

Lecture 20

Reinforcement Learning

(Part II)

Lei Li, Yu-Xiang Wang

Announcements

- 1. Poster printing:** See the poster printing instructions from edstem.
- 2. Project Presentation day:** Small prizes will be given to the best poster presentations. You can invite your friends to come !
- 3. Course evaluation:** Please complete you ESCI surveys if you haven't yet. It takes only a few minutes.

Optional HW4

- For practice problems in convex optimization
 - For a simple coding problem with cvx: Q4 of [here](#) with data [here](#)
 - For practices on convex analysis: Q1,2,3 [here](#)
- Theory / concept practices for RL:
 - Problem 3 and 4 [here](#)
 - Problem 1 and 2 [here](#)
- Coding practice for MDP / RL: [Here](#)
- For more advanced problems in RL:
 - see HW1,2,3 from [my RL theory course](#).
 - These are only useful if you are hoping to do RL research.

Recap: Markov Decision processes (infinite horizon / discounted)

- Infinite horizon / discounted setting

$$\mathcal{M}(\mathcal{S}, \mathcal{A}, P, r, \gamma, \mu)$$

Transition kernel:

$$P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S}) \text{ i.e. } P(s'|s, a)$$

(Expected)

reward function:

$$r: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} / [0, R_{\max}] \quad \mathbb{E}[R_t | S_t=s, A_t=a] =: r(s, a)$$

Initial state distribution

$$\underline{\mu} \in \Delta(\mathcal{S})$$

Discounting factor:

$$\gamma$$

Stationary Policy π : mapping from state to an action (possibly a random action).

Recap: Value functions

- state value function: $V^\pi(s)$
 - **expected long-term** return when starting in s and following π

$$V^\pi(s) = \mathbb{E}_\pi[R_1 + \gamma R_2 + \dots + \gamma^{t-1} R_t + \dots | S_1 = s]$$

- state-action value function: $Q^\pi(s,a)$
 - **expected long-term** return when starting in s , performing a , and following π

$$Q^\pi(s, a) = \mathbb{E}_\pi[R_1 + \gamma R_2 + \dots + \gamma^{t-1} R_t + \dots | S_1 = s, A_1 = a]$$

Recap: Bellman equations

- Bellman consistency equation

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) [r(s, a, s') + \gamma V^\pi(s')]$$

$$V^\pi = r^\pi + \gamma P^\pi V^\pi$$

- Bellman optimality equation

$$V^*(s) = \max_a \sum_{s'} P(s'|s, a) [r(s, a, s') + \gamma V^*(s')]$$

$$Q = r + PV_Q \quad \text{where} \quad V_Q(s) := \max_{a \in \mathcal{A}} Q(s, a).$$

Recap: MDP planning and Value iterations

- MDP planning:

Find π^* such that $V^\pi(s) = V^*(s) \quad \forall s$

π is ϵ -optimal if $V^\pi \geq V^*(s) - \epsilon \mathbf{1}$

- Policy evaluation
 - Solving Bellman consistency equation
- Value iteration
 - Solving Bellman optimality equation

Recap: RL agent needs to learn the underlying MDP model

- **Model-based algorithm**

- Estimates the MDP then do MDP planning

- **Model-free algorithms**

- Monte Carlo Policy evaluation + Policy improvement
- Temporal difference learning = MC + Bellman equations

Recap: TD Learning

- TD-Policy evaluation

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

- TD-Policy optimization

- SARSA (on-policy)

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

Then choose the next A' using Q , e.g., eps-greedy.

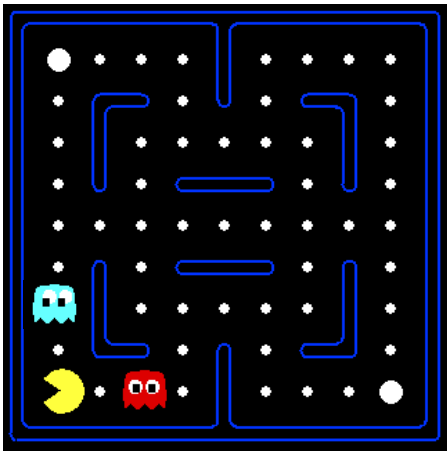
- Q-Learning (off-policy)

