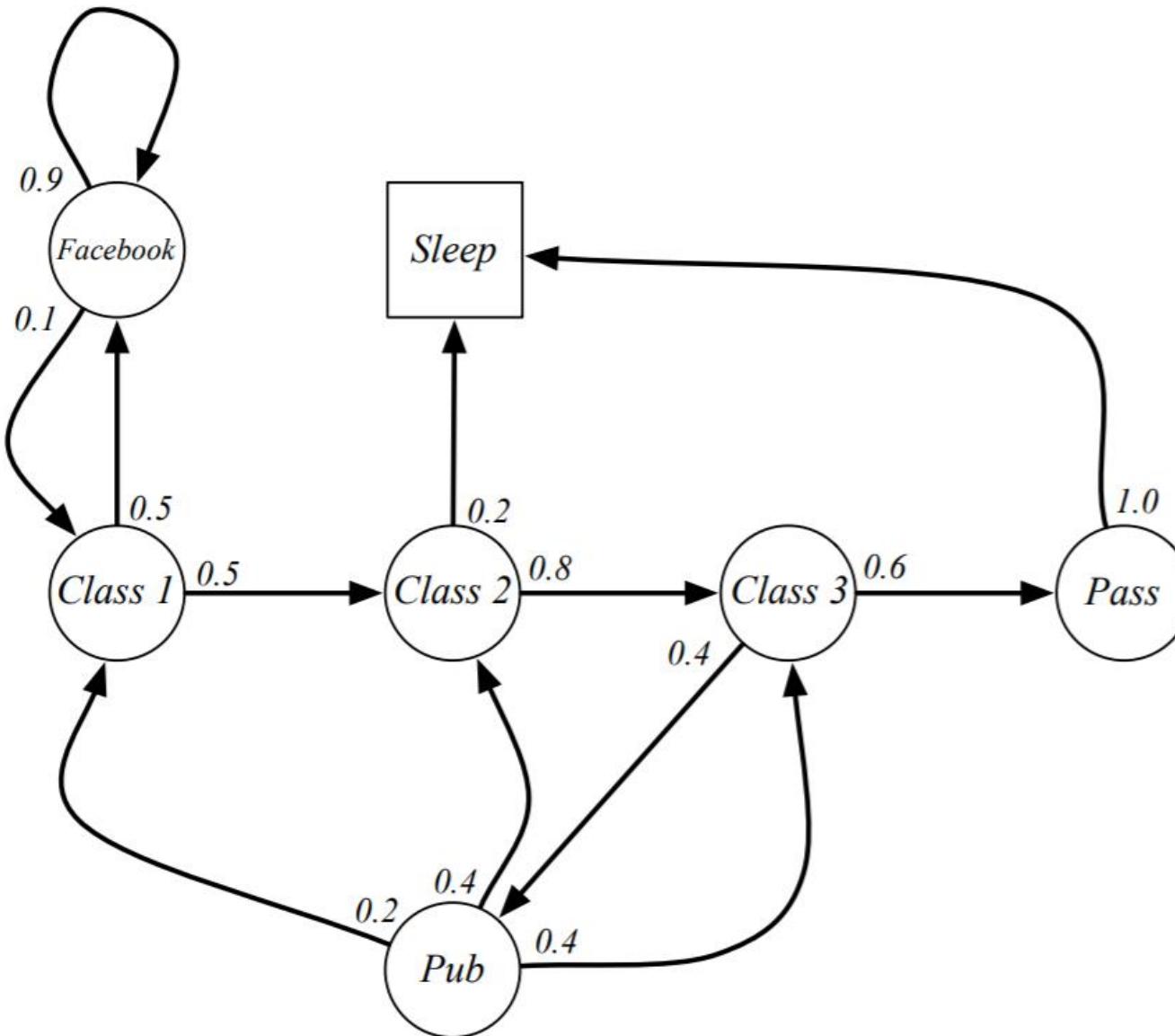
A *Markov Process* (or *Markov Chain*) is a tuple  $\langle \mathcal{S}, \mathcal{P} \rangle$

- $\mathcal{S}$  is a (finite) set of states
- $\mathcal{P}$  is a state transition probability matrix,  
 $\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' | S_t = s]$

**Markovian property:**  $\mathbb{P}[S_{t+1} | S_t] = \mathbb{P}[S_{t+1} | S_1, \dots, S_t]$

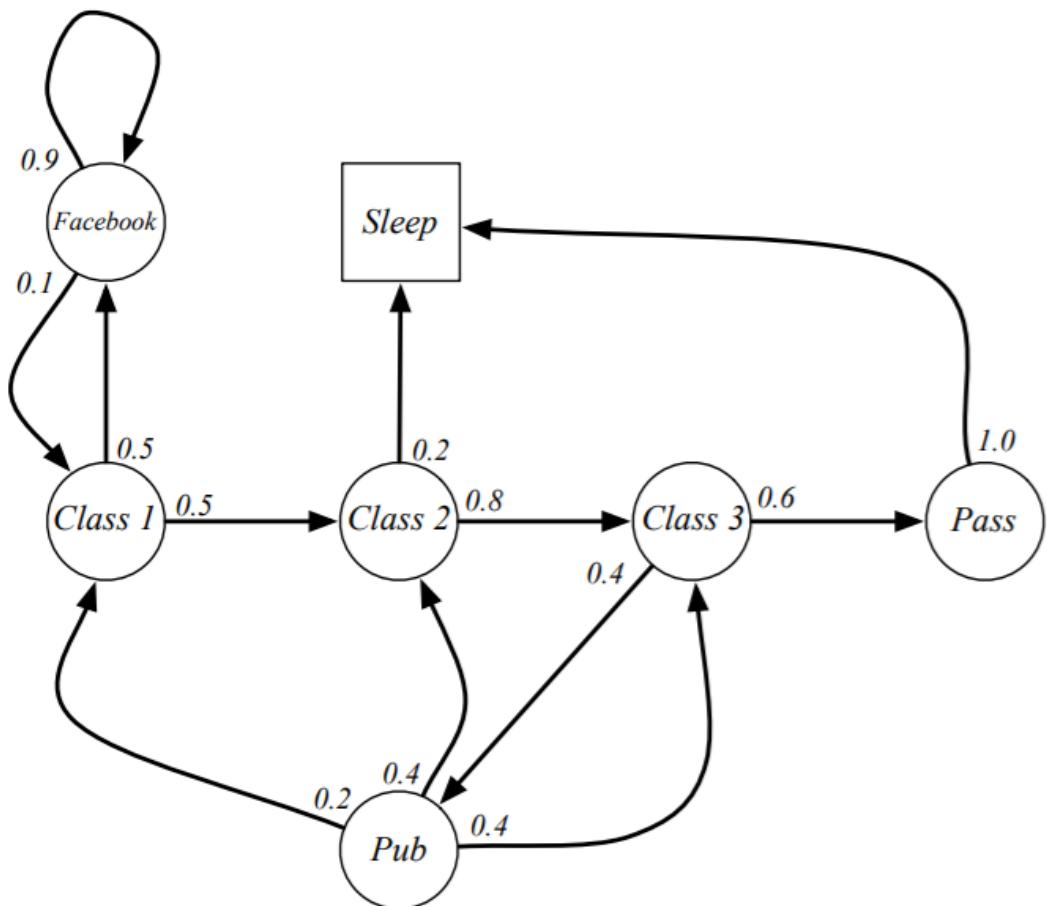
**“The future is independent of the past given the present”**

