

[70240413 Statistical Machine Learning, Spring, 2023]

Reinforcement Learning

Jun Zhu

dcszj@mail.tsinghua.edu.cn

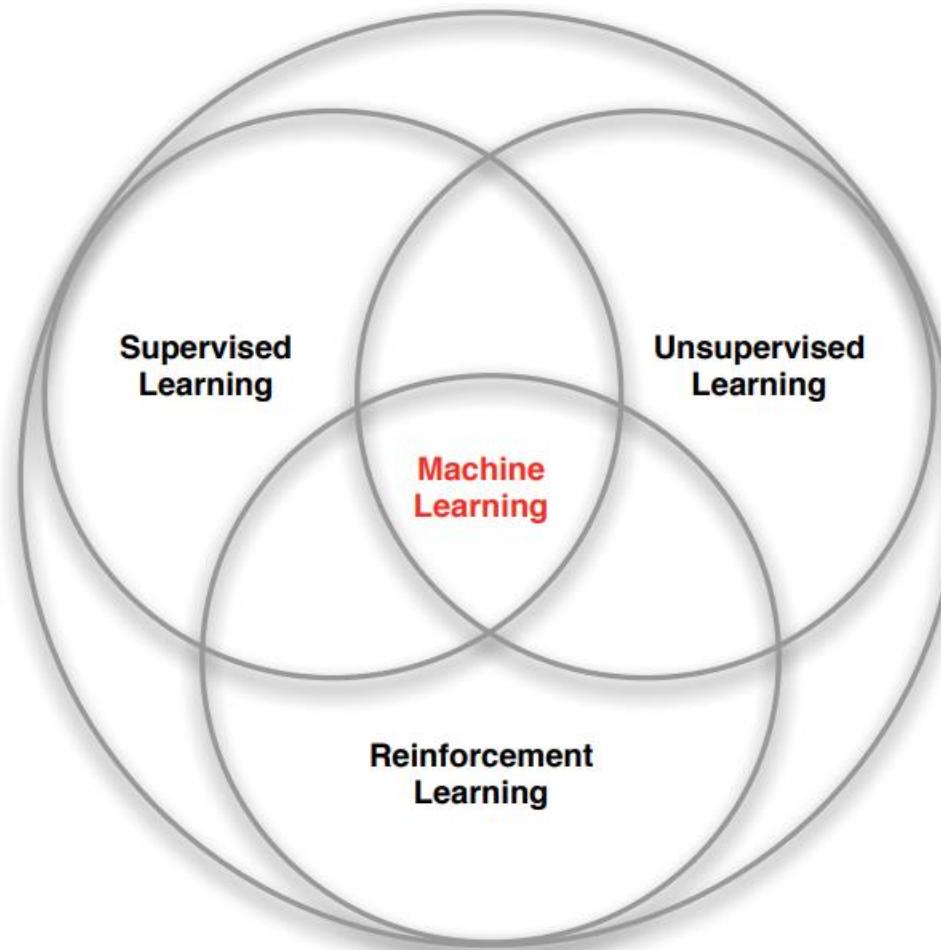
<http://bigml.cs.tsinghua.edu.cn/~jun>

State Key Lab of Intelligent Technology & Systems

Tsinghua University

May 30, 2023

Branches of Machine Learning

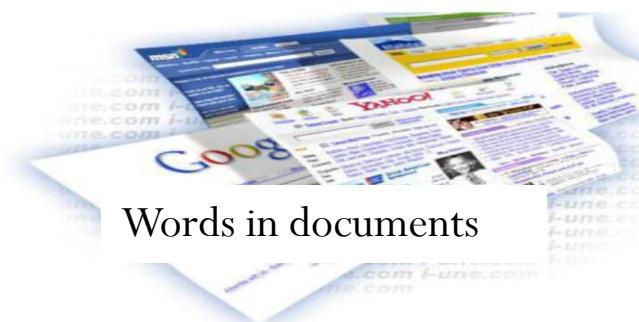


Supervised Learning

- ◆ Task: learn a predictive function

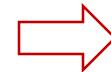
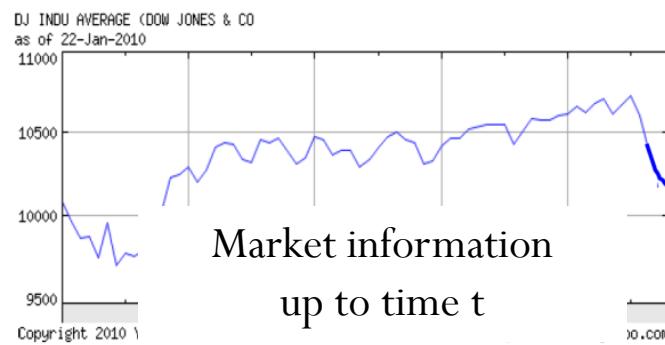
$$h : \mathcal{X} \rightarrow \mathcal{Y}$$

Feature space \mathcal{X}



Label space \mathcal{Y}

“Sports”
“News”
“Politics”
...



Share price
“\$ 20.50”

- ◆ “Experience” or training data:

$$\{\langle x_d, y_d \rangle\}_{d=1}^D, x_d \in \mathcal{X}, y_d \in \mathcal{Y}$$

Unsupervised Learning

- ◆ Task: learn an explanatory function $f(x)$, $x \in \mathcal{X}$
- ◆ Aka “Learning without a teacher”

Feature space \mathcal{X}



Words in documents

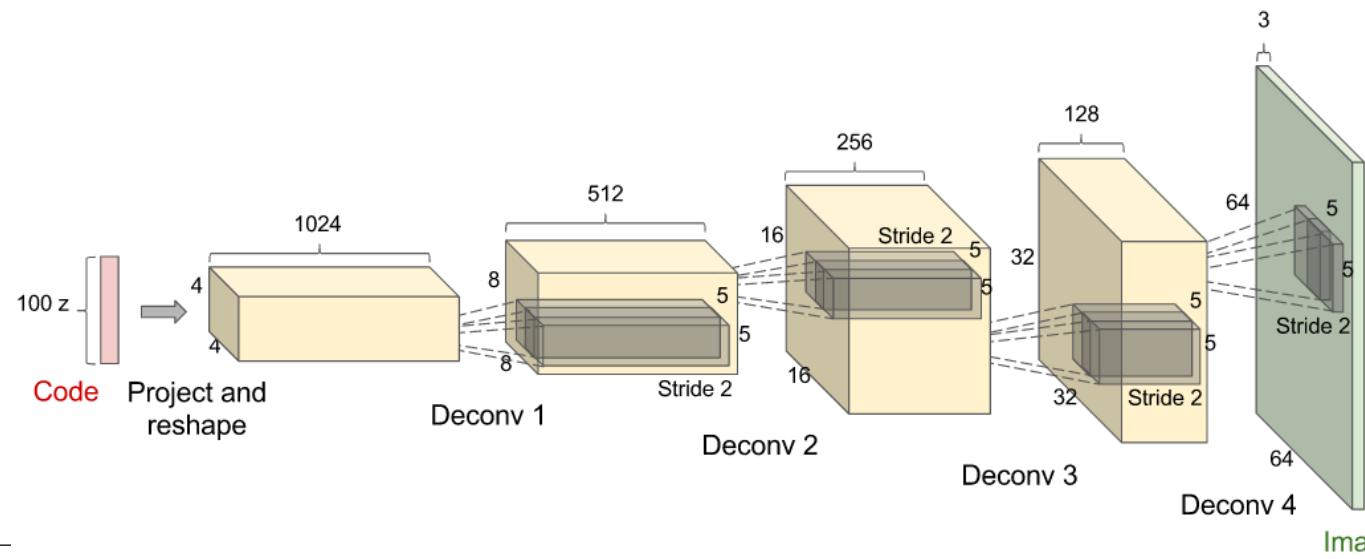
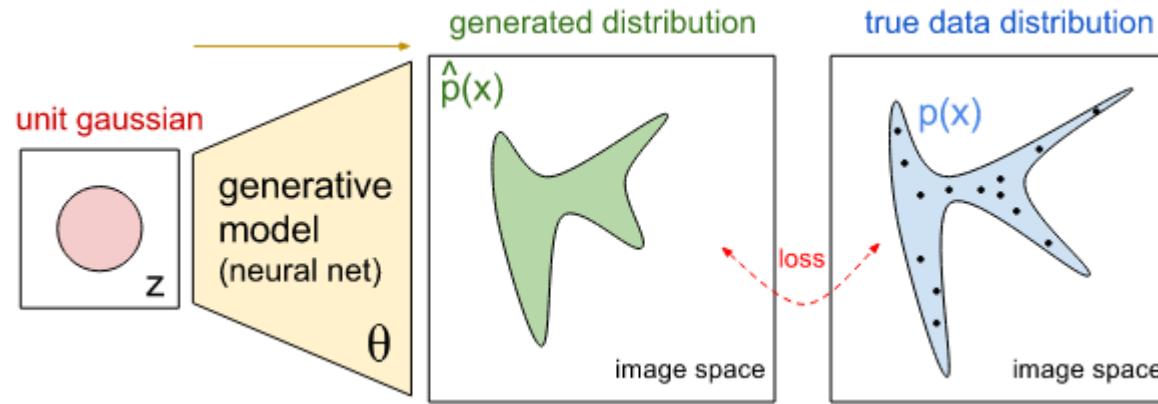


Word distribution
(probability of a word)

- ◆ No training/test split

Deep Generative Models

- ◆ Learn a generative model

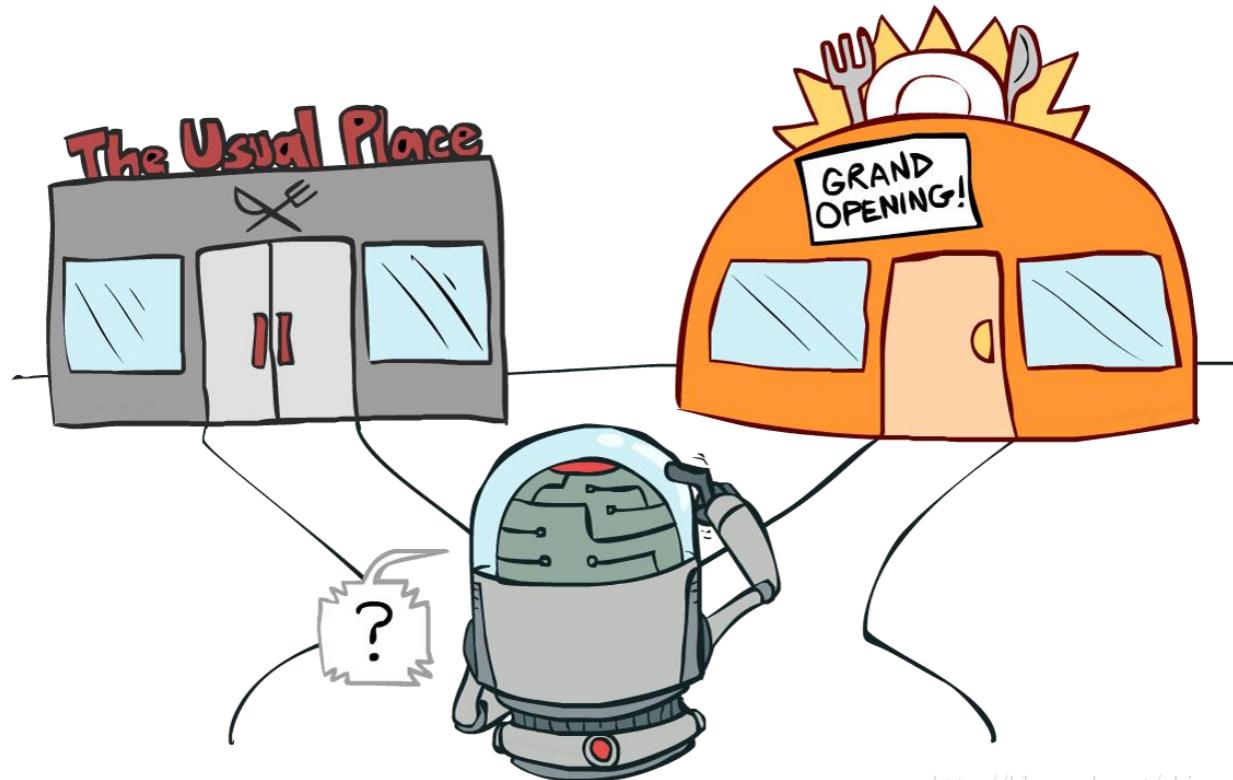


Outline

- ◆ RL preliminaries
- ◆ RL basics
- ◆ RL advanced

Decision Making

- ◆ Make decisions is a key component of intelligence

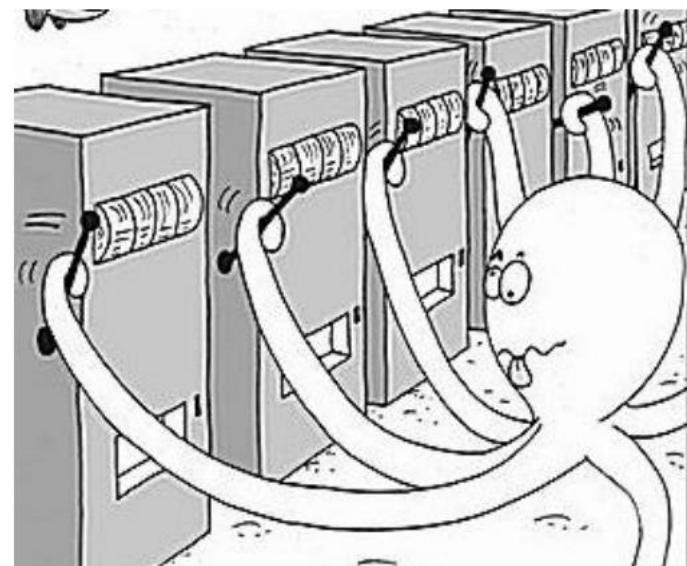


Multi-armed Bandit

- ◆ A gambler walks into a casino
- ◆ A row of slot machines providing random rewards

Objective:

Maximize the sum of rewards (Money)!