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

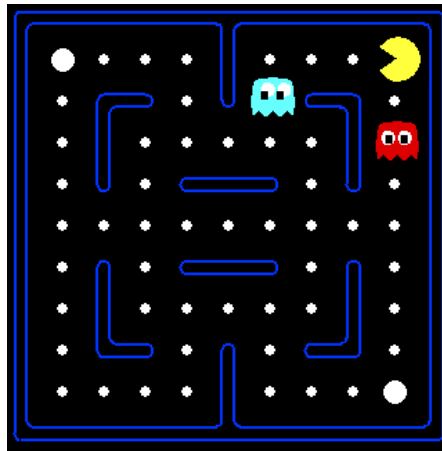
Then choose the next action in your favorite way.

Recap: The problem of large-state space

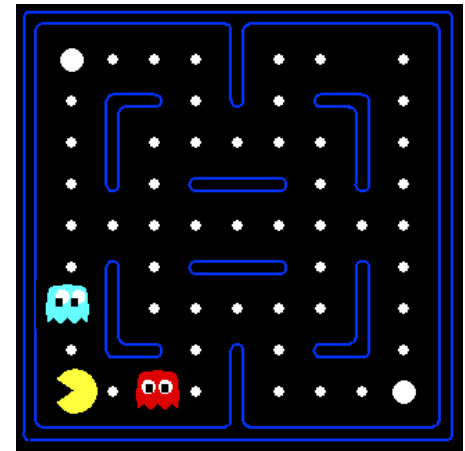
Let's say we discover through experience that this state is bad:



In naïve q-learning, we know nothing about this state:



Or even this one!



(From Dan Klein and Pieter Abbeel)

This lecture

- Solve the problem of large state space with function approximation
- Other RL algorithms: Policy gradient
- Exploration in RL

Video of Demo Q-Learning Pacman – Tiny – Watch All



Video of Demo Q-Learning Pacman – Tiny – Silent Train



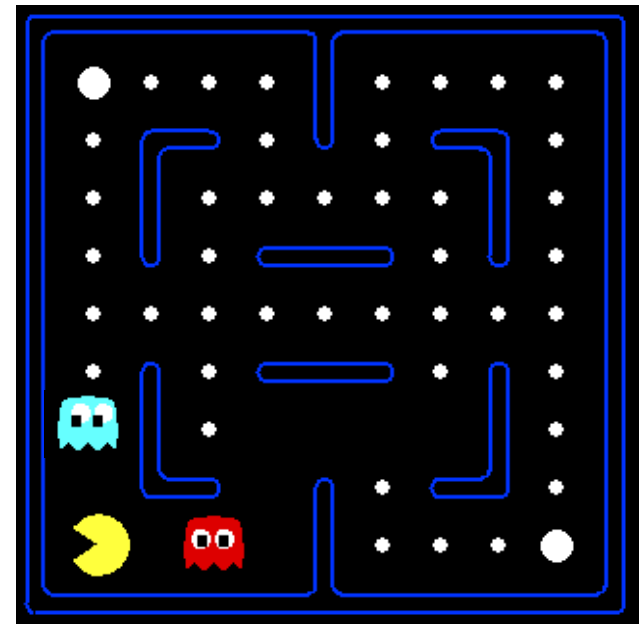
Video of Demo Q-Learning Pacman – Tricky –
Watch All



Why not use an evaluation function?

A Feature-Based Representations

- Solution: describe a state using a vector of features (properties)
 - Features are functions from states to real numbers (often 0/1) that capture important properties of the state
 - Example features:
 - Distance to closest ghost
 - Distance to closest dot
 - Number of ghosts
 - $1 / (\text{dist to dot})^2$
 - Is Pacman in a tunnel? (0/1)
 - etc.
 - Is it the exact state on this slide?
 - Can also describe a q-state (s, a) with features (e.g. action moves closer to food)



Linear Value Functions

- Using a feature representation, we can write a q function (or value function) for any state using a few weights:
 - $V_{\mathbf{w}}(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$
 - $Q_{\mathbf{w}}(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \dots + w_n f_n(s,a)$
- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but actually be very different in value!