# Example: Student Markov Chain



# Example: Student Markov Chain Episodes

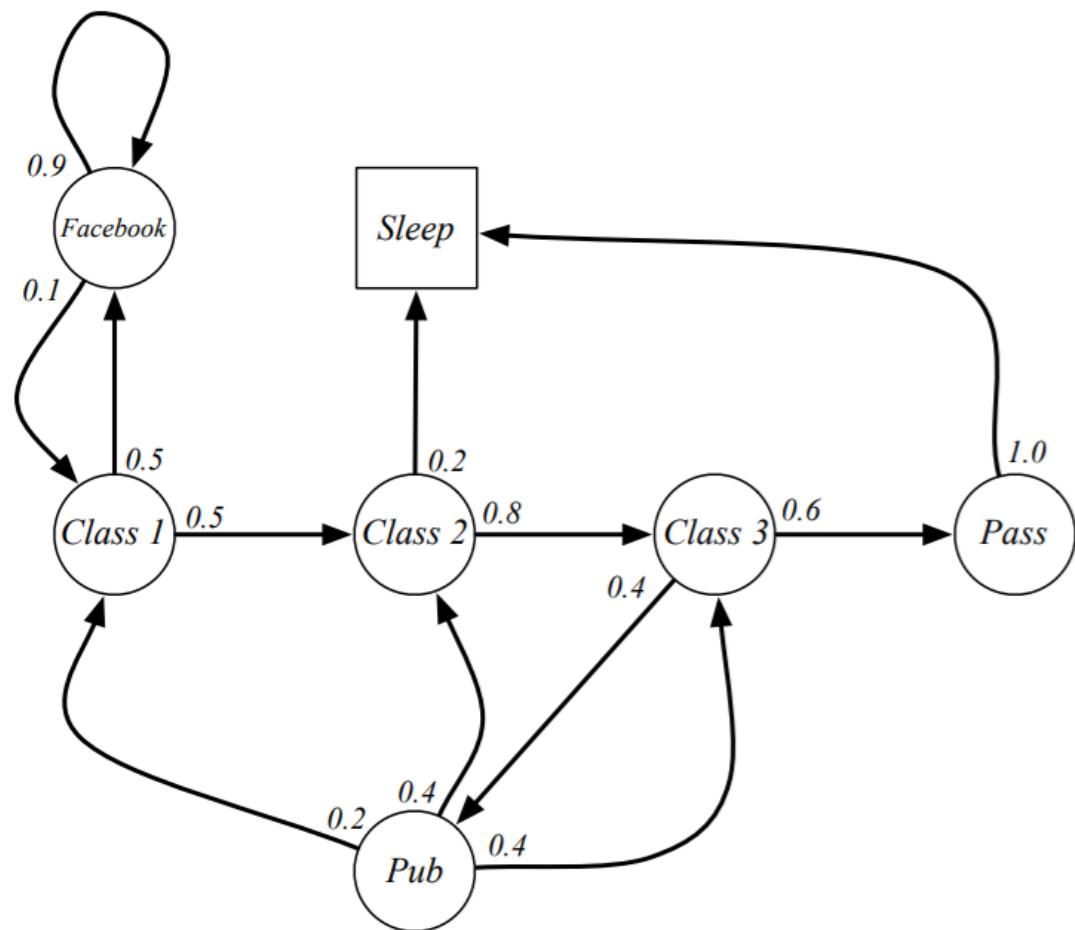
Sample **episodes** for Student Markov Chain starting from  $S_1 = C1$



$S_1, S_2, \dots, S_T$

- C1 C2 C3 Pass Sleep
- C1 FB FB C1 C2 Sleep
- C1 C2 C3 Pub C2 C3 Pass Sleep
- C1 FB FB C1 C2 C3 Pub C1 FB FB FB C1 C2 C3 Pub C2 Sleep

# Example: Transition Matrix



$$\mathcal{P} = \begin{bmatrix} C1 & C2 & C3 & Pass & Pub & FB & Sleep \\ C1 & 0.5 & 0.8 & 0.6 & 0.4 & 0.1 & 0.2 \\ C2 & 0.2 & 0.4 & 0.4 & 0.6 & 0.4 & 0.9 \\ C3 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ Pass & 0.0 & 0.0 & 0.0 & 1.0 & 0.0 & 0.0 \\ Pub & 0.0 & 0.0 & 0.0 & 0.0 & 1.0 & 0.0 \\ FB & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 1.0 \\ Sleep & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 & 0.0 \end{bmatrix}$$

# Markov Reward Process (MRP)

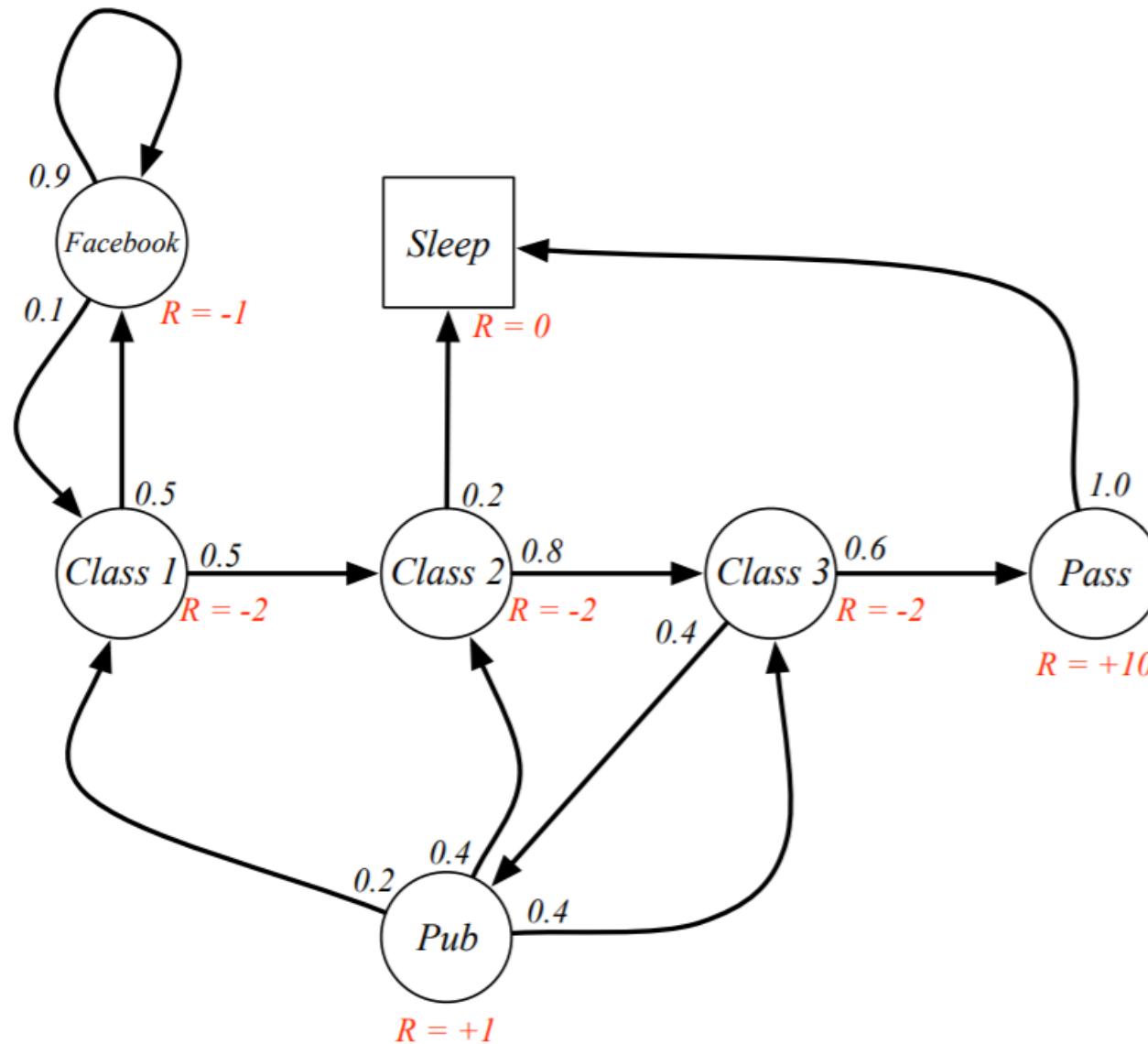
A Markov reward process is a Markov chain with values.

## Definition

A *Markov Reward Process* is a tuple  $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

- $\mathcal{S}$  is a finite set of states
- $\mathcal{P}$  is a state transition probability matrix,  
 $\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' \mid S_t = s]$
- $\mathcal{R}$  is a reward function,  $\mathcal{R}_s = \mathbb{E}[R_{t+1} \mid S_t = s]$
- $\gamma$  is a discount factor,  $\gamma \in [0, 1]$

# Example: Student MRP



# Discounting factor $\gamma$

The **discounting factor** ( $\gamma \in [0, 1]$ ) penalize the rewards in the future.  
Reward at time  $k$  worth only  $\gamma^{k-1}$

## Motivation:

- The future rewards may have higher uncertainty (stock market)
- The future rewards do not provide immediate benefits (As human beings, we might prefer to have fun today rather than 5 years later ;)
- Discounting provides mathematical convenience (we don't need to track future steps infinitely to compute return)
- It is sometimes possible to use undiscounted Markov reward processes (i.e.  $\gamma = 1$ ), e.g. if all sequences terminate.

# Return ( $G_t$ )

## Definition

The *return*  $G_t$  is the total discounted reward from time-step  $t$ .

