

So far... Supervised Learning

Lecture 14: Reinforcement Learning

Data: (x, y)
 x is data, y is label

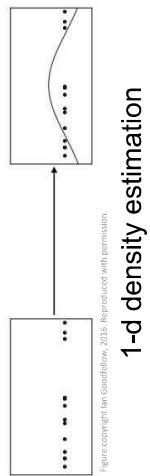


Goal: Learn a function to map $x \rightarrow y$

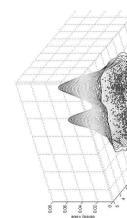
Examples: Classification,
regression, object detection,
semantic segmentation, image
captioning, etc.

So far... Unsupervised Learning

Data: x
Just data, no labels!



1-d density estimation



2-d density estimation

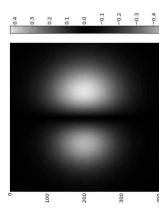
Examples: Clustering,
dimensionality reduction, feature
learning, density estimation, etc.

Today: Reinforcement Learning



Problems involving an **agent**,
interacting with an **environment**,
which provides numeric **reward**
signals

Goal: Learn how to take actions
in order to maximize reward

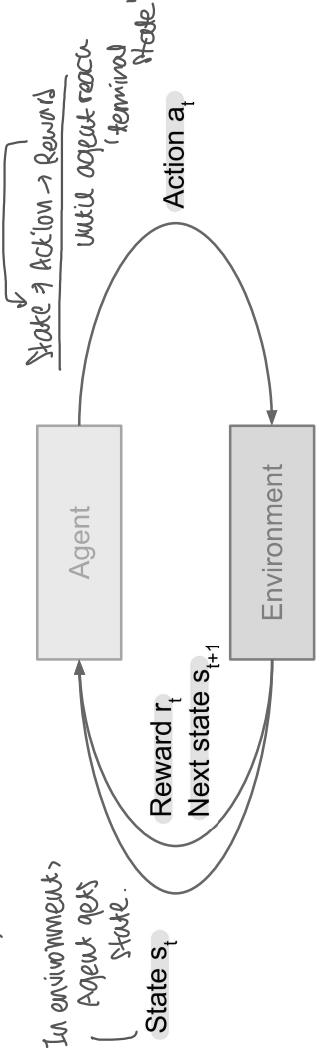


Overview

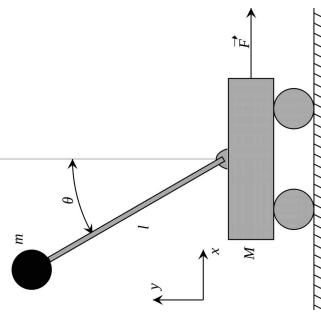
- What is Reinforcement Learning?
- Markov Decision Processes (MDP) – Formalism of Reinforcement Learning.
- Q-Learning [Method]
- Policy Gradients [Method]

Reinforcement Learning

↓ Agent & Environment



Cart-Pole Problem



Objective: Balance a pole on top of a movable cart.

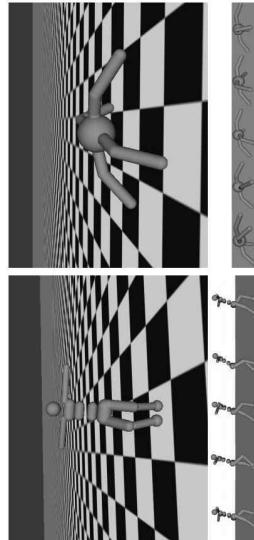
System:

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

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

Atari Games



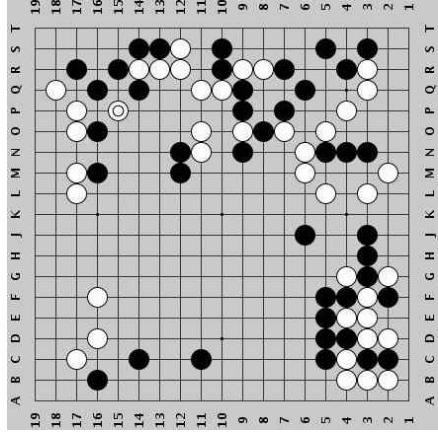
Objective: Complete the game with the highest score

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

Go (Reinforcement learning)



Figures copyright Volodymyr Mnih et al., 2013. Reproduced with permission.

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 16 May 23, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 17 May 23, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 17 May 23, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 17 May 23, 2017

Markov Decision Process

- Mathematical formulation of the RL problem
- **Markov property:** Current state completely characterises the state of the world

Defined by: $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathbb{P}, \gamma)$

\mathcal{S} : set of possible states

\mathcal{A} : set of possible actions

\mathcal{R} : distribution of reward given (state, action) pair

\mathbb{P} : transition probability i.e. distribution over next state given (state, action) pair

γ : discount factor

Markov Decision Process

- At time step t=0, environment samples initial state $s_0 \sim p(s_0)$
- Then, for t=0 until done:
 - Agent selects action a_t
 - Environment samples next state $s_{t+1} \sim P(\cdot | s_t, a_t)$
 - Agent receives reward r_t and next state s_{t+1}
 - ↳ can be deterministic or stochastic
- A policy π is a function from S to A that specifies what action to take in each state
- **Objective:** find policy π^* that maximizes cumulative discounted reward: $\sum_{t \geq 0} \gamma^t r_t$
 - ↳ include any reward that would be received in future