Applications

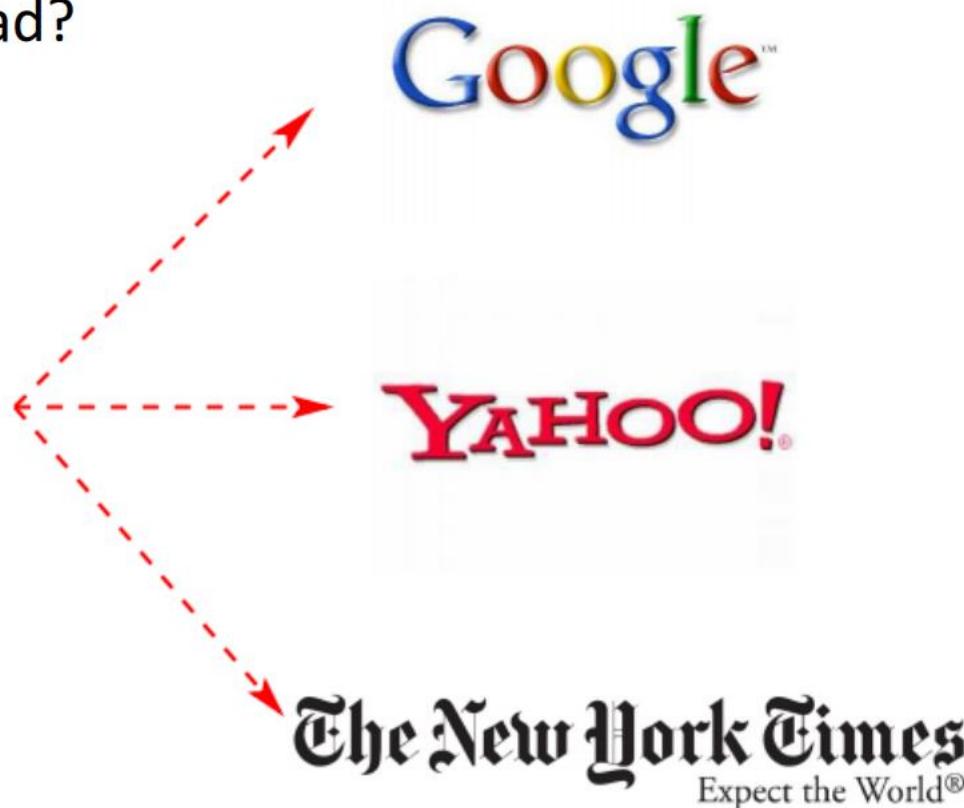
- ◆ Two treatments with unknown effectiveness



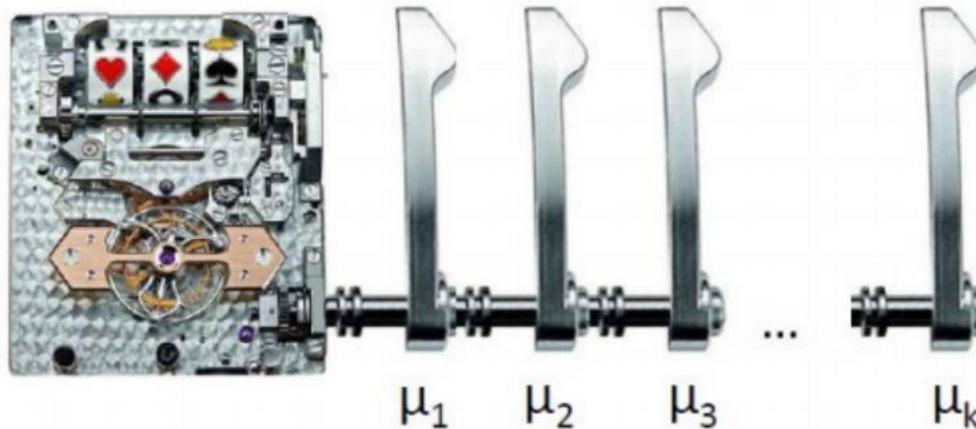
Applications

- ◆ Web advertising

Where to place the ad?



K-armed Bandits



- Each Arm a
 - Wins(reward=1) with fixed(unknown) prob. μ_a
 - Loses(reward=0) with fixed(unknown) prob. $(1 - \mu_a)$
- All draws are independent given $\mu_1 \dots \mu_k$
- How to pull arms to **maximize total reward**? (estimate the arm's prob. of winning μ_a)

Model of K-armed Bandit

- Set of k choices(arms)
- Each choice a is associated with unknown probability distribution P_a in $[0, 1]$
- We play the game for T rounds
- In each round t :
 - We pick some arm j
 - We obtain random sample X_t from P_j (reward is independent of previous draws)
- Goal: maximize $\sum_{t=1}^T X_t$ (without known μ_a)
- However, every time we pull some arm a we get to learn a bit about μ_a .

Performance Metric

- Let be μ_a the mean of P_a
- Payoff/reward **best arm**: $\mu^* = \max\{\mu_a | a = 1, \dots, k\}$
- Let i_1, \dots, i_T be the sequence of arms pulled
- Instantaneous regret at time t: $r_t = \mu^* - \mu_{a_t}$
- Total regret:
 - $R_T = \sum_{t=1}^T r_t$
- Typical goal: arm allocation strategy that guarantees :
 - $\frac{R_T}{T} \rightarrow 0$ as $T \rightarrow \infty$

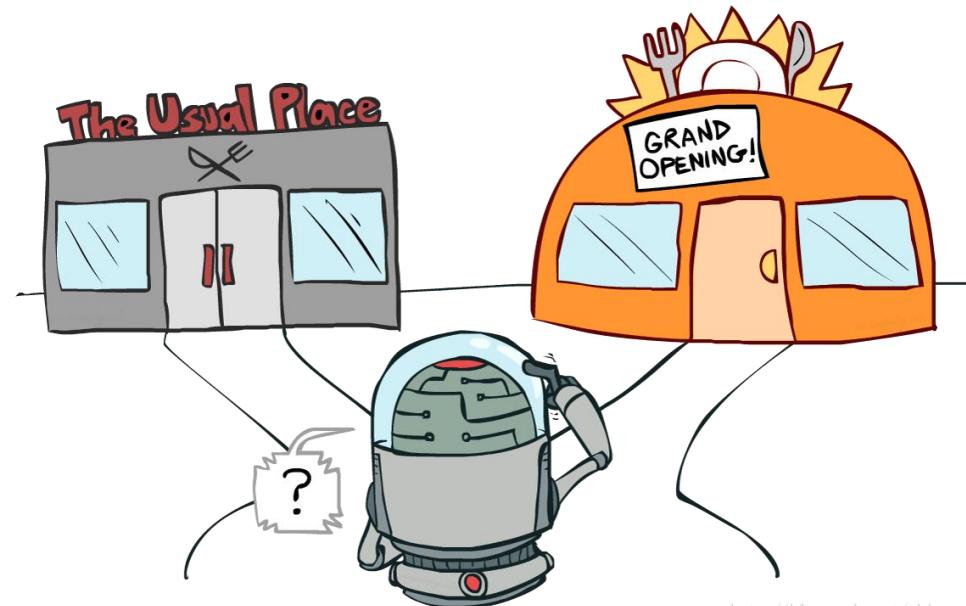
Allocation Strategies

- If we knew the payoffs, which arm should we pull?
 - **best arm:** $\mu^* = \max\{ \mu_a \mid a = 1, \dots, k \}$
- What if we only care about estimating payoff μ_a ?
 - Pick each of k arms equally often : $\frac{T}{k}$
 - **Estimate :** $\widehat{\mu}_a = \sum_{j=1}^k X_{a,j} / (\frac{T}{k}) \rightarrow \frac{k}{T} \sum_{j=1}^{T/k} X_{a,j}$
 - Total regret:
 - $R_T = \frac{T}{k} \sum_{a=1}^k (\mu^* - \mu_a)$

$X_{a,j}$ payoff received when pulling an arm a for j -th time

Need for sufficient exploration

- ◆ Choose the best restaurant
 - 10 restaurants in total
 - You have tried 3 restaurants, and the best one is 清芬园, with score 8
 - What's your choice next time in order to have best dishes?



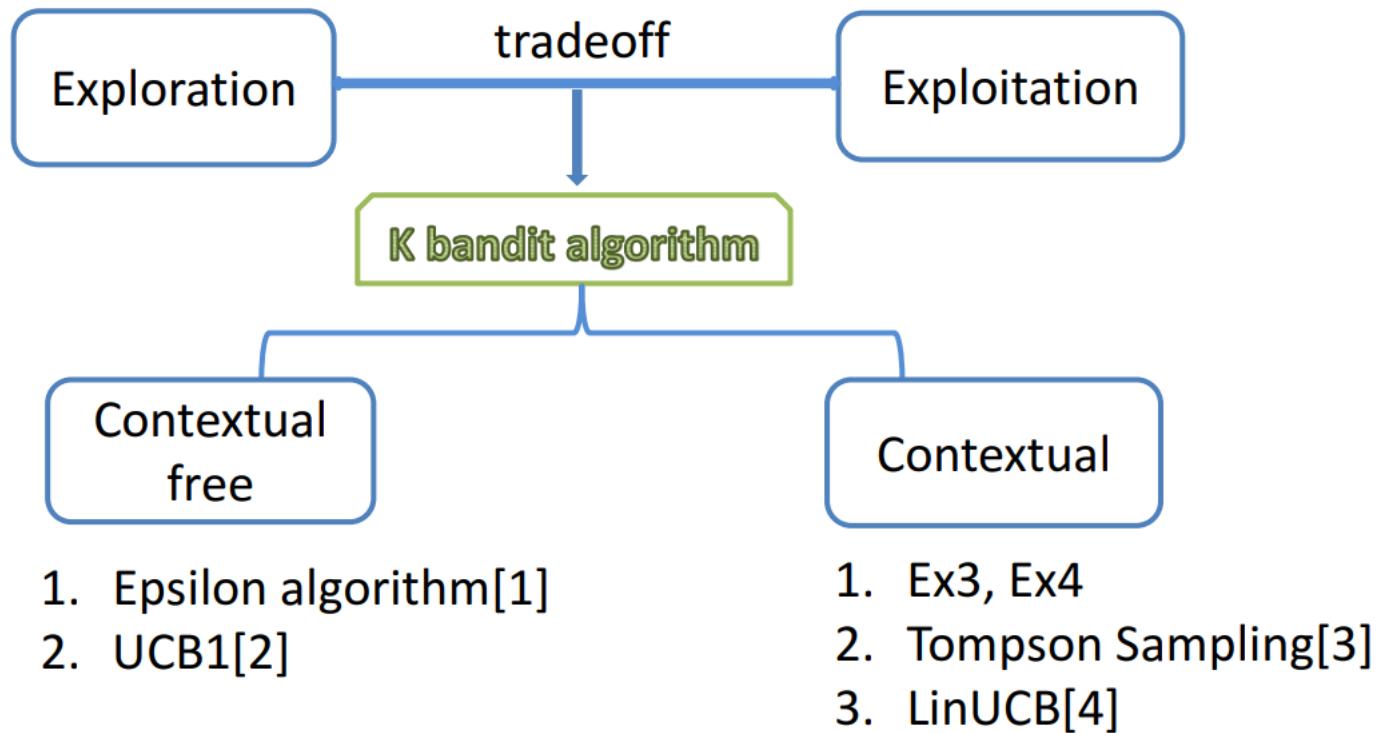
Exploitation vs. Exploration

◆ Tradeoff:

- Only **exploitation**(making decisions based on history data), you will have bad estimation for “best” items.
- Only **exploration**(gathering data about arm payoffs), you will have low user’s engagement.