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- The *discount*  $\gamma \in [0, 1]$  is the present value of future rewards
- The value of receiving reward  $R$  after  $k + 1$  time-steps is  $\gamma^k R$ .
- This values immediate reward above delayed reward.
  - $\gamma$  close to 0 leads to "myopic" evaluation
  - $\gamma$  close to 1 leads to "far-sighted" evaluation

# Bellman Equation for Return

**Expected reward** is a total sum of discounted rewards going forward.

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \\ &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \\ &= R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$

## Example: Returns for Student MRP

Sample **returns** for Student MRP:

Starting from  $S_1 = C1$  with  $\gamma = \frac{1}{2}$

$$G_1 = R_2 + \gamma R_3 + \dots + \gamma^{T-2} R_T$$

C1 C2 C3 Pass Sleep	$v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 10 * \frac{1}{8}$	=	-2.25
C1 FB FB C1 C2 Sleep	$v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16}$	=	-3.125
C1 C2 C3 Pub C2 C3 Pass Sleep	$v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 1 * \frac{1}{8} - 2 * \frac{1}{16} \dots$	=	-3.41
C1 FB FB C1 C2 C3 Pub C1 ...	$v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16} \dots$	=	-3.20
FB FB FB C1 C2 C3 Pub C2 Sleep			

# Value Function $v(s)$

The value function  $v(s)$  gives the long-term value of state  $s$

## Definition

The *state value function*  $v(s)$  of an MRP is the expected return starting from state  $s$

$$v(s) = \mathbb{E}[G_t \mid S_t = s]$$

# Bellman Equation for MRPs

The value function can be decomposed into two parts:

- immediate reward  $R_{t+1}$
- discounted value of successor state  $\gamma v(S_{t+1})$

$$\begin{aligned}v(s) &= \mathbb{E}[G_t \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma G_{t+1} \mid S_t = s] \\&= \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s]\end{aligned}$$

# Markov Decision Process

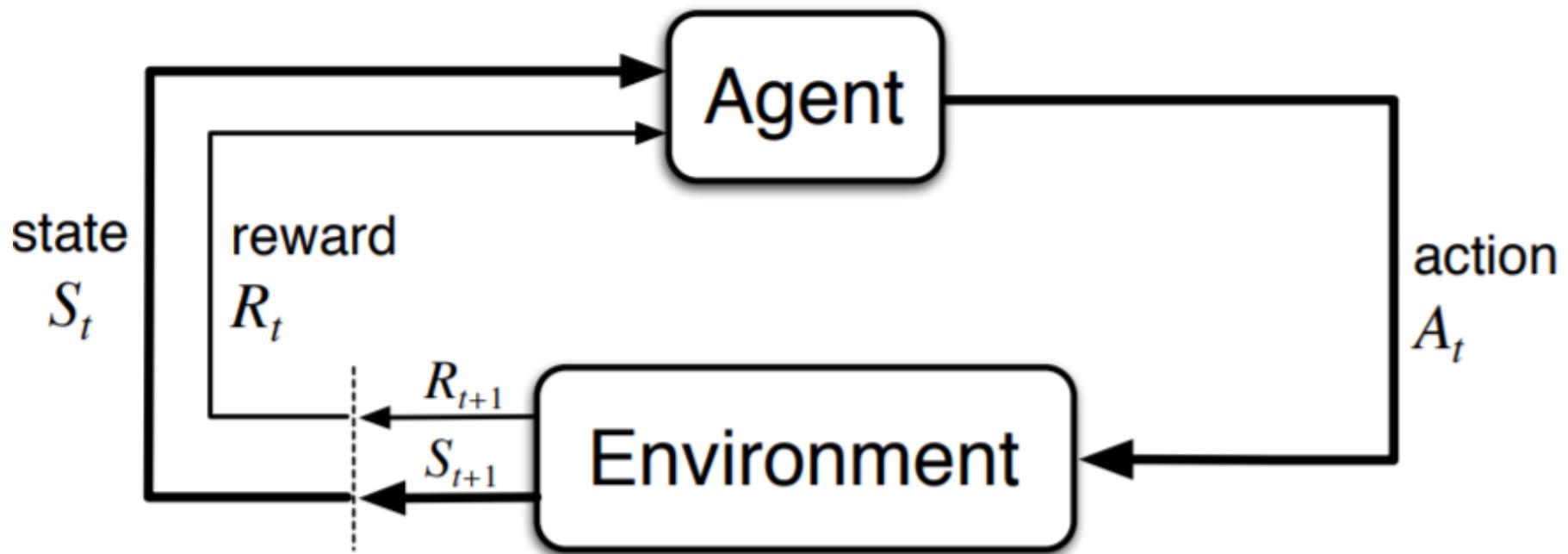
A Markov decision process (MDP) is a Markov reward process with decisions. It is an *environment* in which all states are Markov.

## Definition

A *Markov Decision Process* is a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$

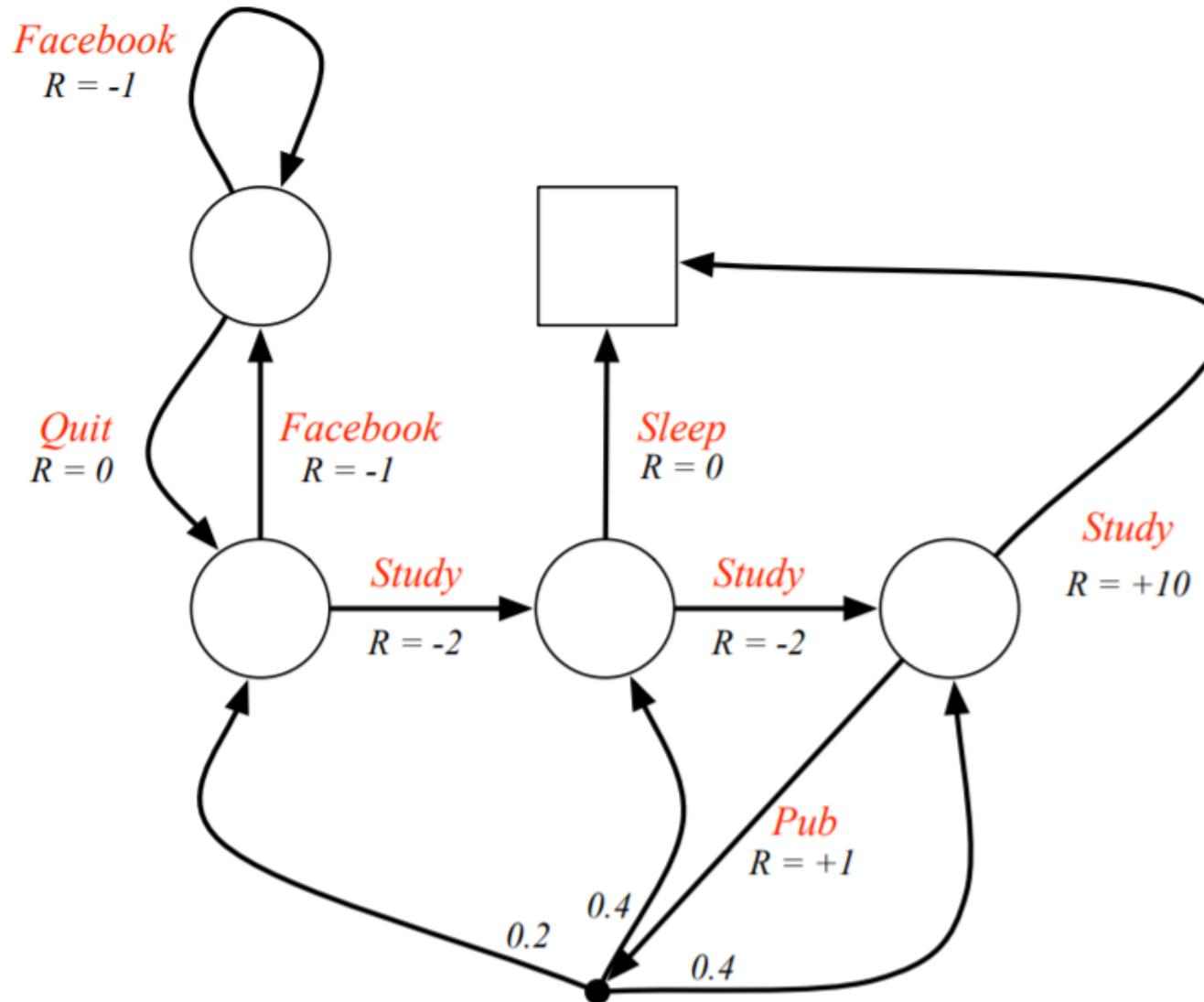
- $\mathcal{S}$  is a finite set of states
- $\mathcal{A}$  is a finite set of actions
- $\mathcal{P}$  is a state transition probability matrix,  
$$\mathcal{P}_{ss'}^{\textcolor{red}{a}} = \mathbb{P}[S_{t+1} = s' \mid S_t = s, A_t = \textcolor{red}{a}]$$
- $\mathcal{R}$  is a reward function,  $\mathcal{R}_s^{\textcolor{red}{a}} = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = \textcolor{red}{a}]$
- $\gamma$  is a discount factor  $\gamma \in [0, 1]$ .