Updating a linear value function

- Original Q learning rule tries to reduce prediction error at s, a :

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

- Instead, we update the weights to try to reduce the error at s, a :

$$\begin{aligned} w_i &\leftarrow w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] \partial Q_w(s,a) / \partial w_i \\ &= w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] f_i(s,a) \end{aligned}$$

Updating a linear value function

- Original Q learning rule tries to reduce prediction error at s, a :

$$Q(s,a) \leftarrow Q(s,a) + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

- Instead, we update the weights to try to reduce the error at s, a :

$$\begin{aligned} w_i &\leftarrow w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] \partial Q_w(s,a) / \partial w_i \\ &= w_i + \alpha \cdot [R(s,a,s') + \gamma \max_{a'} Q(s',a') - Q(s,a)] f_i(s,a) \end{aligned}$$

- Qualitative justification:
 - Pleasant surprise: increase weights on positive features, decrease on negative ones
 - Unpleasant surprise: decrease weights on positive features, increase on negative ones

PACMAN Q-Learning (Linear function approx.)



Deriving the TD via incremental optimization that minimizes Bellman errors

- Mean Square Error and Mean Square Bellman error

So far, in RL algorithms

- Model-based approaches
 - Estimate the MDP parameters.
 - Then use policy-iterations, value iterations.
- Monte Carlo methods:
 - estimating the rewards by empirical averages
- Temporal Difference methods:
 - Combine Monte Carlo methods with Dynamic Programming
- Linear function approximation in Q-learning
 - Similar to SGD
 - Learning heuristic function

Policy class and policy gradient methods

- Policy $\pi \in \Pi$

- Parametric policy class:

$$\Pi = \{\pi_{\theta} \mid \theta \in \mathbb{R}^d\}$$

- Goal: optimize the value
- Policy gradient methods
 - aim at learning the policy parameter by SGD.

Policy gradient

- Objective function to maximize: $J(\boldsymbol{\theta}) \doteq v_{\pi_{\boldsymbol{\theta}}}(s_0)$,
- Do SGD: $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \widehat{\nabla J(\boldsymbol{\theta}_t)}$,
- Policy gradient theorem:

$$\nabla J(\boldsymbol{\theta}) = \sum_s d^{\pi}(s) \sum_a Q^{\pi}(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

*Note how this theorem is non-trivial... The first two terms depends on π , but we did not take the gradient w.r.t. them.

Stochastic approximation in policy gradients

$$\nabla J(\boldsymbol{\theta}) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta})$$

- Sample from running policy π
 - $(S_1, A_1, R_1), \dots, (S_T, A_T, R_T)$
- Idea: Sample s , then the following is an unbiased estimator (finite horizon episodic case)

$$\begin{aligned} & \sum_{t=1}^T \left(\sum_{\ell=t}^T R_\ell \right) \frac{\nabla_{\boldsymbol{\theta}} \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \\ &= \sum_{t=1}^T G_t \nabla_{\boldsymbol{\theta}} \log(\pi(A_t|S_t, \boldsymbol{\theta})) \end{aligned}$$

***Show that this is an unbiased estimator of the gradient.**

Checkpoint for RL

- Model-based methods
- Model-free methods
 - Monte Carlo methods
 - TD-learning: Q-Learning and Sarsa
- Function approximation in RL
 - Approximate the MDP: Model-based
 - Approximate the value function
- Policy gradients
 - Parametrize the policy and run SGD