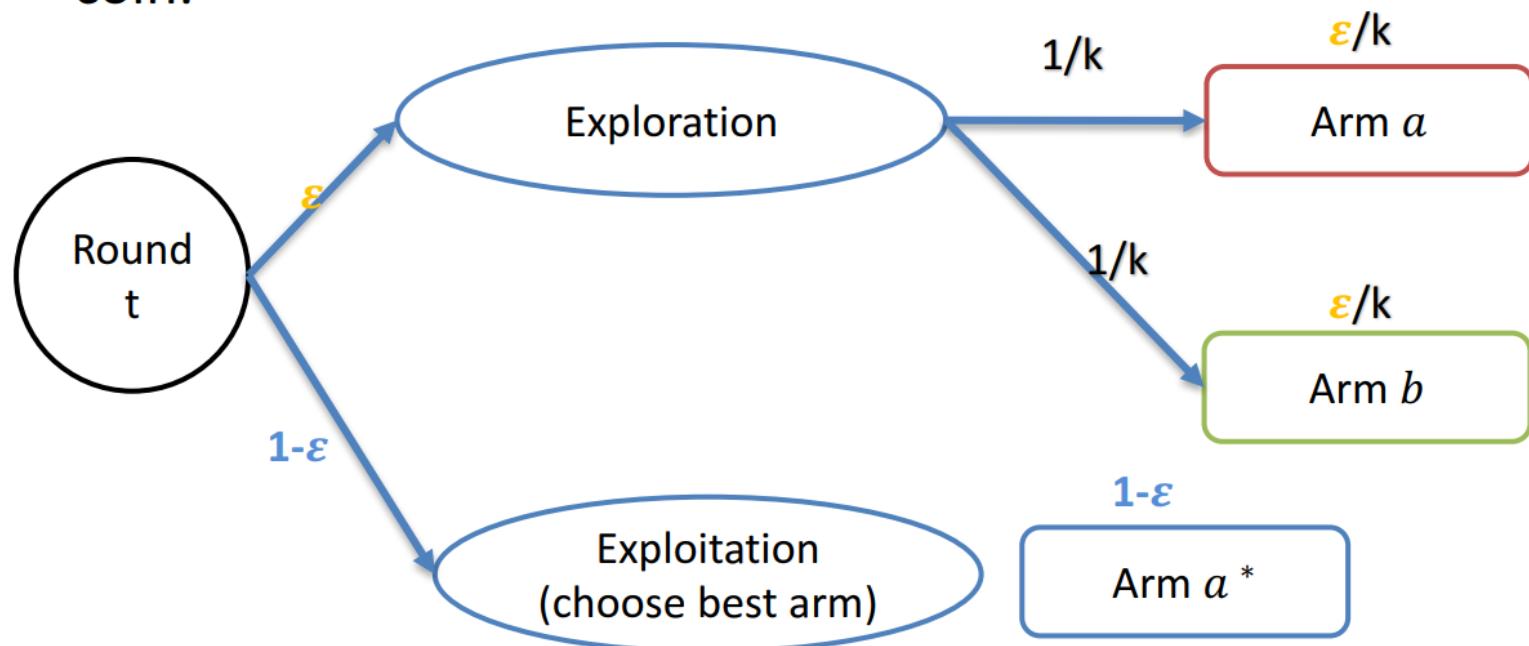
Algorithm to Exploration & Exploitation



- [1] [Wynn P. On the convergence and stability of the epsilon algorithm\[J\]. SIAM Journal on Numerical Analysis, 1966, 3\(1\): 91-122.](#)
- [2] [Auer P, Cesa-Bianchi N, Fischer P. Finite-time analysis of the multi-armed bandit problem\[J\]. Machine learning, 2002, 47\(2-3\): 235-256.](#)
- [3] [Agrawal S, Goyal N. Analysis of Thompson sampling for the multi-armed bandit problem\[J\]. arXiv preprint arXiv:1111.1797, 2011.](#)
- [4] [Li, Lihong, et al. "A contextual-bandit approach to personalized news article recommendation." Proceedings of the 19th international conference on World wide web. ACM, 2010.](#)

Epsilon-Greedy Algorithm

- It tries to be fair to the two opposite goals of exploration (with prob. ε) and exploitation ($1-\varepsilon$) by using a mechanism: flips a coin.



Epsilon-Greedy Algorithm

- For $t=1:T$
 - Set $\varepsilon_t = O\left(\frac{1}{t}\right)$
 - With prob. ε_t : Explore by picking an arm chosen uniformly at random
 - With prob. $1-\varepsilon_t$: Exploit by picking an arm with highest empirical mean payoff
- Theorem [Auer et al. '02]
 - For suitable choice of ε_t it holds that

$$R_T = O(k \log T) \Rightarrow \frac{R_T}{T} = O\left(\frac{k \log T}{T}\right) \rightarrow 0$$

Issues with Epsilon-Greedy Algorithm

- **Not elegant”** : Algorithm explicitly distinguishes between exploration and exploitation
- **More importantly:** Exploration makes **suboptimal choices**(since it picks any arm equally likely)
- Idea: When exploring/exploiting we need to compare arms.

Example: Comparing arms

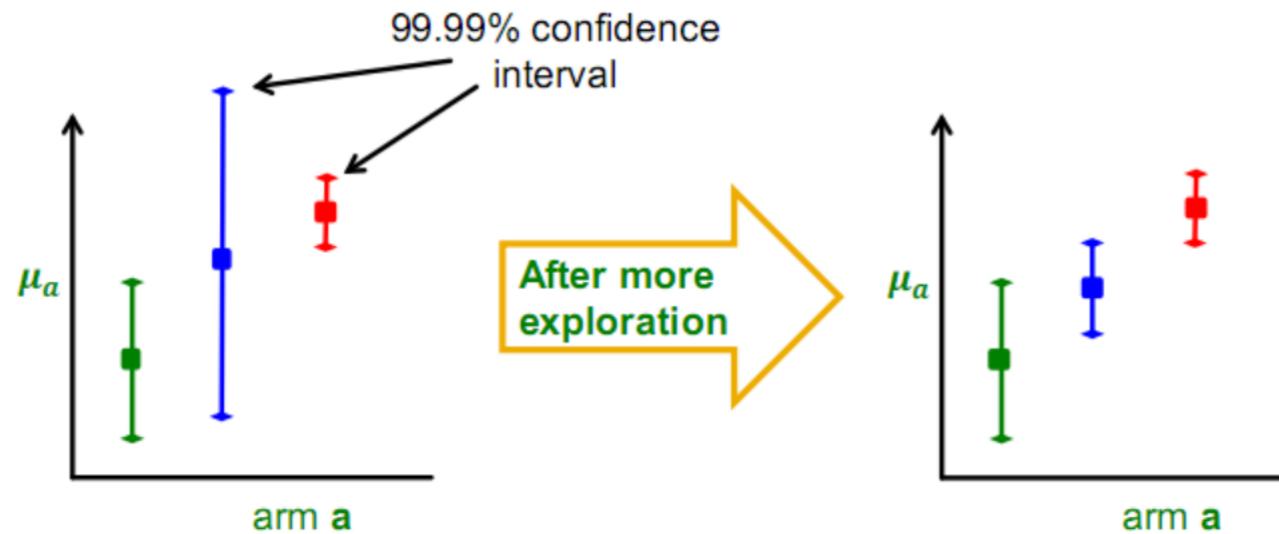
- Suppose we have done experiments :
 - Arm 1: 1 0 0 1 1 1 0 0 0 1
 - Arm 2: 1
 - Arm 3: 1 1 0 1 0 0 1 1 1 1
- Mean arm values:
 - Arm 1: 5/10 Arm 2: 1 Arm 3: 7/10
- Which arm would you choose next?
- Idea: Not only look at the mean but also the confidence!

Confidence Intervals

- A **confidence interval** is a range of values within which we are sure the mean lies with a certain probability
 - We could believe μ_a is within [0.2,0.5] with probability 0.95
 - If we would have tried an action less often, our estimated reward is less accurate so the confidence interval is larger
 - Interval shrinks as we get more information (try the action more often)

Confidence-based Selection

- Assuming we know the confidence intervals
- Then, instead of trying the action with the highest mean we can try the action with the highest upper bound on its confidence interval.



Calculating confidence bounds

- Suppose we fix arm a:
 - Let $r_{a,1} \dots r_{a,m}$ be the payoffs of arm a in the first m trials
 - $r_{a,1} \dots r_{a,m}$ are i.i.d. taking values in [0,1]
 - Our estimate : $\widehat{\mu}_{a,m} = \frac{1}{m} \sum_{j=1}^m r_{a,j}$
 - Want to find b such that with high probability
$$|\mu_a - \widehat{\mu}_{a,m}| \leq b$$
(want b to be as small as possible)
 - Goal : Want to bound $\mathbf{P}(|\mu_a - \widehat{\mu}_{a,m}| \leq b)$

UCB Algorithm

- **UCB1 (Upper confidence sampling) algorithm**
 - Let $\widehat{\mu}_1 = \dots = \widehat{\mu}_k = 0$ and $m_1 = \dots = m_k = 0$
 - $\widehat{\mu}_a$ is our estimate of payoff of arm i
 - m_a is the number of pulls of arm i so far.
 - For $t = 1 : T$
 - For each arm a calculate $UCB(a) = \widehat{\mu}_a + \alpha \sqrt{\frac{2 \ln t}{m_a}}$
 - Pick arm $j = \operatorname{argmax}_a UCB(a)$
 - Pull arm j and observe y_t
 - $m_j = m_j + 1$ and $\widehat{\mu}_j = \frac{1}{m_j} (y_t + (m_j - 1) \widehat{\mu}_j)$

Hoeffding's Inequality

UCB Algorithm

- Theorem [Auer et al. 2002]
 - Suppose optimal mean payoff is
 - And for each arm let $\mu^* = \max_a \mu_a$
 - Then it holds that $\Delta_a = \mu^* - \mu_a$

$$E[R_T] = \left[8 \sum_{a: \mu_a < \mu^*} \frac{\ln T}{\Delta_a} \right] + \left(1 + \frac{\pi^2}{3} \right) \left(\sum_{i=a}^k \Delta_a \right)$$

$O(k \ln T)$ $O(k)$

- So, we get $O\left(\frac{R_T}{T}\right) = k \frac{\ln T}{T}$

Thompson Sampling

- A simple natural Bayesian heuristic
 - Maintain a belief(distribution) for the unknown parameters
 - Each time, pull arm a and observe a reward r
- Initialize priors using belief distribution
 - For $t=1:T$:
 - Sample random variable X from each arm's belief distribution
 - Select the arm with largest X
 - Observe the result of selected arm
 - Update prior belief distribution for selected arm

A Simple Example

- Coin toss: $x \sim \text{Bernoulli}(\theta)$
- Let's assume that
 - $\theta \sim \text{Beta}(\alpha_H, \alpha_T)$
 - $P(\theta) \propto \theta^{\alpha_H-1} (1-\theta)^{\alpha_T-1}$
- $$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{\sum_{\theta} P(X|\theta)}$$

Beta distribution

Prior

The prior is conjugate!

Posterior

Thompson Sampling using Beta belief

- Theorem [Emilie et al. 2012]
 - Initially assumes arm i with prior Beta(1,1) on μ_i
 - $S_i = \#$ “Success”, $F_i = \#$ “Failure”

Algorithm 1: Thompson Sampling for Bernoulli bandits

$S_i = 0, F_i = 0.$

foreach $t = 1, 2, \dots$, **do**

 For each arm $i = 1, \dots, N$, sample $\theta_i(t)$ from the Beta($S_i + 1, F_i + 1$) distribution.

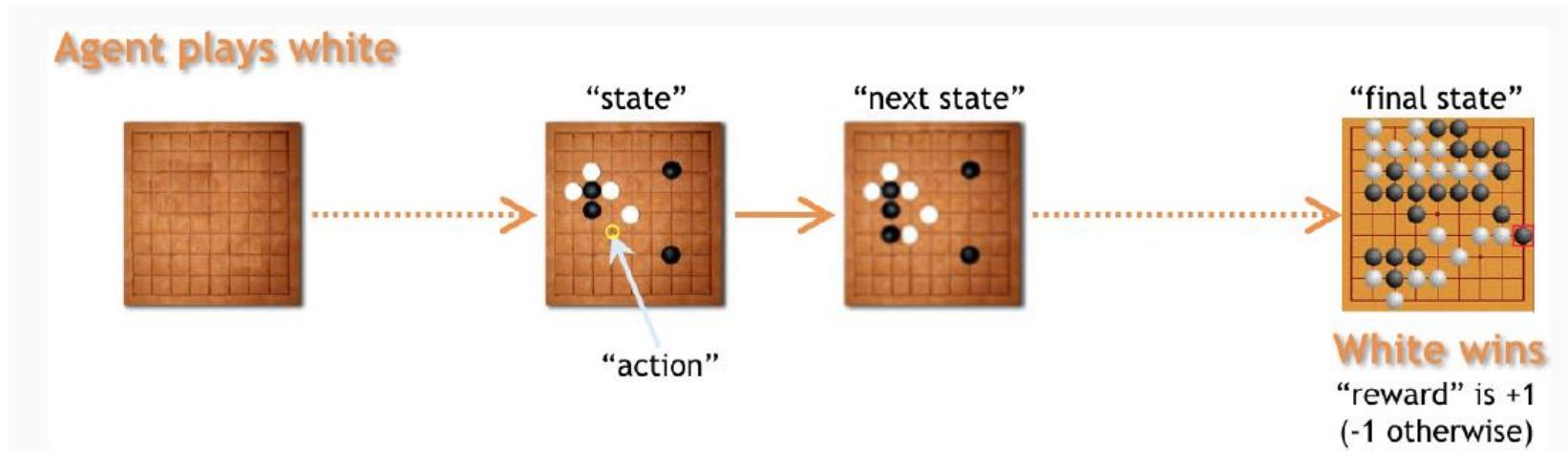
 Play arm $i(t) := \arg \max_i \theta_i(t)$ and observe reward r_t .