# Finite Markov Decision Processes (MDP)



**Daily life trajectory:**  $S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, S_3, A_3, \dots, S_T$

# Example: Student MDP



# Policy ( $\pi$ )

## Definition

A *policy*  $\pi$  is a distribution over actions given states,

$$\pi(a|s) = \mathbb{P}[A_t = a \mid S_t = s]$$

- A policy fully defines the behaviour of an agent
- MDP policies depend on the current state (not the history)
- i.e. Policies are *stationary* (time-independent),  
 $A_t \sim \pi(\cdot|S_t), \forall t > 0$

# Value Functions $v_\pi(s)$ , $q_\pi(s, a)$

## Definition

The *state-value function*  $v_\pi(s)$  of an MDP is the expected return starting from state  $s$ , and then following policy  $\pi$

$$v_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s]$$

## Definition

The *action-value function*  $q_\pi(s, a)$  is the expected return starting from state  $s$ , taking action  $a$ , and then following policy  $\pi$

$$q_\pi(s, a) = \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a]$$

# Bellman Expectation Equation

The state-value function can again be decomposed into immediate reward plus discounted value of successor state,

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

The action-value function can similarly be decomposed,

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

# Optimal Value Function

## Definition

The *optimal state-value function*  $v_*(s)$  is the maximum value function over all policies

$$v_*(s) = \max_{\pi} v_{\pi}(s)$$

The *optimal action-value function*  $q_*(s, a)$  is the maximum action-value function over all policies

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$$

# Finding an Optimal Policy

An optimal policy can be found by maximising over  $q_*(s, a)$ ,

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \underset{a \in \mathcal{A}}{\operatorname{argmax}} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

- There is always a deterministic optimal policy for any MDP
- If we know  $q_*(s, a)$ , we immediately have the optimal policy

# POMDPs

A Partially Observable Markov Decision Process is an MDP with hidden state. It is a hidden Markov model with actions.

## Definition

A *POMDP* is a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{P}, \mathcal{R}, \mathcal{Z}, \gamma \rangle$

- $\mathcal{S}$  is a finite set of states
- $\mathcal{A}$  is a finite set of actions
- $\mathcal{O}$  is a finite set of observations
- $\mathcal{P}$  is a state transition probability matrix,  
$$\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' \mid S_t = s, A_t = a]$$
- $\mathcal{R}$  is a reward function,  $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$
- $\mathcal{Z}$  is an observation function,  
$$\mathcal{Z}_{s'o}^a = \mathbb{P}[O_{t+1} = o \mid S_{t+1} = s', A_t = a]$$
- $\gamma$  is a discount factor  $\gamma \in [0, 1]$ .

# MDP: Definitions

Episode  $(s, a, s', r)$  - agent in state  $s$  takes action  $a$  to arrive in the next state  $s'$  and obtain reward  $r$

**Transition probability function  $P(s', r|s, a)$ :**

$$P(s', r | s, a) = \mathbb{P}[S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a]$$

$$\sum_{s' \in \mathcal{S}} \sum_{r \in \mathcal{R}} p(s', r | s, a) = 1$$

**State-transition function  $P(s'|s, a)$ :**

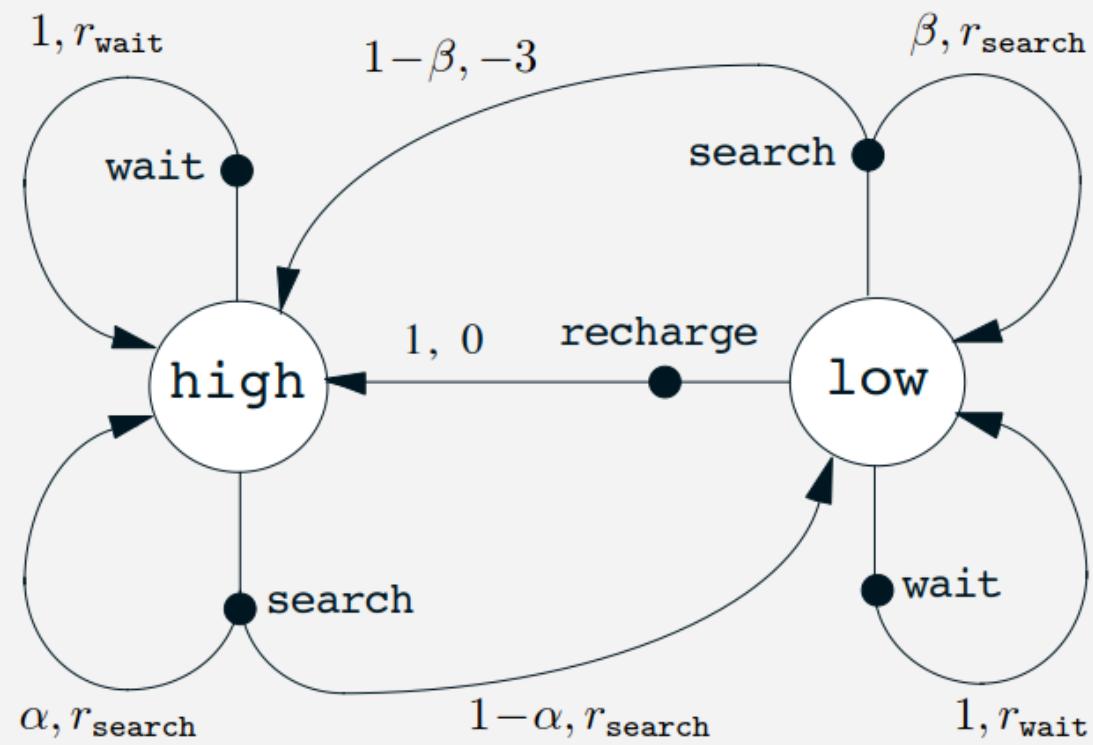
$$P(s'|s, a) = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} P(s', r | s, a)$$

**Expected reward function  $R(s, a)$ :**

$$R(s, a) = E [R_{t+1} | S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r \sum_{s' \in \mathcal{S}} P(s', r | s, a)$$

# Case study: Recycling Robot

$s$	$a$	$s'$	$p(s'   s, a)$	$r(s, a, s')$
high	search	high	$\alpha$	$r_{\text{search}}$
high	search	low	$1 - \alpha$	$r_{\text{search}}$
low	search	high	$1 - \beta$	$-3$
low	search	low	$\beta$	$r_{\text{search}}$
high	wait	high	1	$r_{\text{wait}}$
high	wait	low	0	-
low	wait	high	0	-
low	wait	low	1	$r_{\text{wait}}$
low	recharge	high	1	0
low	recharge	low	0	-



# Policy

- A policy is the agent's behavior
- It is a map from state to action
- It can be either deterministic or stochastic:

**Deterministic:**

$$\pi(s) = a$$

**Stochastic:**

$$\pi(a|s) = \mathbb{P}_\pi[A = a | S = s]$$

# Action-value method (Q)

$$Q_t(a) \doteq \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t} = \frac{\sum_{i=1}^{t-1} R_i \cdot \mathbb{1}_{A_i=a}}{\sum_{i=1}^{t-1} \mathbb{1}_{A_i=a}}$$

As the denominator goes to infinity, by the law of large numbers,  $Q_t(a)$  converges to  $q^*(a)$ - *sample-average method* for estimating action values because each estimate is an average of the sample of relevant rewards.