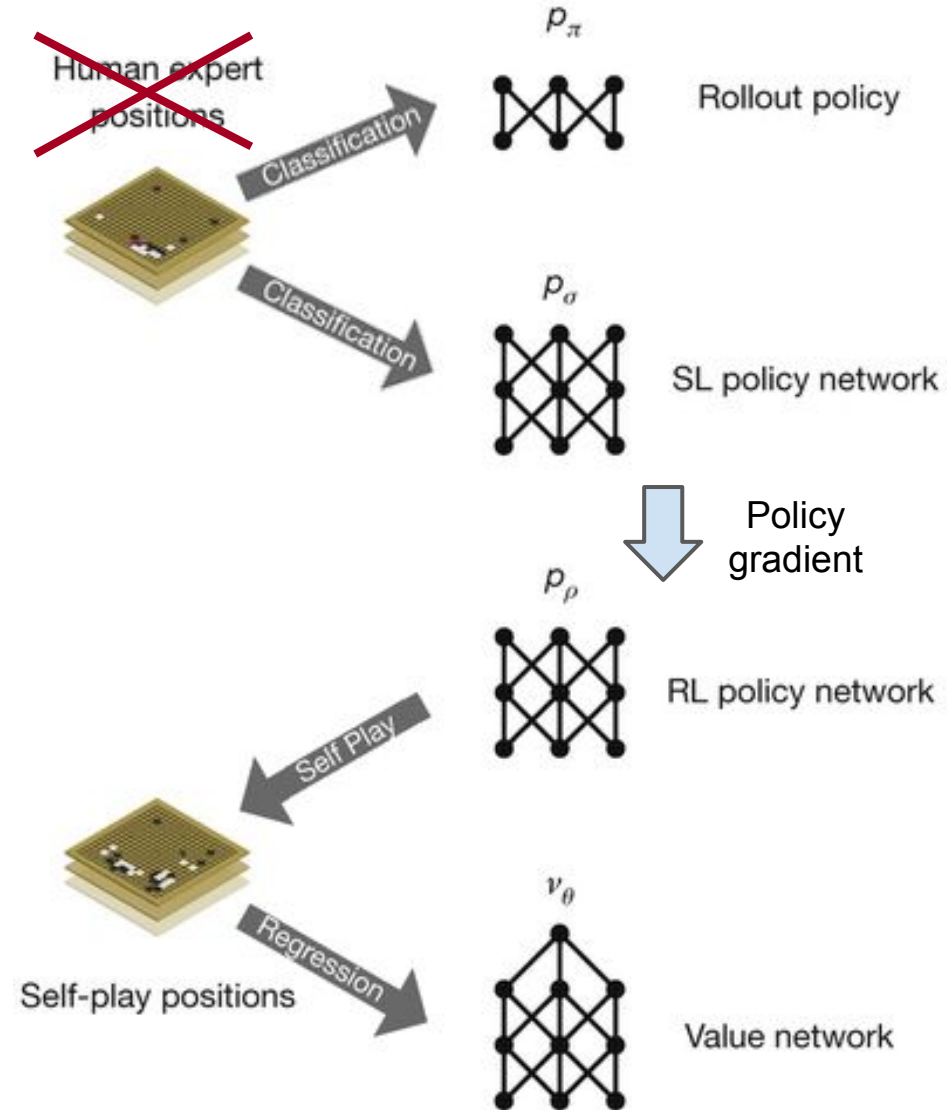
Elements of State-of-the-Art Reinforcement Learning

- Use a deep neural network to parameterize Q-function
- Use a deep neural network to parameterize the policy π
- Run a combination of Q-learning and Policy Gradient.
 - Actor-Critics, A3C, etc...
- Heuristic-based exploration: curiosity, reward shaping, etc..
- Experience replay to generate more data from existing data.
- Multi-agent RL: modeling your opponents

Alpha-Go and Alpha-Zero

- Parameterize the policy networks with CNN
- ~~• Supervised learning initialization~~
- RL using Policy gradient
- Fit Value Network (This is a value function approximation)
- Monte-Carlo Tree Search

<https://www.youtube.com/watch?v=4D5yGiYe8p4>



What I did not cover

- Useful results in RL for both theory and alg design
 - Simulation lemma
 - Advantage function and performance difference lemma
- Exploration
 - “Optimism in the face of uncertainty”
- Offline RL
 - “Pessimism in the face of uncertainty”
- How to start research in RL ?
 - Take my RL course (email me to ask for the videos)
 - Solve homework problems, implement RL algorithms from scratch.

Final words to students

- If you are doing theoretical research
 - It's useful have an empirical mind set
 - implement your algorithm, try it on examples (even toy examples would work)
 - These help you to challenge your assumptions and define theoretical problems that are useful
- If you are doing empirical research
 - Don't just chase SOTA in benchmarks
 - Think deeply about the problems you are working on
 - ML theory helps you to avoid pitfalls and design better algorithms.

Thank you! Looking forward to
your project presentations!