 If $r = 1$, then $S_i = S_i + 1$, else $F_i = F_i + 1$.

end

Sequential Decision Making with RL

- ◆ Many tasks are sequentially making decisions



Problem Statement



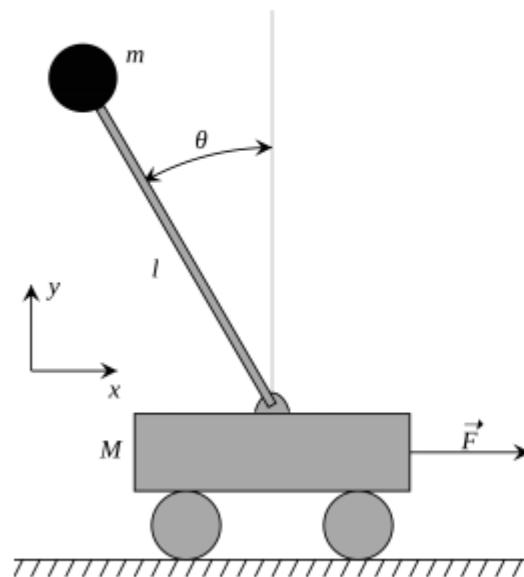
Terminology

- Trajectory/episode: $s_1, a_1, r_1, s_2, a_2, r_2, \dots$
- Policy: $\pi(s_t) \rightarrow a_t$
- Objective: choose a_t to maximize return

$$R_t := r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots$$

Some Examples

Cart-Pole Problem



Objective: Balance a pole on top of a movable cart

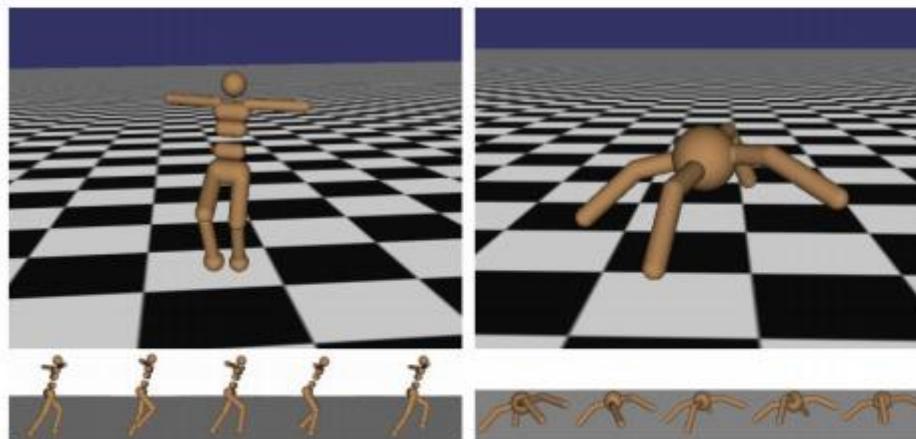
State: angle, angular speed, position, horizontal velocity

Action: horizontal force applied on the cart

Reward: 1 at each time step if the pole is upright

Some Examples

Robot Locomotion



Objective: Make the robot move forward

State: Angle and position of the joints

Action: Torques applied on joints

Reward: 1 at each time step upright + forward movement

Some Examples

Atari Games

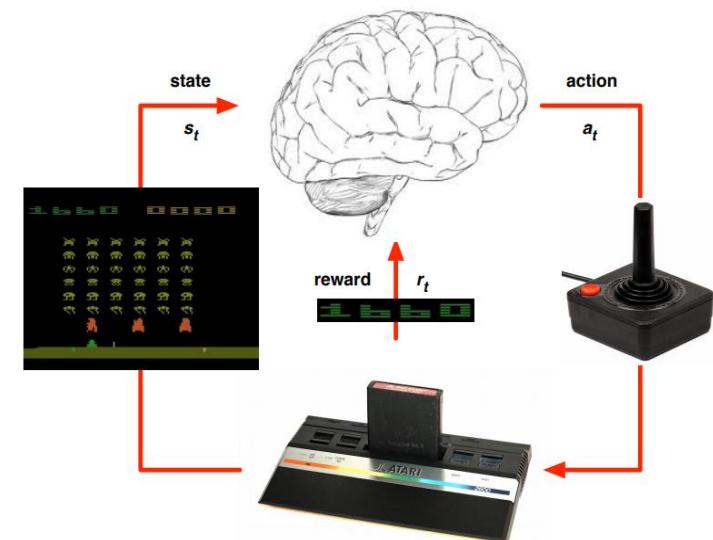


Objective: Complete the game with the highest score

State: Raw pixel inputs of the game state

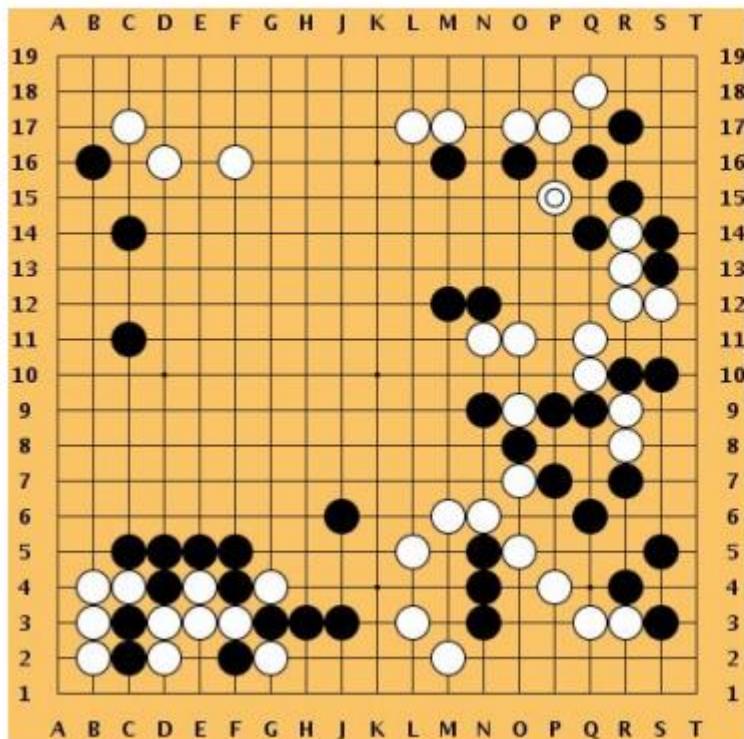
Action: Game controls e.g. Left, Right, Up, Down

Reward: Score increase/decrease at each time step



Some Examples

Go



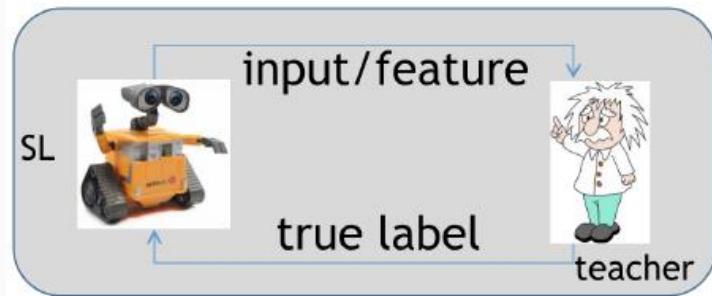
Objective: Win the game!

State: Position of all pieces

Action: Where to put the next piece down

Reward: 1 if win at the end of the game, 0 otherwise

RL vs. SL (supervised learning)



Differences from SL

- Learn by trial-and-error (“experimenting”)
 - Need **exploration/exploitation trade-off**
- Optimize long-term reward
$$(r_1 + \gamma r_2 + \dots)$$
 - Need **temporal credit assignment**

Similarities to SL

- Representation
- Generalization
- Hierarchical problem solving
- ...

RL or SL?



Problem: detect whether an email is spam or not.



Labeled training data

Your favorite
ML algorithm



Supervised learning

RL or SL?



Credit: caradisiac.com

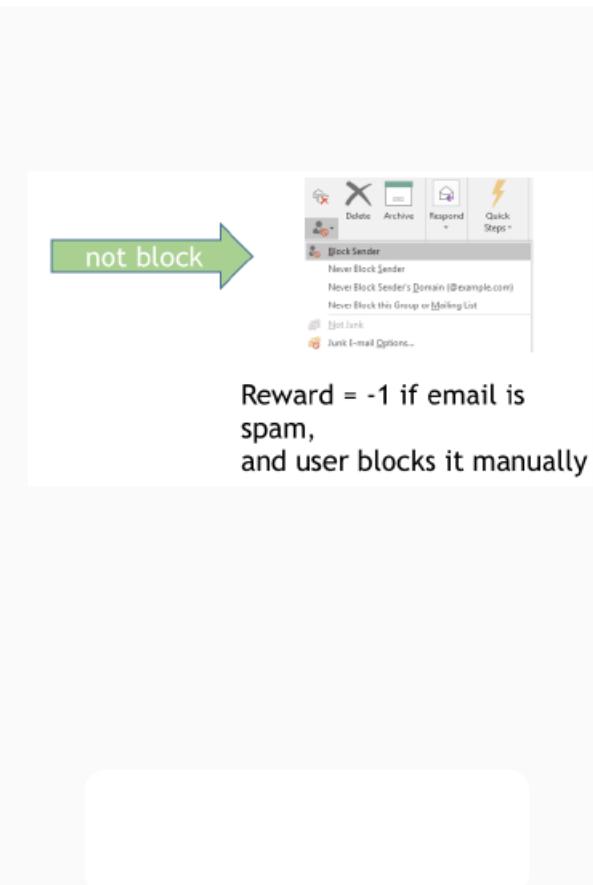
Problem: build a self-driving car

Reinforcement learning

Supervised learning

RL or SL?

- ◆ Email spam (again)



Categories of RL (1)

♦ Model-free vs. Model-based

- A model of the environment allows inferences to be made about how the environment will behave
 - Example: Given a state and an action to be taken while in that state, the model could predict the next state and the next reward
- Models are used for planning, which means deciding on a course of action by considering possible future situations before they are experienced
- Model-based methods use models and planning. Think of this as modelling the dynamics $p(s' | s, a)$
- Model-free methods learn exclusively from trial-and-error (i.e. no modelling of the environment)
- This presentation focuses on model-free methods

Categories of RL (2)

- ◆ On-policy vs. Off-policy
 - An on-policy agent learns only about the policy that it is executing
 - An off-policy agent learns about a policy or policies different from the one that it is executing

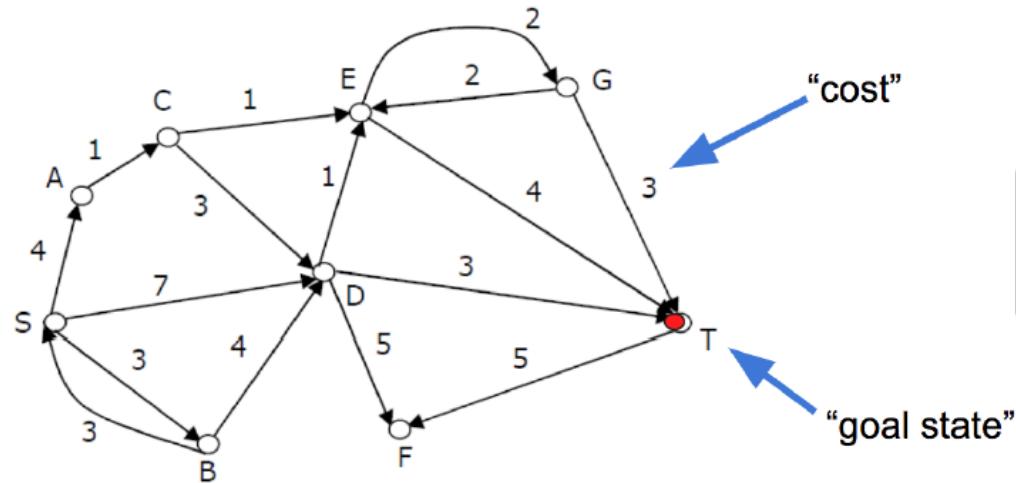
Markov Decision Process (MDP)

- Described by $M = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$:
 - \mathcal{S} : set of states
 - \mathcal{A} : set of actions
 - $P(s'|s, a)$: state-transition probabilities
 - $R(s, a) \in [0, 1]$: average immediate (one-step) reward
 - $\gamma \in [0, 1]$: discount factor
- Policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ (deterministic) or $\mathcal{S} \rightarrow \Delta(\mathcal{A})$ (stochastic)
($\Delta(\mathcal{A})$ is the set of distributions over \mathcal{A})
- Goal of RL: find optimal policy π^* to maximize long-term reward:

$$\pi^* \in \arg \max_{\pi} \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t | \pi \right]$$

- States are **Markovian**
 - States are sufficient statistics of history
 - Future is independent from history, conditioned on current state

Simple MDP: Shortest path problem



node → state
edge → action
cost → (negative) reward

$$\forall i : \text{CostToGo}(i) = \min_{j \in \text{Neighbor}(i)} \{ \text{cost}(i \rightarrow j) + \text{CostToGo}(j) \}$$

Principle of Optimality (Richard Bellman, 1957)

An optimal policy has the property that whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision.



Bellman Equations for MDPs

- ◆ General case:

Deterministic

shortest
path

$$\text{CostToGo}(i) = \min_{j \in \text{Neighbor}(i)} \{\text{cost}(i \rightarrow j) + \text{CostToGo}(j)\}$$



Markov
decision
process

$$V^*(s) = \max_{a \in \mathcal{A}} \{R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[V^*(s')]\}$$

(maximum long-term reward starting from s)

$$Q^*(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[\max_{a' \in \mathcal{A}} Q^*(s', a') \right]$$

(maximum long-term reward after choosing a from s)

V^* and Q^* are called **optimal value functions**

Bellman Operator

- Bellman operator \mathcal{B} transforms Q to another function $\mathcal{B}Q$ on $\mathcal{S} \times \mathcal{A}$:

$$\mathcal{B}Q(s, a) := R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [\max_{a'} Q(s', a')]$$

Special case with a fixed policy π :

$$\mathcal{B}^\pi Q(s, a) := R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [Q(s', \pi(s'))]$$

Similarly for $V(s)$

- Bellman equations re-expressed as: $Q^* = \mathcal{B}Q^*$ (or $Q^\pi = \mathcal{B}^\pi Q^\pi$)
(Q^* and Q^π are called “fixed points” of \mathcal{B} and \mathcal{B}^π)
- Some key properties of \mathcal{B} :
 - Q^* leads to an optimal policy

$$\pi^*(s) := \max_a Q^*(s, a)$$

- More generally, near-optimal Q leads to a near-optimal policy [4]
- Monotonicity: $Q_1 \geq Q_2$ implies $\mathcal{B}Q_1 \geq \mathcal{B}Q_2$
- It is a contraction
- Its fixed point exists and is unique

Value Iteration (VI)

Algorithm

- Initialize Q_0 arbitrarily (e.g., $Q_0 \equiv 0$)
- For $k = 1, 2, \dots$ (until **termination condition**)
 - $Q_k \leftarrow \mathcal{B}Q_{k-1}$

◆ Convergence:

- So V.I. converges to Q^* **exponentially** fast
- Can be used to decide **termination condition** for VI

Policy Iteration (PI)

Algorithm

- Start with arbitrary policy $\pi_0 : \mathcal{S} \rightarrow \mathcal{A}$
- For $k = 0, 1, 2, \dots$
 - Policy **evaluation**: solve for Q_k that satisfies

$$\forall (s, a) : Q_k(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) Q_k(s', \pi_k(s'))$$

(Just solving Bellman equation $\mathcal{B}^{\pi_k} Q = Q$)

(Just a system of linear equations)

- Policy **improvement**:

$$\pi_{k+1}(s) \leftarrow \arg \max_a Q_k(s, a)$$

(π_{k+1} is called a **greedy policy** w.r.t. Q_k)

Policy Improvement Theorem

Theorem

In PI, either π_{k+1} is strictly better than π_k , or π_k is optimal.