**Greedy action:**  $A_t \doteq \arg \max_a Q_t(a)$

# Definitions

***State-value* function:**

$$V_\pi(s) \doteq \mathbb{E}_\pi[G_t | S_t = s]$$
$$= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right], \text{ for all } s \in S$$

**Action-value function:**

$$Q_\pi(s, a) \doteq \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$$
$$= \mathbb{E}_\pi \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

State-value function w.r.t. probability distribution over possible actions and the Q-values:

$$V_\pi(s) = \sum_{a \in \mathcal{A}} Q_\pi(s, a) \pi(a | s)$$

**Action advantage function:**

$$A_\pi(s, a) = \bar{Q}_\pi(s, a) - V_\pi(s)$$

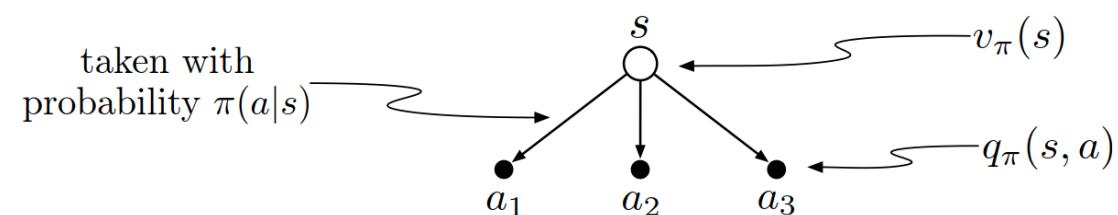
# Bellman Equation for state-value function

Bellman equation decomposes the value function into the immediate reward plus the discounted future values.

$$\begin{aligned} V_\pi(s) &\doteq \mathbb{E}_\pi[G_t | S_t = s] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] \\ &= \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) \left[ r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s'] \right] \\ &= \sum_a \pi(a|s) \sum_{s',r} p(s', r | s, a) \left[ r + \gamma V_\pi(s') \right], \text{ for all } s \in S \end{aligned}$$

Similarly for Q-value:

$$Q(s, a) = \mathbb{E} [R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$



# Optimal Value and Policy

- The optimal value function produces the maximum return:

$$V_*(s) = \max_{\pi} V_{\pi}(s)$$

$$Q_*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

- The optimal policy achieves optimal value functions:

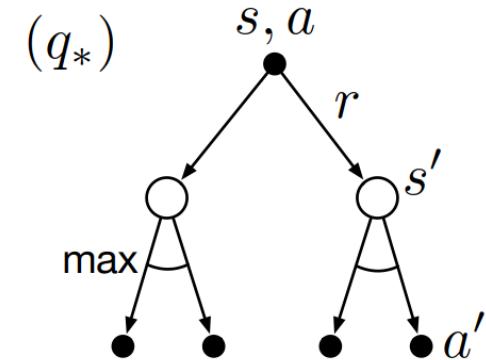
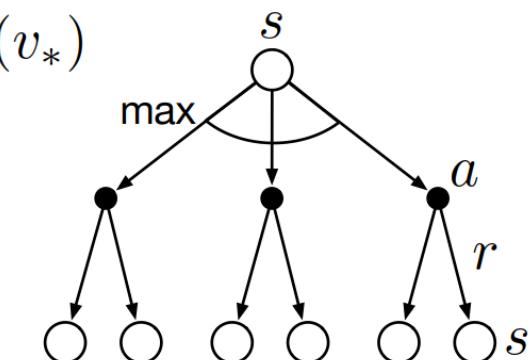
$$\pi_* = \arg \max_{\pi} V_{\pi}(s)$$

$$\pi_* = \arg \max_{\pi} Q_{\pi}(s, a)$$

- Thus we have

$$V_{\pi_*}(s) = V_*(s)$$

$$Q_{\pi_*}(s, a) = Q_*(s, a)$$



# Bellman optimality equations

**State-Value:**

$$\begin{aligned} V_\pi(s) &= \max_{a \in A(s)} Q_{\pi_*}(s, a) \\ &= \max_a \mathbb{E}_{\pi_*}[G_t | S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \\ &= \max_a \mathbb{E}_{\pi_*}[R_{t+1} + \gamma V_*(S_{t+1}) | S_t = s, A_t = a] \\ &= \max_a \sum_{s',r} p(s', r | s, a) [r + \gamma V_*(s')] \end{aligned}$$

**State-Action:**

$$\begin{aligned} Q_\pi(s, a) &= \mathbb{E} \left[ R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = s, A_t = a \right] \\ &= \sum_{s',r} p(s', r, a) [r + \gamma Q_*(s', a')] \end{aligned}$$

# Dynamic Programming (DP)

- DP is a model-based approach
- A method for solving complex problems, by breaking them down into subproblems:
  - Solve the subproblems
  - Combine solutions to subproblems
- Iteratively evaluates value functions and improving policy following Bellman equations.
- Assumes full knowledge of the MDP
- Bellman equations give recursive decomposition

# DP: Iterative Policy Evaluation

- Problem: evaluate a given policy  $\pi$
- Solution: iterative application of Bellman expectation backup
- $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_\pi$
- Using *synchronous* backups,
  - At each iteration  $k + 1$
  - For all states  $s \in \mathcal{S}$
  - Update  $v_{k+1}(s)$  from  $v_k(s')$
  - where  $s'$  is a successor state of  $s$

# DP: Policy Evaluation & Improvement

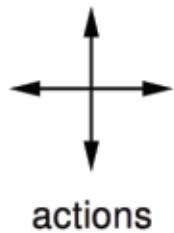
**Iterative policy evaluation:**

$$\begin{aligned} V_{k+1}(s) &\doteq \mathbb{E}_\pi[R_{t+1} + \gamma V_k(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s',r} P(s',r|s,a) (r + \gamma V_k(s')) \end{aligned}$$

**Policy Improvement:**

$$\begin{aligned} Q_\pi(s, a) &\doteq \mathbb{E}[R_{t+1} + \gamma V_\pi(S_{t+1}) | S_t = s, A_t = a] \\ &= \sum_{s',r} p(s',r|s,a) \left[ r + \gamma V_\pi(s') \right] \end{aligned}$$

# DP in Small Gridworld



	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$r = -1$   
on all transitions

- Undiscounted episodic MDP ( $\gamma = 1$ )
- Nonterminal states 1, ..., 14
- One terminal state (shown twice as shaded squares)
- Actions leading out of the grid leave state unchanged
- Reward is  $-1$  until the terminal state is reached
- Agent follows uniform random policy

$$\pi(n|\cdot) = \pi(e|\cdot) = \pi(s|\cdot) = \pi(w|\cdot) = 0.25$$

# Convergence of iterative policy

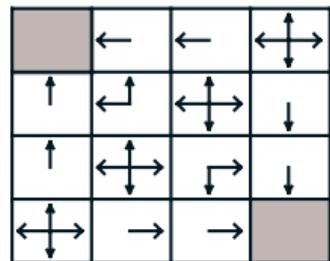
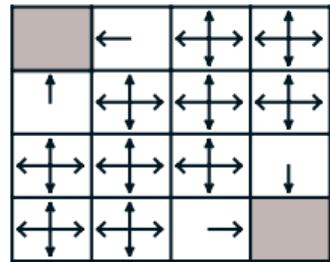
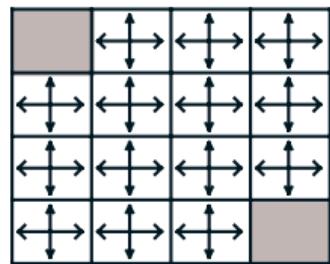
$v_k$  for the random policy

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0

greedy policy  
w.r.t.  $v_k$



$k = 3$

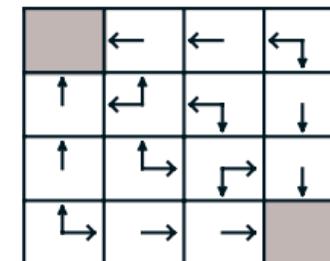
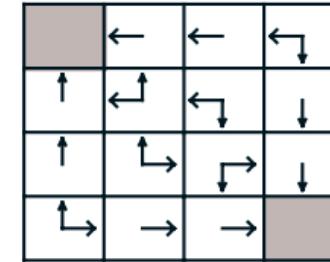
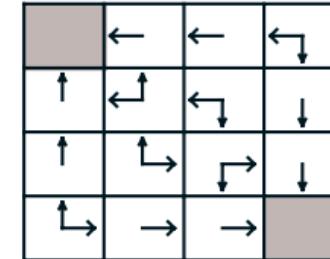
0.0	-2.4	-2.9	-3.0
-2.4	-2.9	-3.0	-2.9
-2.9	-3.0	-2.9	-2.4
-3.0	-2.9	-2.4	0.0