VI vs. PI

- VI works in value function space (asymptotic convergence)

$$Q_0 \rightarrow Q_1 \rightarrow Q_2 \rightarrow \dots \rightarrow Q^*$$

- PI does hill-climbing in policy space (finite-time convergence)

$$\pi_0 \xrightarrow{Q_0} \pi_1 \xrightarrow{Q_1} \pi_2 \dots \xrightarrow{Q_{K-1}} \pi_K = \pi^*$$

RL Basics

RL in MDP

Typical life of an RL agent

- Observe initial state s_1
- For $t = 1, 2, 3, \dots$
 - Choose action a_t based on s_t and current policy
 - Observe reward r_t and next state s_{t+1}
 - Update policy using new information (s_t, a_t, r_t, s_{t+1})
- Episode length may be finite or infinite
- Agent can have multiple episodes starting from new initial states

What we have learned so far

- Given MDP model ($R(s, a)$ and $P(s'|s, a)$), we can compute an optimal policy, using value iteration, policy iteration, etc.
- What if R and P are **unknown**?
 - This is what reinforcement **learning** is about!

Q-Learning

- Value iteration assumes knowledge of R and P :

$$Q(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q(s', a')$$

- We may use data to approximate **unknown** quantities
- Given $\mathcal{D} = (s_1, a_1, r_1, s_2, a_2, r_2, \dots)$, we want Q so that

$$\forall t : Q(s_t, a_t) \approx \underbrace{r_t + \gamma \max_{a'} Q(s_{t+1}, a')}_{\text{"bootstrapping": use own estimate to create learning target}}$$

- This inspires Q-learning update

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t \underbrace{\left(r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right)}_{\text{"temporal difference" or "TD error"}}$$

where $\alpha_t \in (0, 1)$ is a stepsize.

- Bootstrapping: enables learning **without a teacher** (unlike SL)

Convergence of Q-Learning

Observations

- Q-learning updates are **local**: only $Q(s_t, a_t)$ changes at step t
- Bootstrapped value is a locally **unbiased** estimate of \mathcal{B}

$$\mathbb{E}_{r_t, s_{t+1} | s_t, a_t} \left[r_t + \gamma \max_{a'} Q(s_{t+1}, a') \right] = \mathcal{B}Q(s_t, a_t)$$

Theorem (from Stochastic Approximation theory [5])

Q-learning running on trajectory $\mathcal{D} = (s_1, a_1, r_1, s_2, a_2, r_2, s_3, \dots)$ converges to Q^ almost surely, if the following holds for all (s, a) :*

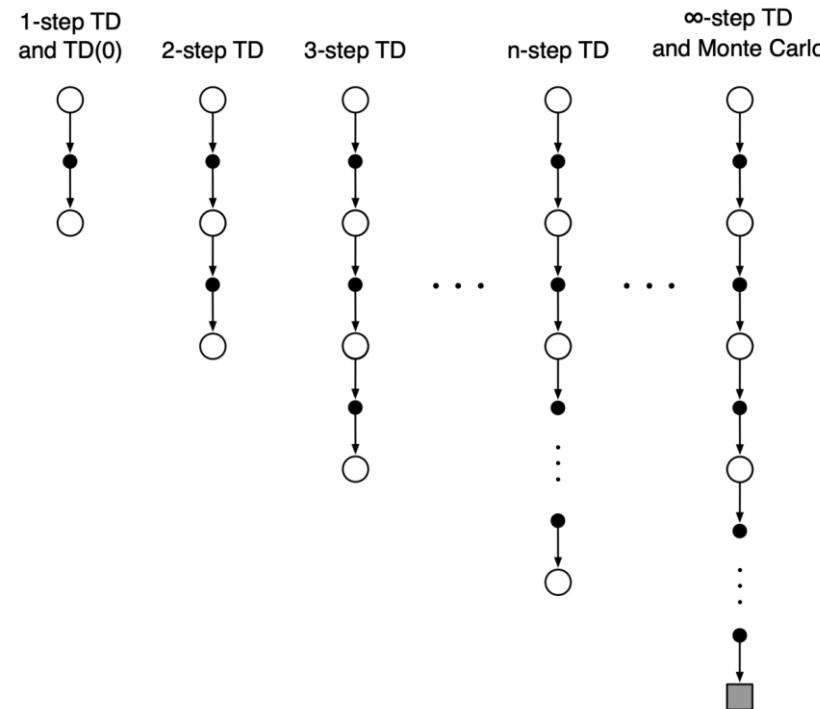
$$\sum_{t=1}^{\infty} \alpha_t \cdot \mathbb{I}\{s_t = s, a_t = a\} = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 \cdot \mathbb{I}\{s_t = s, a_t = a\} < \infty.$$

*The condition implies that every (s, a) must occur **infinitely often** in \mathcal{D} .*

n-step TD

- ◆ n-step TD methods that generalize one-step TD
- ◆ We can shift from one to the other smoothly as needed to meet the demands of a particular task.



n-step TD prediction

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

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

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

⋮

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

- ◆ Thus the n-step return is defined as

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

- ◆ n -step TD:

$$v(S_t) \leftarrow v(S_t) + \alpha(G_t^n - v(S_t))$$

Q-learning with ϵ -greedy exploration

- 1: **Input:** $\epsilon \in (0, 1)$, $\{\alpha_t\}$
- 2: Initialize $Q(s, a) \leftarrow 0$ for all (s, a)
- 3: Observe initial state s_1
- 4: **for** $t = 1, 2, 3, \dots$ **do**
- 5: Choose action with ϵ -greedy exploration:

$$a_t \leftarrow \begin{cases} \arg \max_a Q(s_t, a), & \text{with prob. } 1 - \epsilon \\ \text{random action,} & \text{with prob. } \epsilon \end{cases}$$

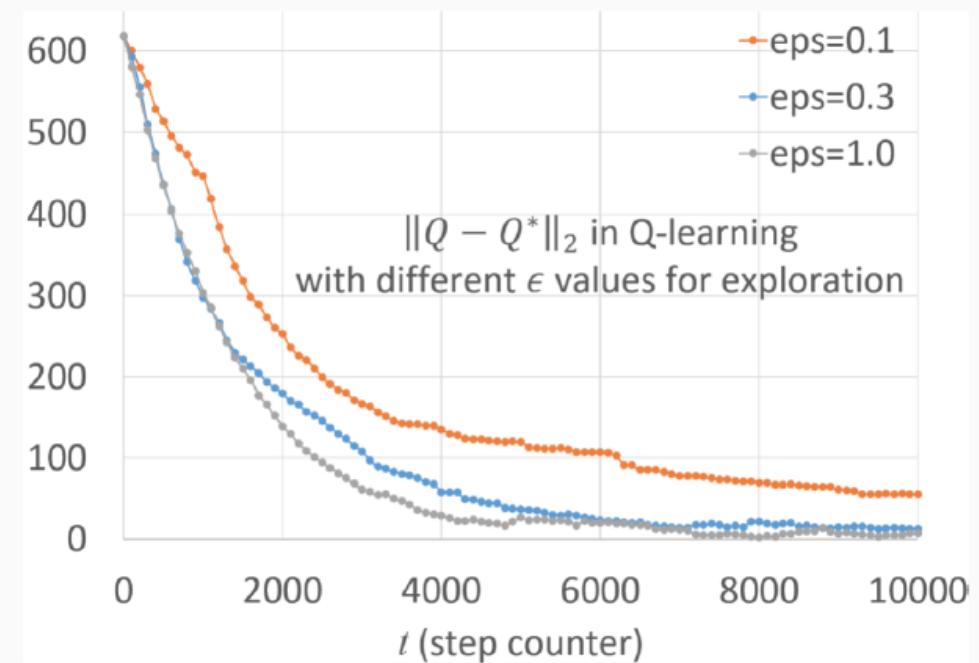
- 6: Observe reward r_t and next-state s_{t+1}
- 7: Update Q-function

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t \left(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right)$$

- 8: **end for**

Some Examples

- Slowly decaying learning rate from 0.05 to 0.0001
- Different ϵ values
- Convergence is always guaranteed, but at different speed



Main lesson

Need efficient exploration/exploitation trade-off!

Exploration: ϵ -greedy

Selection rule

At every step, select a random action with small probability, and the current best-guess otherwise:

$$a_t = \begin{cases} \arg \max_a Q(s_t, a), & \text{with prob. } 1 - \epsilon \\ \text{random action,} & \text{with prob. } \epsilon \end{cases}$$

- Larger ϵ , more exploration
 - Uniform random exploration when $\epsilon = 1$
- Very simple to implement; can be quite effective in practice
 - Often use a sequence of decaying ϵ values, reducing amount of exploration over time.
- Related strategies: ϵ -first, softmax/Boltzmann exploration
- Theoretically flawed; can get stuck in certain cases

Exploration Bonus

Selection rule

$$a_t = \arg \max_a \{ Q(s_t, a) + B(s_t, a) \}$$

for some pre-defined exploration bonus function $B(s_t, a)$.

- $B(s, a)$ measures the need for exploration, e.g., some *decreasing* function of $C(s, a)$ — number of times a has been chosen in s
- Popular choice: $B(s, a) = \alpha / \sqrt{C(s, a)}$ with parameter $\alpha > 0$.
 - Also known as Upper Confidence Bound (UCB)
- An example of guided exploration (as opposed to ϵ -greedy)

Other Exploration Strategies

- Unguided exploration: ϵ -greedy, ϵ -first, forced exploration, ...
- Boltzmann/softmax exploration:
- Exponentiated exploration rules [1]
- Thompson sampling [6, 26]
- Exploration bonus: UCB, pseudocounts, curiosity driven, ...
- ...

RL Advanced

Approximation and Generalization

So far we have focused on **finite** state/action MDPs where all calculations (such as Bellman operator) can be done **exactly**.

- **Computationally**, how to handle large MDPs?
 - Number of legal states in Go is 2×10^{170}
 - In some problems (e.g., autonomous driving), state/actions are even continuous

No relation is assumed between two distinct states $s', s'' \in \mathcal{S}$.

- **Statistically**, how to generalize information to unseen states to reduce sample complexity?
 - No one is able to exhaust all states before mastering Go

That is why we need good **representation** in RL

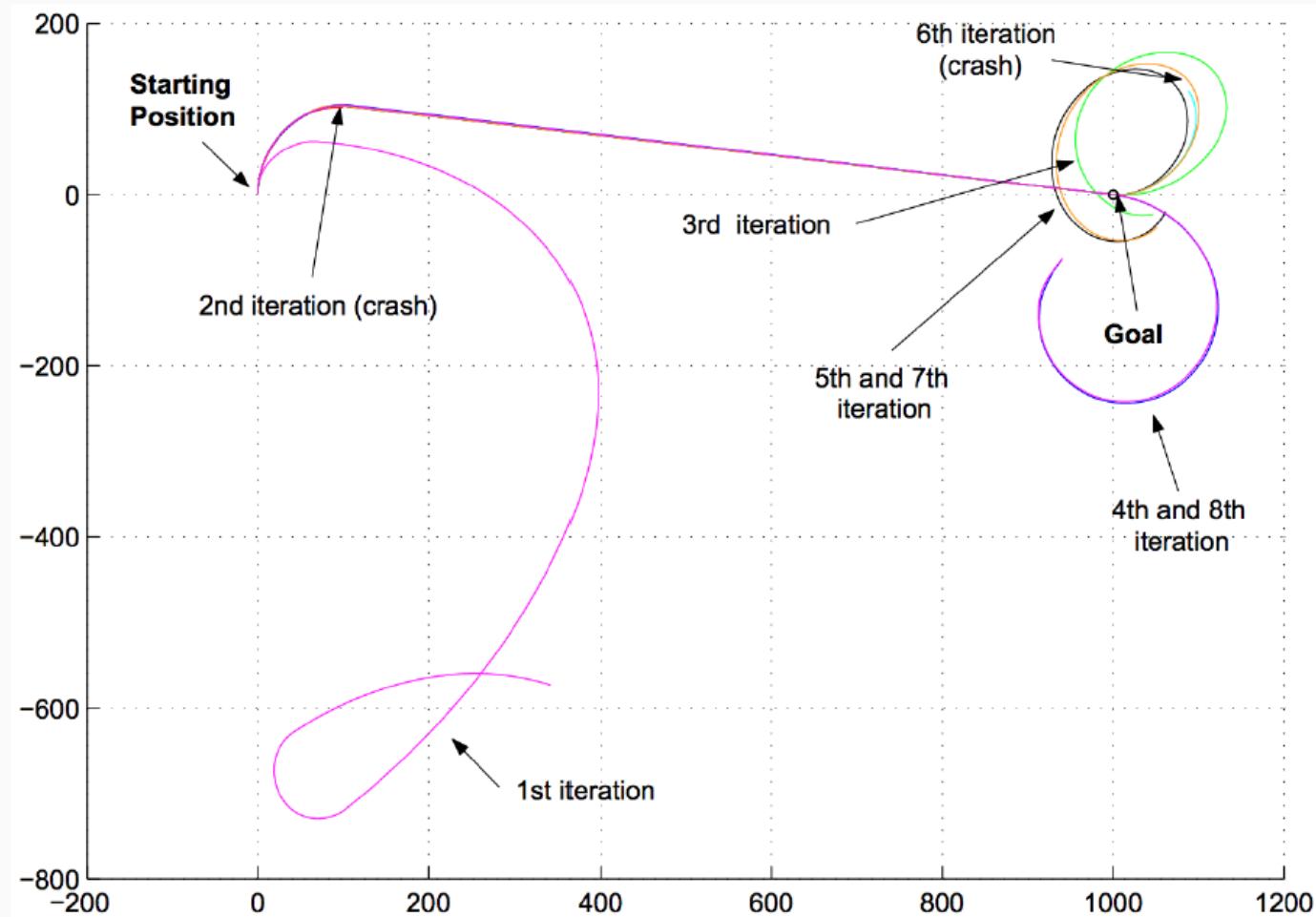
Function Approximation (FA) in RL

- Tabular/exact representations (covered so far); do not scale well
 - Q , V and π represented as $|\mathcal{S}|$ - or $|\mathcal{S} \times \mathcal{A}|$ -dimensional vectors
- Aggregation/abstraction/discretization: treating multiple states/actions as the same [20]
- Linearly parametric approximation with given features (a.k.a. basis functions):