$k = 10$

0.0	-6.1	-8.4	-9.0
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9.0	-8.4	-6.1	0.0

$k = \infty$

0.0	-14.	-20.	-22.
-14.	-18.	-20.	-20.
-20.	-20.	-18.	-14.
-22.	-20.	-14.	0.0



optimal policy

# How to Improve a Policy

- Given a policy  $\pi$ 
  - Evaluate the policy  $\pi$

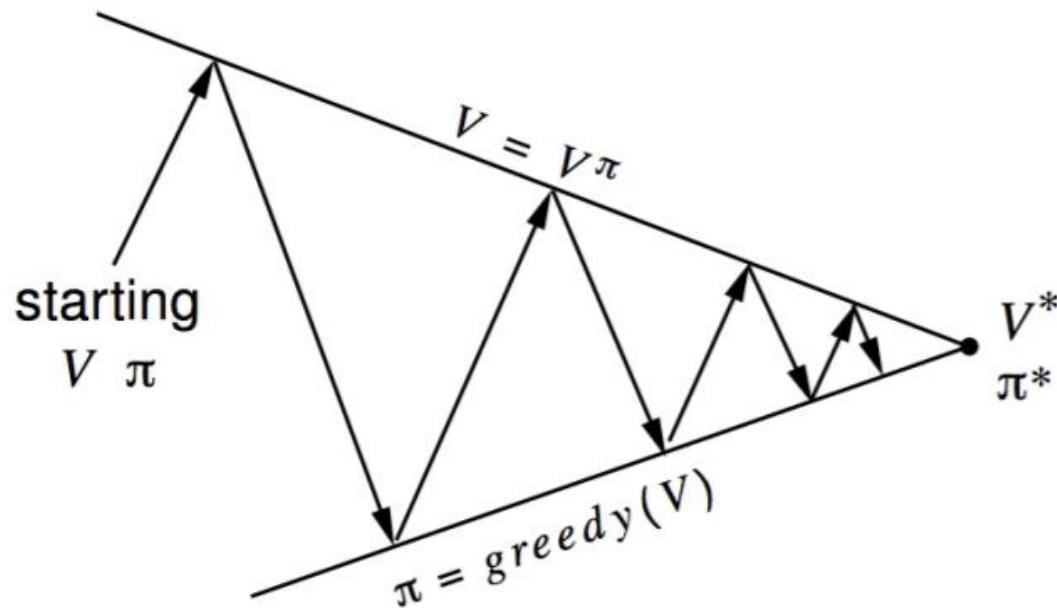
$$v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | S_t = s]$$

- Improve the policy by acting greedily with respect to  $v_\pi$

$$\pi' = \text{greedy}(v_\pi)$$

- In Small Gridworld improved policy was optimal,  $\pi' = \pi^*$
- In general, need more iterations of improvement / evaluation
- But this process of **policy iteration** always converges to  $\pi^*$

# Policy Iteration

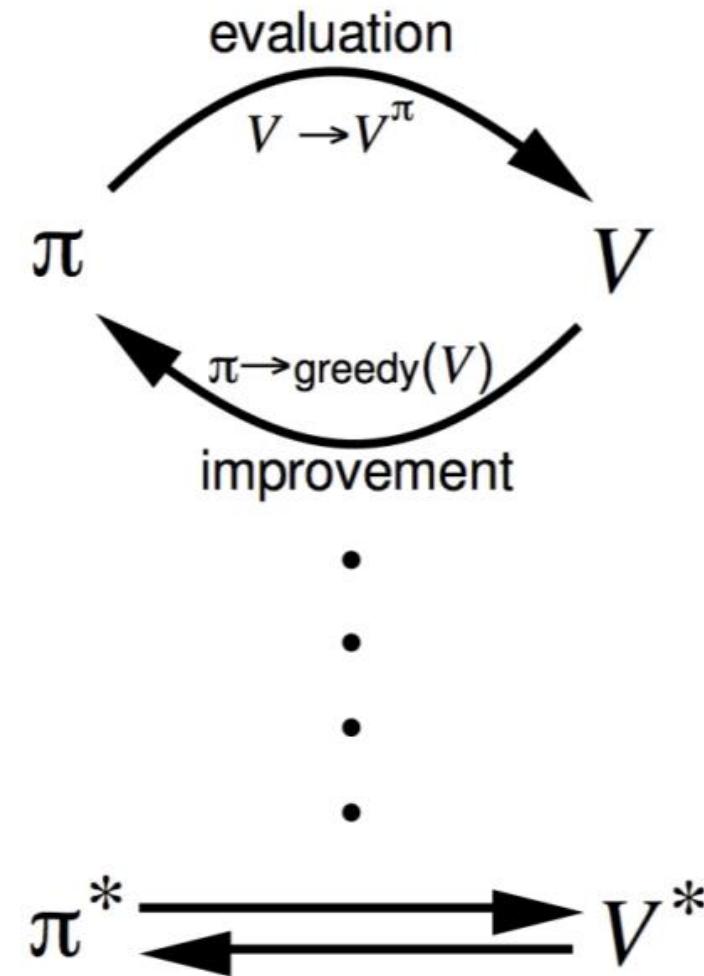


Policy evaluation Estimate  $v_\pi$

Iterative policy evaluation

Policy improvement Generate  $\pi' \geq \pi$

Greedy policy improvement



# Policy Improvement

- Consider a deterministic policy,  $a = \pi(s)$
- We can *improve* the policy by acting greedily

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} q_\pi(s, a)$$

- This improves the value from any state  $s$  over one step,

$$q_\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_\pi(s, a) \geq q_\pi(s, \pi(s)) = v_\pi(s)$$

- It therefore improves the value function,  $v_{\pi'}(s) \geq v_\pi(s)$

$$\begin{aligned} v_\pi(s) &\leq q_\pi(s, \pi'(s)) = \mathbb{E}_{\pi'} [R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 q_\pi(S_{t+2}, \pi'(S_{t+2})) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \dots \mid S_t = s] = v_{\pi'}(s) \end{aligned}$$

# Policy Improvement

- If improvements stop,

$$q_{\pi}(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_{\pi}(s, a) = q_{\pi}(s, \pi(s)) = v_{\pi}(s)$$

- Then the Bellman optimality equation has been satisfied

$$v_{\pi}(s) = \max_{a \in \mathcal{A}} q_{\pi}(s, a)$$

- Therefore  $v_{\pi}(s) = v_*(s)$  for all  $s \in \mathcal{S}$
- so  $\pi$  is an optimal policy

# Model-Free Predictions

- Monte-Carlo Learning
- Temporal-Difference Learning
- SARSA
- Q-learning
- n-step Bootstrapping

# Monte-Carlo Method

- Monte-Carlo (MC) learns from episodes of experience
- MC is model-free: no knowledge of MDP transition/rewards
- MC learns from complete episodes  $S_1, A_1, R_2, \dots, S_T$  (all the episodes must eventually terminate): no bootstrapping
- MC uses simple idea: value = mean return

# Monte-Carlo Method

- Goal: learn  $v_\pi$  from episodes of experience under policy  $\pi$

$$S_1, A_1, R_2, \dots, S_k \sim \pi$$

- Recall that the *return* is the total discounted reward:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

- Recall that the value function is the expected return:

$$v_\pi(s) = \mathbb{E}_\pi [G_t \mid S_t = s]$$

- Monte-Carlo policy evaluation uses *empirical mean* return instead of *expected* return

# Monte-Carlo Method

**Empirical mean return for state  $s$ :**

$$V(s) = \frac{\sum_{t=1}^T 1[S_t = s]G_t}{\sum_{t=1}^T 1[S_t = s]}$$

**Action-value function:**

$$Q(s, a) = \frac{\sum_{t=1}^T 1[S_t = s, A_t = a]G_t}{\sum_{t=1}^T 1[S_t = s, A_t = a]}$$