$$Q(s, a; \theta) = \theta^\top \phi(s, a), \quad \text{where } \theta, \phi(s, a) \in \mathbb{R}^d$$

- Nonlinear parametric approximation, e.g., neural nets
- Non-parametric representations: kernel, decision trees, ...

Learning to Ride Bicycle



Results of Least-Squares Policy Iteration (linear FA) [16]

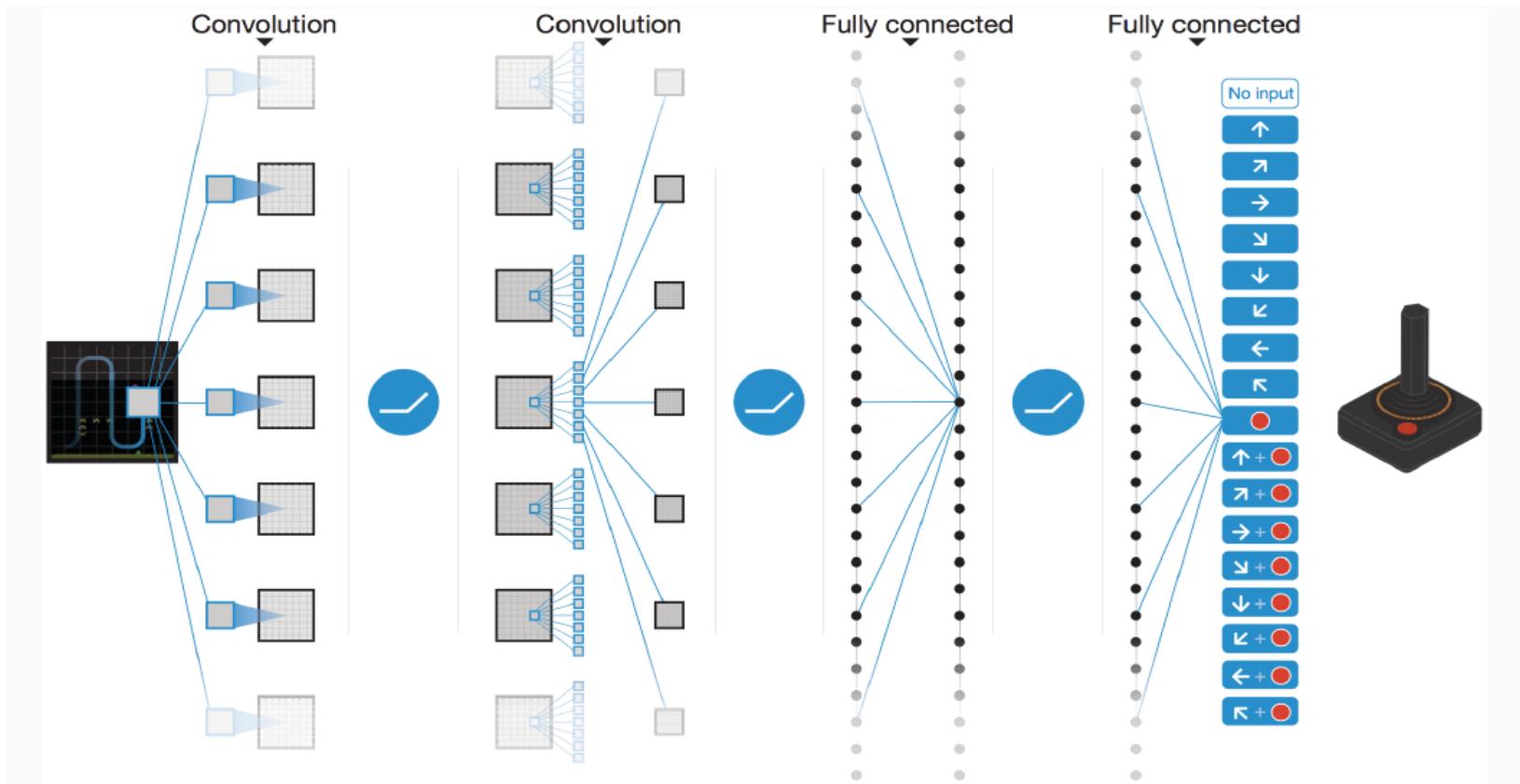
Learning to Ride Bicycle: States & Features

- State $s = (\theta, \dot{\theta}, \omega, \dot{\omega}, \ddot{\omega}, \psi) \in \mathbb{R}^6$
 - θ : angle of handlebar
 - ω : vertical angle of the bicycle
 - ψ : angle of the bicycle to the goal
- Action $a = (\tau, \nu) \in \mathbb{R}^2$
 - $\tau \in \{0, \pm 2\}$: torque applied to the handlebar
 - $\nu \in \{0, \pm 0.02\}$: displacement of the rider
- Popular featurization techniques: polynomial basis, radial basis, Fourier basis, coarse coding, ...
- Previous slide's results are based on degree-2 polynomial features:

$$Q(s, a; \theta) = \theta_a^\top \phi(s), \quad \text{where } \theta = (\theta_1, \dots, \theta_{|\mathcal{A}|}) \in \mathbb{R}^{|\phi| \times |\mathcal{A}|}$$

- Challenge: how to construct features automatically [24]
- Use neural networks to learn representations (“deep RL”)

Nonlinear FA with Neural Nets



The [DQN](#) architecture that enables an RL agent to excel at a wide range of video games [23]. The use of neural networks in RL dates back to the 80s, with a great success of TD-Gammon that beat human champion in the game of back-gammon [36].

Q-Learning with Function Approximation

- Recall tabular Q-learning update on transition (s_t, a_t, r_t, s_{t+1})

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha_t \underbrace{\left(r_t + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right)}_{\text{TD error, denoted by } \delta_t}$$

regression target

- Q-learning with differentiable parametric representation $Q(s, a; \theta)$

$$\theta \leftarrow \theta + \alpha_t \cdot \delta_t \cdot \nabla_\theta Q(s_t, a_t; \theta)$$

- Linear FA special case with $Q(s, a; \theta) = \theta^\top \phi(s, a)$

$$\theta \leftarrow \theta + \alpha_t \cdot \delta_t \cdot \phi(s_t, a_t)$$

Q-learning with FA Pseudo-code

- 1: **Input:** $\theta_0, \epsilon \in (0, 1), \{\alpha_t\}$
- 2: Initialize $\theta \leftarrow \theta_0 \in \Re^d$
- 3: Observe initial state s_1
- 4: **for** $t = 1, 2, 3, \dots$ **do**
- 5: Choose action with ϵ -greedy exploration:

$$a_t \leftarrow \begin{cases} \arg \max_a Q(s_t, a; \theta), & \text{with prob. } 1 - \epsilon \\ \text{random action,} & \text{with prob. } \epsilon \end{cases}$$

- 6: Observe reward r_t and next-state s_{t+1}
- 7: Update Q-function

$$\theta \leftarrow \theta + \alpha_t \left(r_t + \gamma \max_a Q(s_{t+1}, a; \theta) - Q(s_t, a_t; \theta) \right) \nabla_\theta Q(s_t, a_t; \theta)$$

- 8: **end for**

Heuristics that Sometimes Work

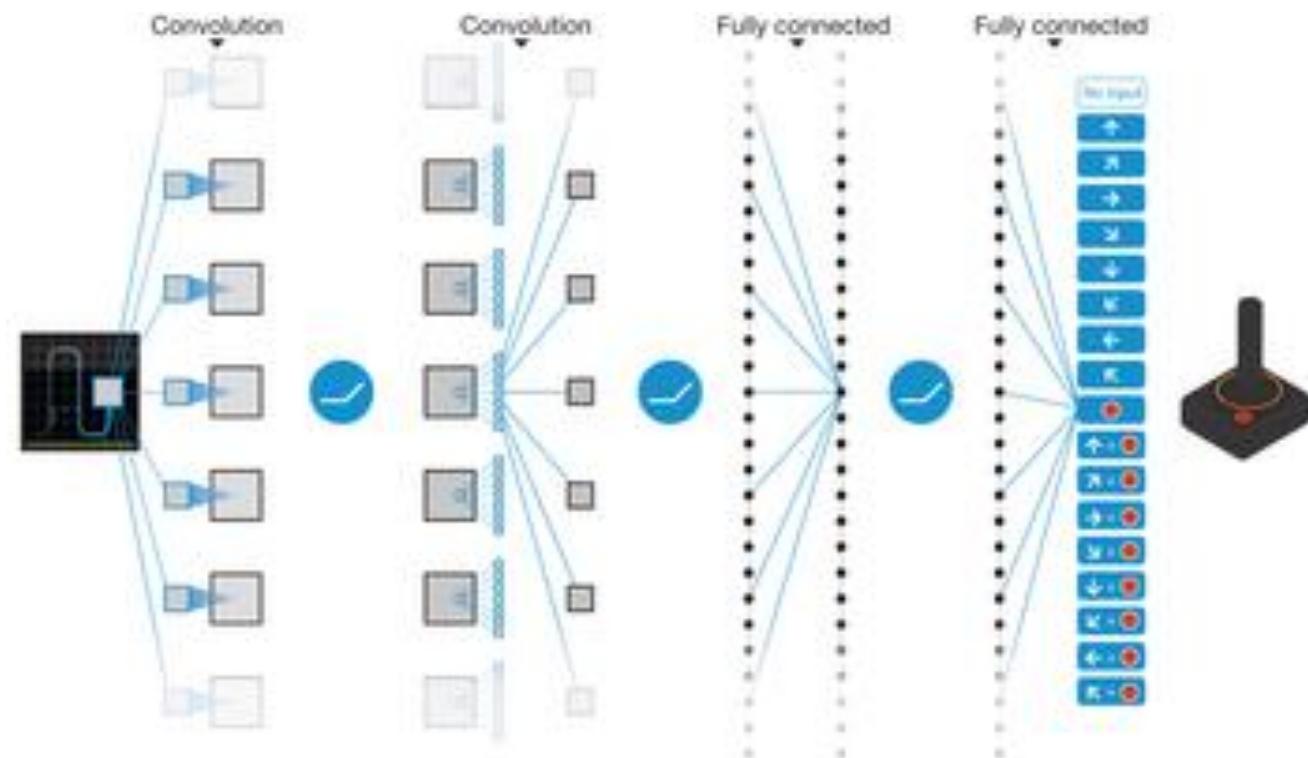
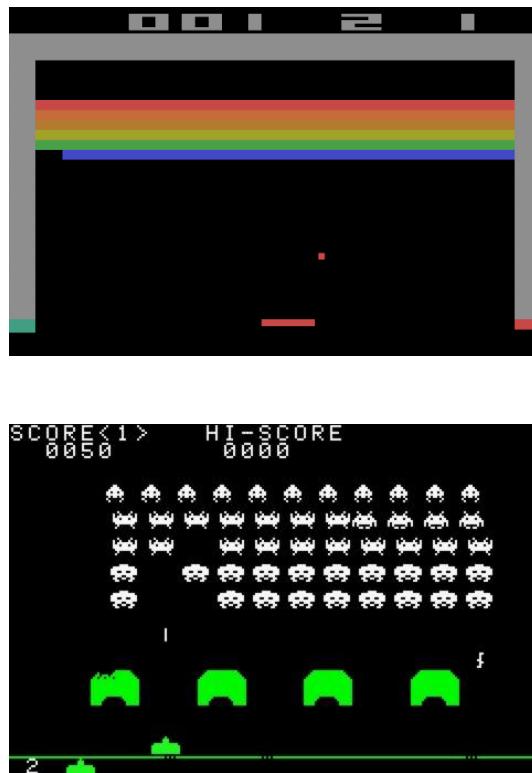
Q-learning with FA can easily be unstable in practice.
Some potentially helpful heuristics:

- Richer/safer representation [30]
- Two-network implementation [23]
- Experience replay [21, 23]
- Use an on-policy alternative such as Sarsa [30]
- ...