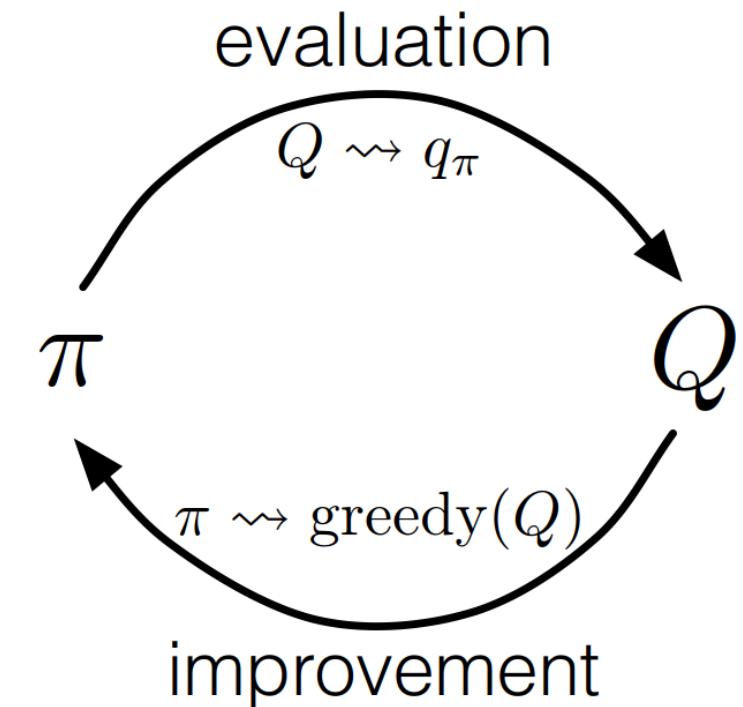
# Learning the optimal policy by MC

1. Improve the policy greedily with respect to the current value function:

$$\pi'(s) \arg \max_{a \in \mathcal{A}} Q(s, a)$$

2. Generate a new episode with the new policy  $\pi'$  (using  $\epsilon$ -greedy algorithm)
3. Estimate  $Q$  using the new episode:

$$q_\pi(s, a) = \frac{\sum_{t=1}^T \left( \mathbf{1}[S_t=s, A_t=a] \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1} \right)}{\sum_{t=1}^T \mathbf{1}[S_t=s, A_t=a]}$$



# Temporal-Difference(TD) Learning

- Temporal-Difference (TD) Learning is model-free and learns from episodes of experience.
- TD is *model-free*: no knowledge of MDP transition/rewards
- TD learning can learn from **incomplete** episodes, by bootstrapping.
- Updates targets with regard to existing estimates rather than exclusively relying on actual rewards and complete returns

# MC and TD

- Goal: learn  $v_\pi$  online from experience under policy  $\pi$
- Incremental every-visit Monte-Carlo
  - Update value  $V(S_t)$  toward *actual* return  $G_t$

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$

- Simplest temporal-difference learning algorithm: TD(0)
  - Update value  $V(S_t)$  toward *estimated* return  $R_{t+1} + \gamma V(S_{t+1})$

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

- $R_{t+1} + \gamma V(S_{t+1})$  is called the *TD target*
- $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$  is called the *TD error*

# TD Algorithm

## Tabular TD(0) for estimating $v_\pi$

Input: the policy  $\pi$  to be evaluated

Initialize  $V(s)$  arbitrarily (e.g.,  $V(s) = 0$ , for all  $s \in \mathcal{S}^+$ )

Repeat (for each episode):

    Initialize  $S$

    Repeat (for each step of episode):

$A \leftarrow$  action given by  $\pi$  for  $S$

        Take action  $A$ , observe  $R, S'$

$V(S) \leftarrow V(S) + \alpha [R + \gamma V(S') - V(S)]$

$S \leftarrow S'$

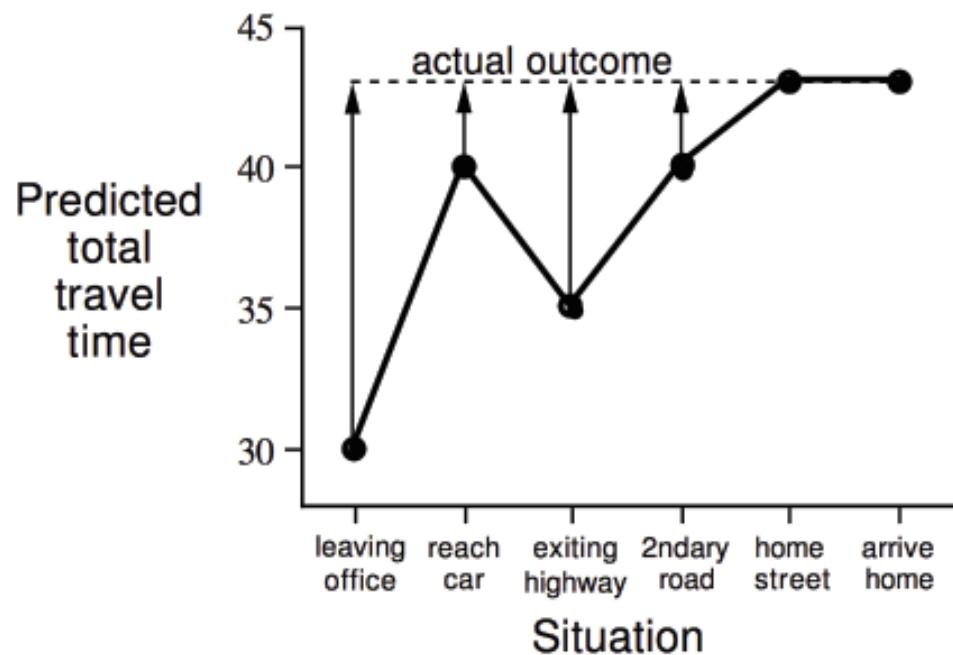
    until  $S$  is terminal

## Example: Driving Home

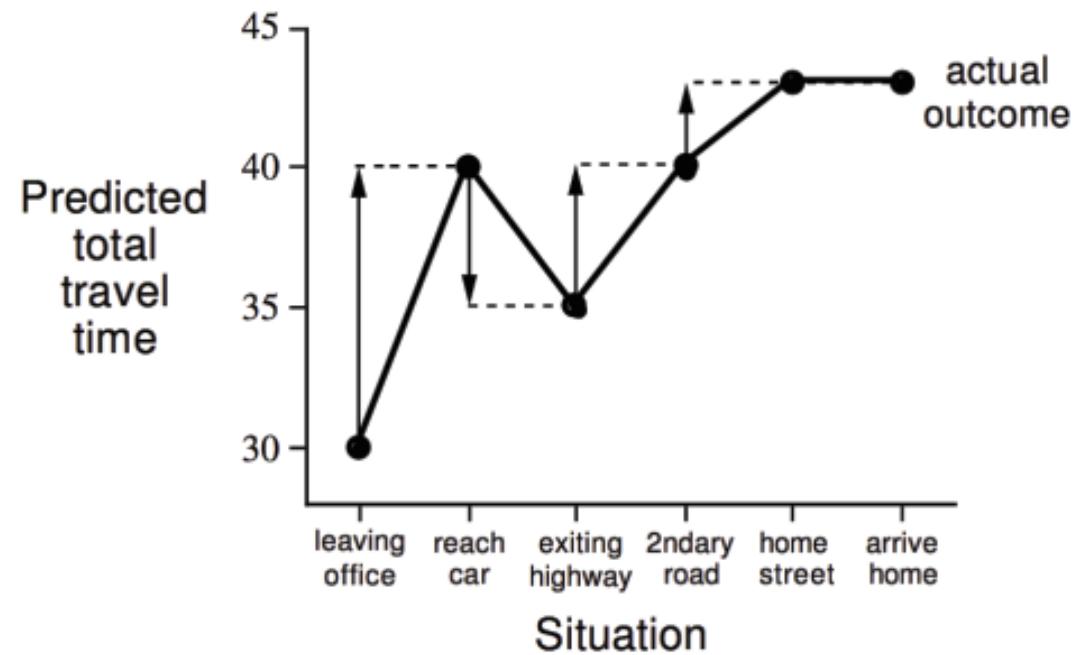
<b>State</b>	<b>Elapsed Time (minutes)</b>	<b>Predicted Time to Go</b>	<b>Predicted Total Time</b>
leaving office	0	30	30
reach car, raining	5	35	40
exit highway	20	15	35
behind truck	30	10	40
home street	40	3	43
arrive home	43	0	43

# Example: Driving Home

Changes recommended by  
Monte Carlo methods ( $\alpha=1$ )



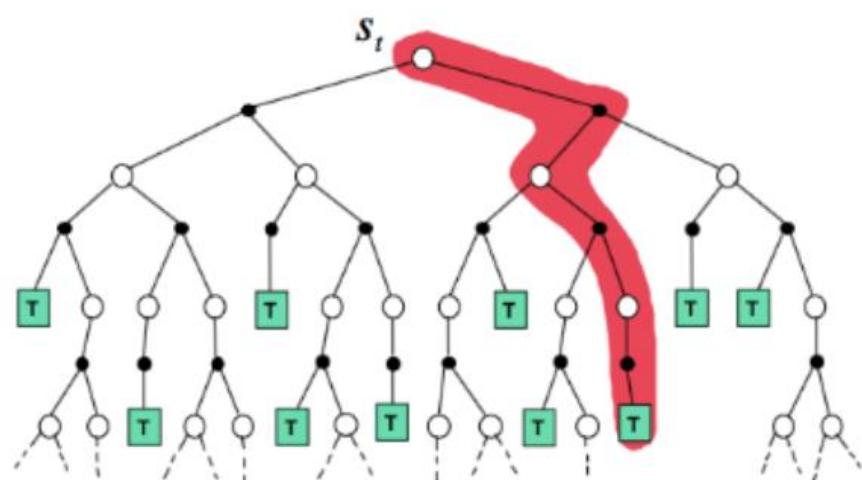
Changes recommended  
by TD methods ( $\alpha=1$ )



# Comparison of state value functions

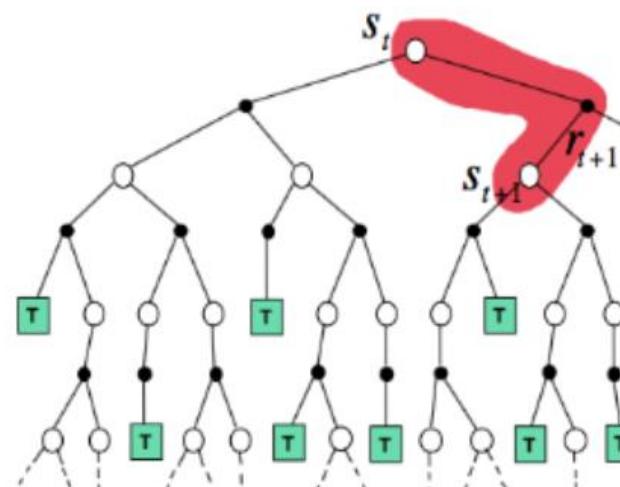
Monte-Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



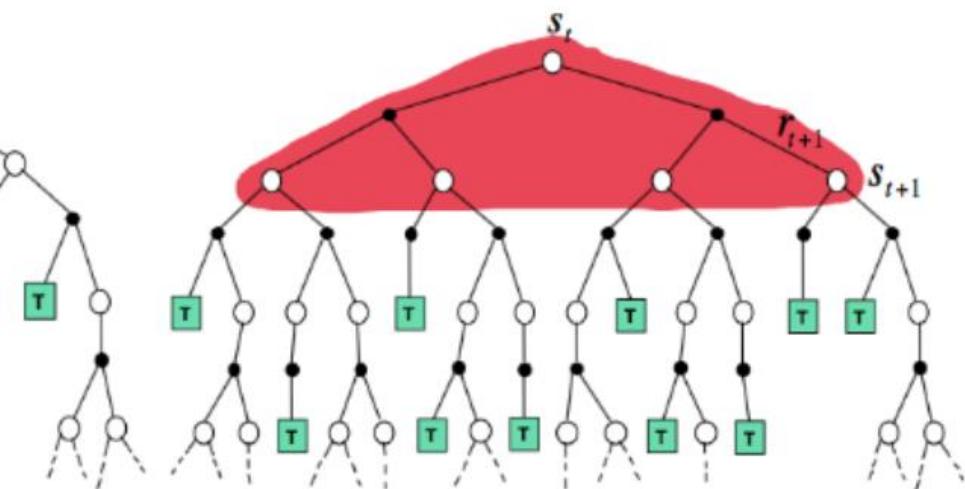
Temporal-Difference

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



Dynamic Programming

$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$



# SARSA: On-Policy TD control

Updating Q-value by following a sequence of  $S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1} \dots, S_T$

Sarsa (on-policy TD control) for estimating  $Q \approx q_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

    Initialize  $S$

    Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  $A_t = \arg \max_{a \in \mathcal{A}} Q(S_t, a)$

    Repeat (for each step of episode):

        Take action  $A$ , observe  $R, S'$

        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$$

$S \leftarrow S'; A \leftarrow A'$ ;

    until  $S$  is terminal

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t)]$$

**Expected SARSA:**  $\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right]$ .

# Q-Learning: Off-policy TD control

Q-learning (off-policy TD control) for estimating  $\pi \approx \pi_*$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

    Initialize  $S$

    Repeat (for each step of episode):

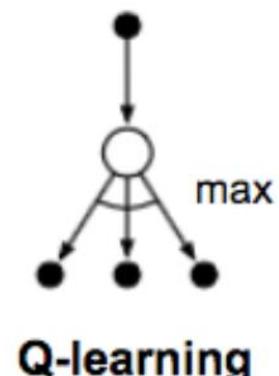
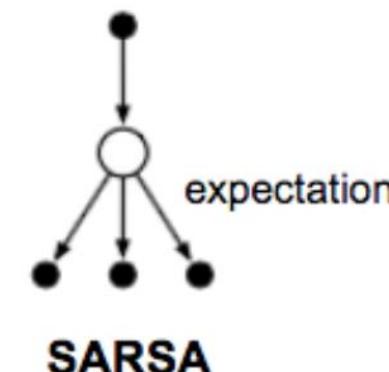
        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)  $A_t = \arg \max_{a \in \mathcal{A}} Q(S_t, a)$

        Take action  $A$ , observe  $R, S'$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

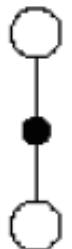
$$S \leftarrow S'$$

    until  $S$  is terminal



# n-Step Algorithm

TD (1-step)



2-step



3-step



$n$ -step



Monte Carlo



# n-Step Return

- Consider the following  $n$ -step returns for  $n = 1, 2, \infty$ :

$$n = 1 \quad (\text{TD}) \quad G_t^{(1)} = R_{t+1} + \gamma V(S_{t+1})$$

$$n = 2 \quad G_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$$

⋮

$$n = \infty \quad (\text{MC}) \quad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

- Define the  $n$ -step return

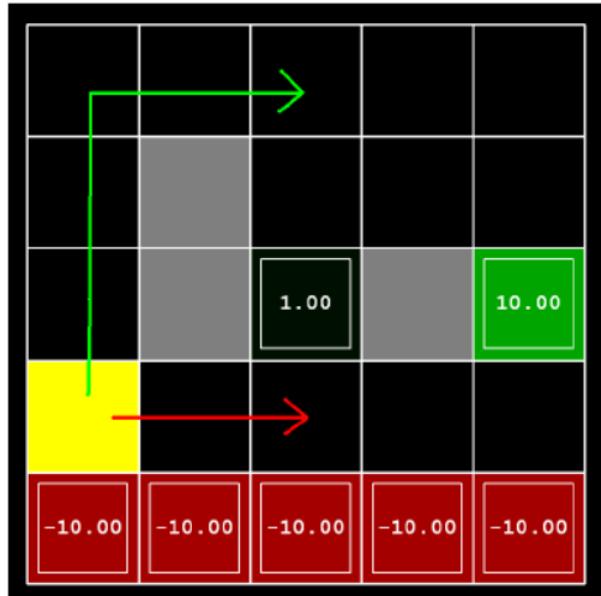
$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

- $n$ -step temporal-difference learning

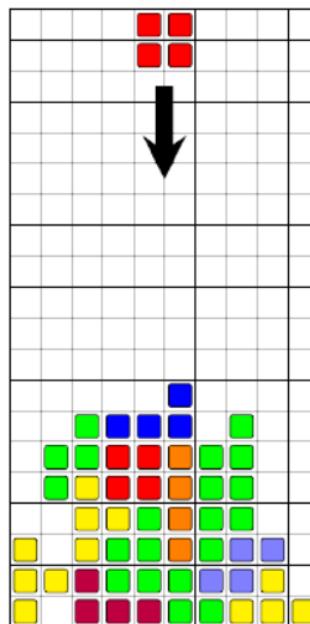
$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t^{(n)} - V(S_t) \right)$$

# Can tabular methods scale?

We need a huge amount of storage for tabular methods. For instance, for discrete environments such as Atari breakout, this could be  $10^{308}$



Gridworld  
 $10^1$



Tetris  
 $10^{60}$



Atari  
 $10^{308}$  (ram)    $10^{16992}$  (pixels)

# References

- David Silver's course slides “Reinforcement Learning”
- Sutton, Richard S., and Andrew G. Barto. "Reinforcement Learning: An introduction (2 edition)" (2018)