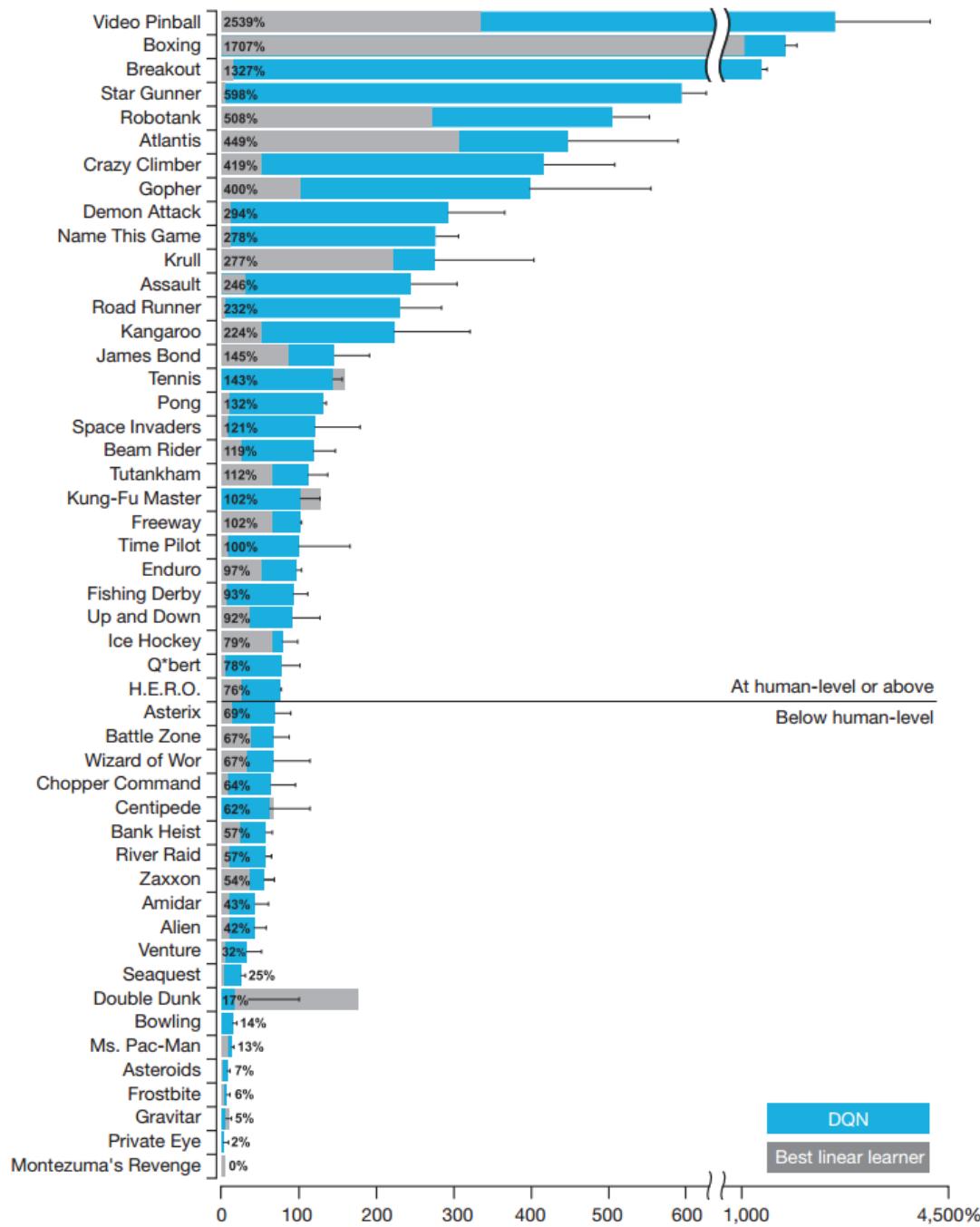
Are there provably stable algorithmic alternatives?

Human-Level Control via Deep RL

- ◆ Deep Q-network with human-level performance on Atari games



[Mnih et al., Nature 518, 529–533, 2015]



Demo

- ◆ Google DeepMind's DQN playing Atari Breakout

Demo

- ◆ Google DeepMind's DQN playing Atari Pacman



Google Deepmind DQN playing
Atari Pacman

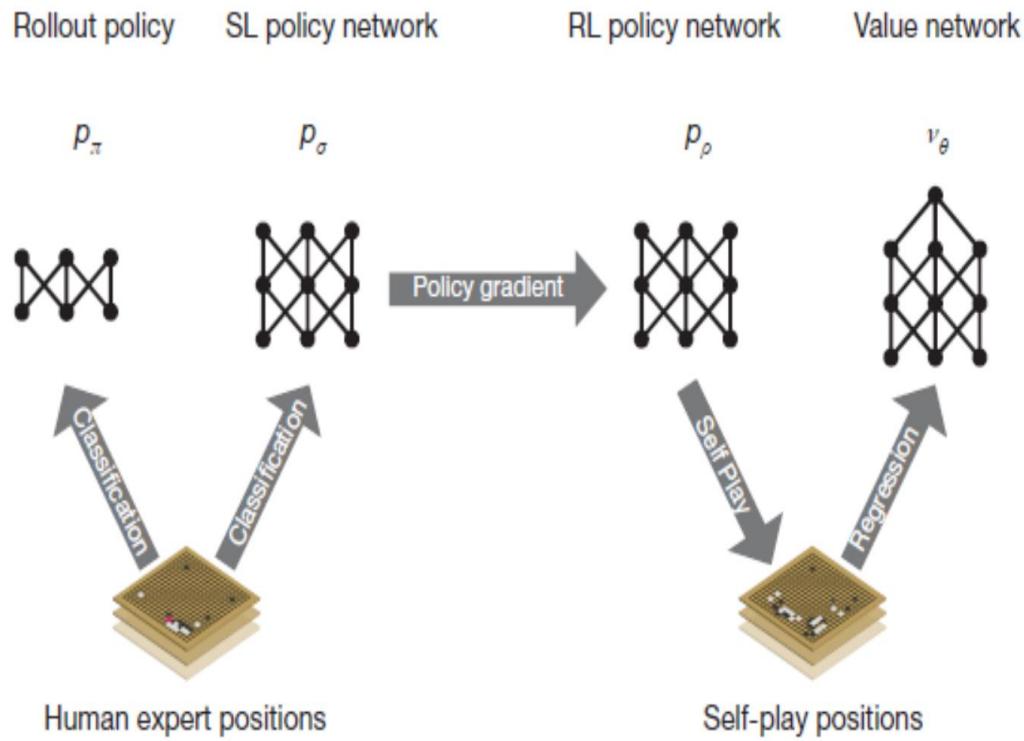
Setup:

NVIDIA GTX 690
i7-3770K - 16 GB RAM
Ubuntu 16.04 LTS
Google Deepmind DQN

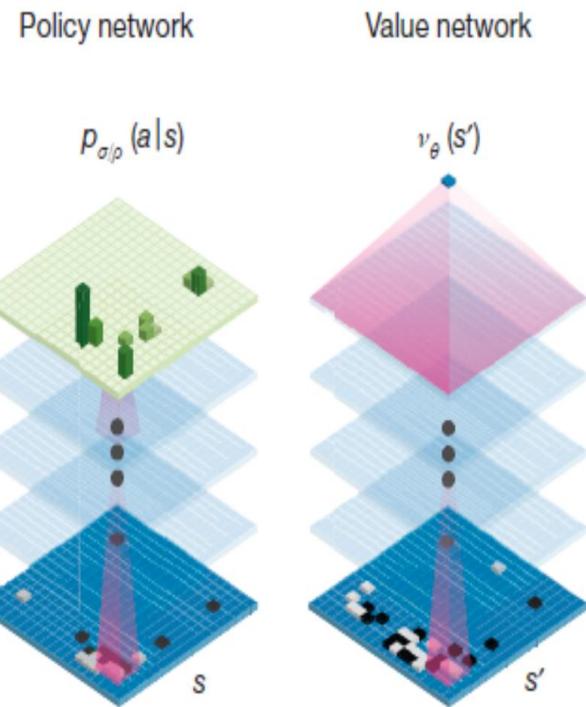
AlphaGo

- ◆ Neural network training pipeline and architecture

a

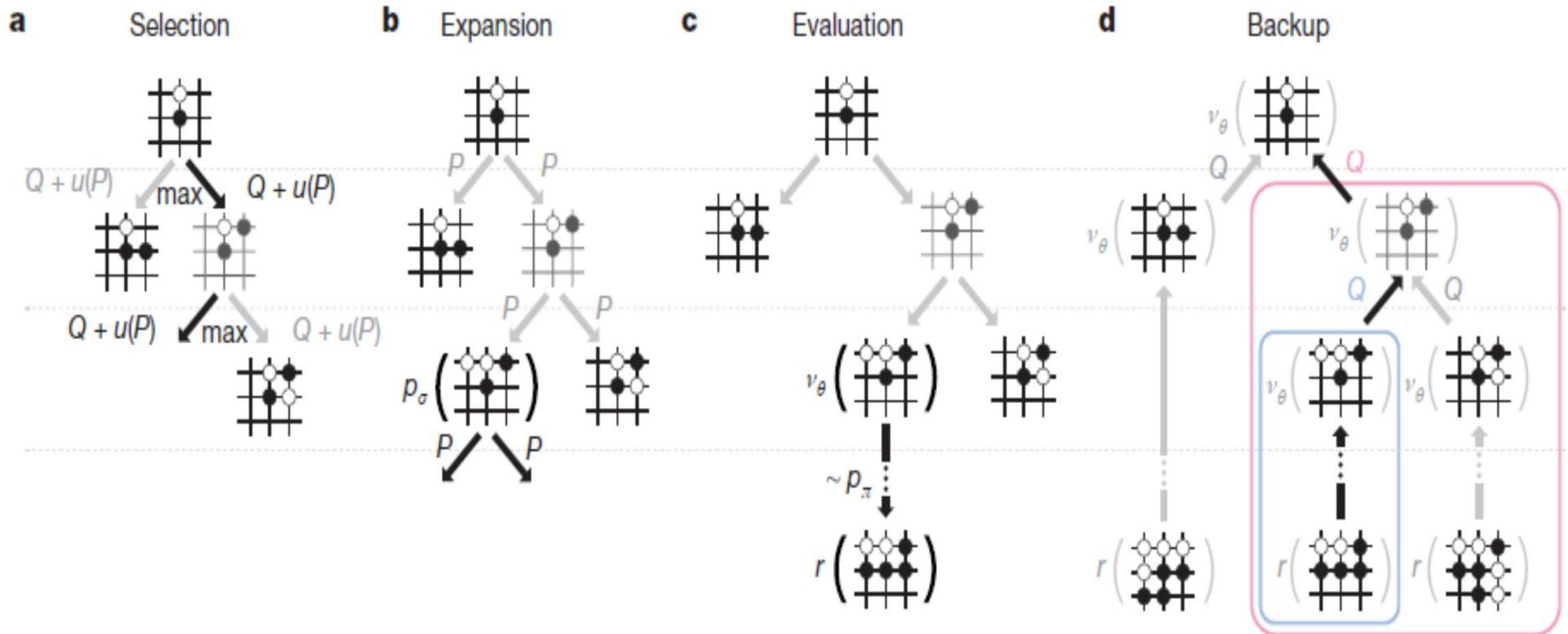


b



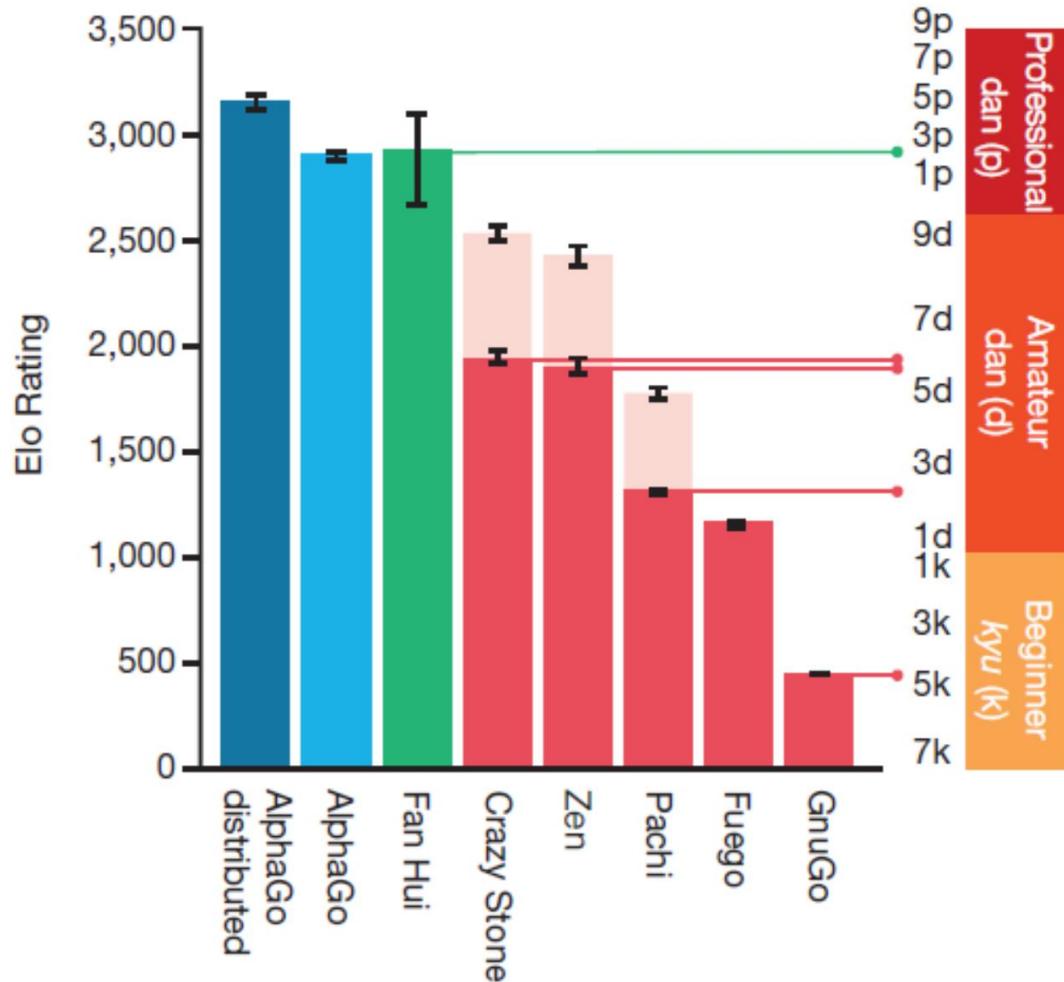
AlphaGo

◆ Monte Carlo tree search



[Silver et al., Mastering the game of Go with deep neural networks and tree search. Nature, 484(529), 2016]

AlphaGo



[Silver et al., Mastering the game of Go with deep neural networks and tree search. Nature, 484(529), 2016]



RL Algorithms: A Simplified Overview

| | Model (P & R) is known (aka “planning”) | Model is unknown (full RL) |
|------------------------|---|---|
| Value-function methods | value iteration policy iteration linear program | approx. VI (TD, Q-learning) approx. PI (LSPI, ...) approx. LP Monte Carlo estimation |
| Policy-search methods | Monte Carlo tree search Pegasus | Policy gradient (actor-critic) |

- A third dimension — representation
 - Tabular, aggregation, linear, kernel, decision trees, neural nets, ...
- Other topics not covered so far
 - Model-based RL
 - Hierarchical RL
 - Multi-task RL, transfer RL, continual RL
 - Multi-agent RL

Further Readings

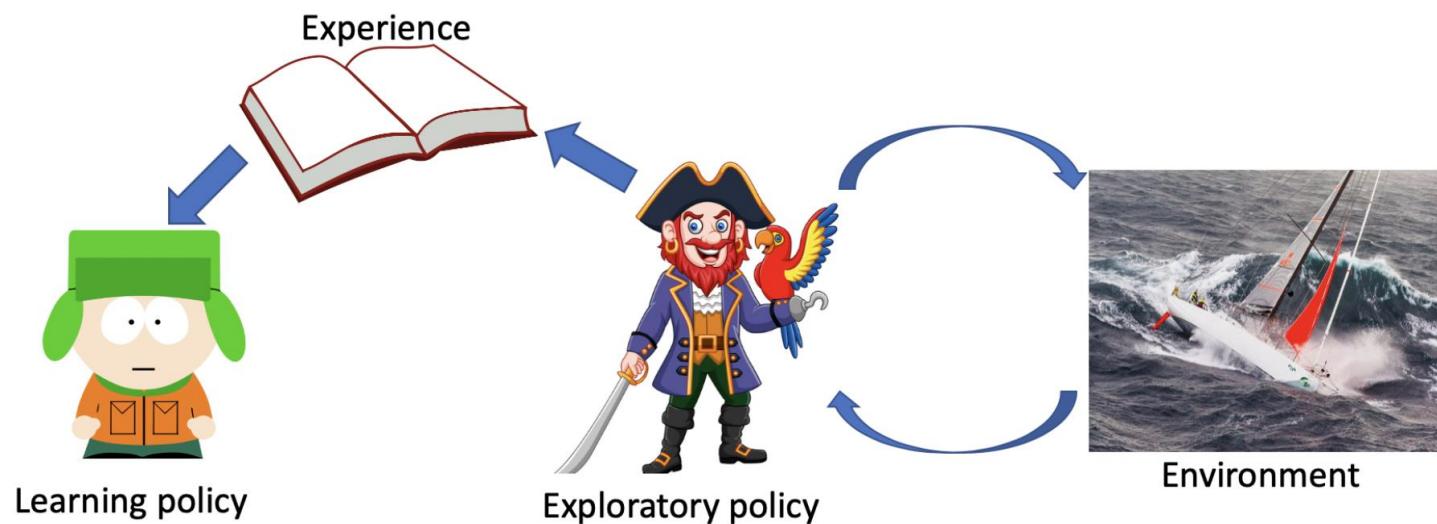
- Textbooks/monographs on RL in general: [4, 25, 31, 35, 40]
- Value function approximation: [5, 9, 37]
- Exploration: [3, 18, 26]
- Off-policy value estimation: [13, Chapter 4] [14]
- Hierarchical RL: [34]

Off-Policy Learning

- ◆ On-policy learning: Learn about policy π from the experience collected from π
 - Behave non-optimally in order to explore all actions, then reduce the exploration. e.g., ϵ -greedy
- ◆ Another important approach is **off-policy learning** which essentially uses **two different policies**:
 - the one which is being learned about and becomes the optimal policy
 - the other one which is more exploratory and is used to generate trajectories

Off-Policy Learning

- ◆ Off-policy learning: Learn about policy π from the experience sampled from another policy μ
 - π : target policy
 - μ : behavior policy



Off-Policy Control with Q-Learning

- ◆ We allow both behavior and target policies to improve
- ◆ The target policy π is greedy on $Q(s; a)$

$$\pi(S_{t+1}) = \arg \max_{a'} Q(S_{t+1}, a')$$

- ◆ The behavior policy μ could be totally random, but we let it improve following ϵ -greedy on $Q(s; a)$
- ◆ Thus Q-learning target:

$$\begin{aligned} R_{t+1} + \gamma Q(S_{t+1}, A') &= R_{t+1} + \gamma Q(S_{t+1}, \arg \max Q(S_{t+1}, a')) \\ &= R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') \end{aligned}$$

- ◆ Thus the Q-Learning update

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

References

-  Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire.
The nonstochastic multiarmed bandit problem.
SIAM Journal on Computing, 32(1):48–77, 2002.
-  Leemon C. Baird.
Residual algorithms: Reinforcement learning with function approximation.
In *Proceedings of the Twelfth International Conference on Machine Learning (ICML-95)*, pages 30–37, 1995.
-  Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos.
Unifying count-based exploration and intrinsic motivation.
In *Advances in Neural Information Processing Systems (NIPS-16)*, pages 1471–1479, 2016.
-  Dimitri P. Bertsekas and John N. Tsitsiklis.
Neuro-Dynamic Programming.
Athena Scientific, September 1996.
-  Vivek S. Borkar.
Stochastic Approximation: A Dynamical Systems Viewpoint.
Cambridge University Press, 2008.
-  Olivier Chapelle and Lihong Li.
An empirical evaluation of Thompson sampling.
In *Advances in Neural Information Processing Systems 24 (NIPS-11)*, pages 2249–2257, 2012.
-  Yichen Chen, Lihong Li, and Mengdi Wang.
Scalable bilinear π -learning using state and action features, 2018.
In preparation.

References



Bo Dai, Albert Shaw, Niao He, Lihong Li, and Le Song.

Boosting the actor with dual critic.

In *Proceedings of the Sixth International Conference on Learning Representations (ICLR-18)*, 2018.



Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song.

SBEED learning: Convergent control with nonlinear function approximation, 2018.

In preparation.



Daniela Pucci de Farias and Benjamin Van Roy.

The linear programming approach to approximate dynamic programming.

Operations Research, 51(6):850–865, 2003.



Miroslav Dudík, John Langford, and Lihong Li.

Doubly robust policy evaluation and learning.

In *Proceedings of the Twenty-Eighth International Conference on Machine Learning (ICML-11)*, pages 1097–1104, 2011.

CoRR abs/1103.4601.



Geoffrey J. Gordon.

Stable function approximation in dynamic programming.

In *Proceedings of the Twelfth International Conference on Machine Learning (ICML-95)*, pages 261–268, 1995.



Katja Hofmann, Lihong Li, and Filip Radlinski.

Online evaluation for information retrieval.

Foundations and Trends in Information Retrieval, 10(1):1–117, 2016.



Nan Jiang and Lihong Li.

Doubly robust off-policy evaluation for reinforcement learning.

In *Proceedings of the Thirty-Third International Conference on Machine Learning (ICML-16)*, 2016.

References

-  Levente Kocsis and Csaba Szepesvári.
Bandit based Monte-Carlo planning.
In *Proceedings of the Seventeenth European Conference on Machine Learning (ECML-06)*, pages 282–293, 2006.
-  Michail G. Lagoudakis and Ronald Parr.
Least-squares policy iteration.
Journal of Machine Learning Research, 4:1107–1149, 2003.
-  Chandrashekhar Lakshminarayanan, Shalabh Bhatnagar, and Csaba Szepesvári.
A linearly relaxed approximate linear program for Markov decision processes.
IEEE Transactions on Automatic Control, 63(4):1185–1191, 2017.
-  Lihong Li.
Sample complexity bounds of exploration.
In Marco Wiering and Martijn van Otterlo, editors, *Reinforcement Learning: State of the Art*, pages 175–204. Springer Verlag, 2012.
-  Lihong Li, Wei Chu, John Langford, and Robert E. Schapire.
A contextual-bandit approach to personalized news article recommendation.
In *Proceedings of the Nineteenth International Conference on World Wide Web (WWW-10)*, pages 661–670, 2010.
-  Lihong Li, Thomas J. Walsh, and Michael L. Littman.
Towards a unified theory of state abstraction for MDPs.
In *Proceedings of the Ninth International Symposium on Artificial Intelligence and Mathematics (ISAIM-06)*, 2006.
-  Long-Ji Lin.
Self-improving reactive agents based on reinforcement learning, planning and teaching.
Machine Learning, 8(3–4):293–321, 1992.

References



Hamid Reza Maei, Csaba Szepesvári, Shalabh Bhatnagar, and Richard S. Sutton.

Toward off-policy learning control with function approximation.

In *ICML*, pages 719–726, 2010.



Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis.

Human-level control through deep reinforcement learning.

Nature, 518:529–533, 2015.



Ronald Parr, Christopher Painter-Wakefield, Lihong Li, and Michael L. Littman.

Analyzing feature generation for value-function approximation.

In *Proceedings of the Twenty-Fourth International Conference on Machine Learning (ICML-07)*, pages 737–744, 2007.



Martin L. Puterman.

Markov Decision Processes: Discrete Stochastic Dynamic Programming.

Wiley-Interscience, New York, 1994.



Daniel J. Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen.

A tutorial on Thompson sampling, 2018.



Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever.

Evolution strategies as a scalable alternative to reinforcement learning, 2017.

arXiv:1703.03864.



Paul J. Schweitzer and Abraham Seidmann.

Generalized polynomial approximations in Markovian decision processes.

Journal of Mathematical Analysis and Applications, 110(2):568–582, 1985.

References

- 
- David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis.
Mastering the game of Go with deep neural networks and tree search.
Nature, 529:484–489, 2016.
- 
- Richard S. Sutton.
Generalization in reinforcement learning: Successful examples using sparse coarse coding.
In *Advances in Neural Information Processing Systems 8 (NIPS-95)*, pages 1038–1044, 1996.
- 
- Richard S. Sutton and Andrew G. Barto.
Reinforcement Learning: An Introduction.
MIT Press, Cambridge, MA, March 1998.
- 
- Richard S. Sutton and Andrew G. Barto.
Reinforcement Learning: An Introduction.
MIT Press, 2nd (draft) edition, 2018.
<https://drive.google.com/file/d/1xeUDVGWGUUv1-ccUMAZHJLej2C7aAFWY/view?usp=sharing>.
- 
- Richard S. Sutton, Hamid Maei, Doina Precup, Shalab Bhatnagar, Csaba Szepesvári, and Eric Wiewiora.
Fast gradient-descent methods for temporal-difference learning with linear function approximation.
In *Proceedings of the Twenty-Sixth International Conference on Machine Learning (ICML-09)*, pages 993–1000, 2009.
- 
- Richard S. Sutton, Doina Precup, and Satinder P. Singh.
Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning.
Artificial Intelligence, 112(1–2):181–211, 1999.
An earlier version appeared as Technical Report 98-74, Department of Computer Science, University of Massachusetts, Amherst, MA 01003. April, 1998.

References

-  Csaba Szepesvri.
Algorithms for Reinforcement Learning.
Morgan & Claypool, 2010.
-  Gerald Tesauro.
Temporal difference learning and TD-Gammon.
Communications of the ACM, 38(3):58–68, March 1995.
-  John N. Tsitsiklis and Benjamin Van Roy.
An analysis of temporal-difference learning with function approximation.
IEEE Transactions on Automatic Control, 42:674–690, 1997.
-  Mengdi Wang.
Primal-dual π learning: Sample complexity and sublinear run time for ergodic Markov decision problems, 2017.
CoRR abs/1710.06100.
-  Shimon Whiteson and Peter Stone.
Evolutionary function approximation for reinforcement learning.
Journal of Machine Learning Research, 7:877–917, 2006.
-  Marco Wiering and Martijn van Otterlo.
Reinforcement Learning: State of the Art.
Springer, 2012.
-  Yinyu Ye.
The simplex and policy-iteration methods are strongly polynomial for the Markov decision problem with a fixed discount rate.
Mathematics of Operations Research, 36(4):593–603, 2011.

- ◆ Thanks to Dr. Lihong Li for the slides content.