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 19 May 23, 2017

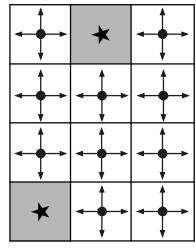
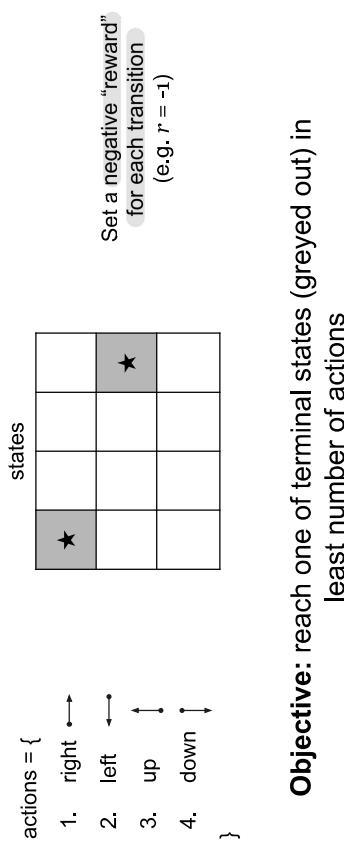
Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 20 May 23, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 20 May 23, 2017

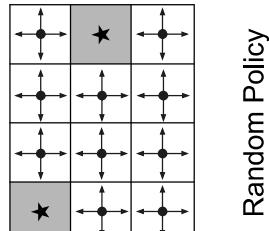
Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 20 May 23, 2017

A simple MDP: Grid World

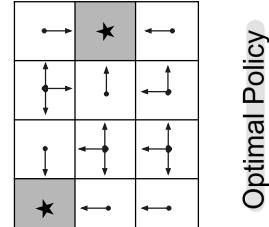
A simple MDP: Grid World



Objective: reach one of terminal states (greyed out) in
least number of actions



Objective: reach one of terminal states (greyed out) in
least number of actions



Optimal Policy

Random Policy

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

↙ How do we handle the randomness (initial state, transition probability...)?

→ Maximize the **expected sum of rewards!**

The optimal policy π^*

We want to find optimal policy π^* that maximizes the sum of rewards.

↙ How do we handle the randomness (initial state, transition probability...)?

Maximize the **expected sum of rewards!**

Formally: $\pi^* = \arg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi \right]$ with $s_0 \sim p(s_0)$, $a_t \sim \pi(\cdot | s_t)$, $s_{t+1} \sim p(\cdot | s_t, a_t)$

Definitions: Value function and Q-value function

Following a policy produces sample trajectories (or paths) $s_0, a_0, r_0, s_1, a_1, r_1, \dots$

\rightarrow چگونه از یک سکسیون

How good is a state?

The **value function** at state s , is the expected cumulative reward from following the policy from state s :

$$V^\pi(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi \right]$$

How good is a state-action pair?

The **Q-value function** at state s and action a , is the expected cumulative reward from taking action a in state s and then following the policy:

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

Intuition: if the optimal state-action values for the next time-step $Q^*(s', a')$ are known, then the optimal strategy is to take the action that maximizes the expected value of $(r + \gamma Q^*(s', a'))$.
این دستورات کدام قاعده را پیدا کردند
این دستورات کدام قاعده را پیدا کردند
این دستورات کدام قاعده را پیدا کردند
این دستورات کدام قاعده را پیدا کردند

Fei-Fei Li & Justin Johnson & Serena Yeung	Lecture 14 - 27	May 23, 2017	Fei-Fei Li & Justin Johnson & Serena Yeung	Lecture 14 - 30	May 23, 2017
--	-----------------	--------------	--	-----------------	--------------

Bellman equation

The optimal Q-value function Q^* is the maximum expected cumulative reward achievable from a given (state, action) pair:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

Q^* satisfies the following **Bellman equation**:

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

Value iteration algorithm: Use Bellman equation as an iterative update

$$Q_{i+1}(s, a) = \mathbb{E} \left[r + \gamma \max_{a'} Q_i(s', a') | s, a \right]$$

Q_i will converge to Q^* as $i \rightarrow \infty$

What's the problem with this?

Not scalable. Must compute $Q(s, a)$ for every state-action pair. If state is e.g. current game state pixels, computationally infeasible to compute for entire state space!

Solution: use a function approximator to estimate $Q(s, a)$. E.g. a neural network!
این روش

Solving for the optimal policy: Q-learning

Q-learning: Use a function approximator to estimate the action-value function

$$Q(s, a; \theta) \approx Q^*(s, a)$$

function parameters (weights)

If the function approximator is a deep neural network \Rightarrow **deep q-learning**!

Bellman function بلمن
 \rightarrow بلمن این را می بینیم

Fei-Fei Li & Justin Johnson & Serena Yeung	Lecture 14 - 34	May 23, 2017	Fei-Fei Li & Justin Johnson & Serena Yeung	Lecture 14 - 37	May 23, 2017
--	-----------------	--------------	--	-----------------	--------------

Solving for the optimal policy: Q-learning

Remember: want to find a Q-function that satisfies the Bellman Equation:

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

Forward Pass \rightarrow calculate loss function. (minimize \mathbb{V})

$$\text{Loss function: } L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$

Backward Pass \rightarrow update parameter θ (based on calculated y_i)

Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right] \nabla_{\theta_i} Q(s, a; \theta_i)$$

$\nabla_{\theta_i} \rightarrow$ target for Q-function or $\nabla_{\theta_i} \rightarrow$ gradient of Q-function

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 41 May 23, 2017 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 42 May 23, 2017 Fei-Fei Li & Justin Johnson & Serena Yeung



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

Mnih et al. NIPS Workshop 2013; Nature 2015. Reproduced with permission.

Case Study: Playing Atari Games

Q-network Architecture

Q, ∇_{θ_i} , FC layer \rightarrow $\nabla_{\theta_i} Q$.

Output of FC: Q values.



Current state s_t : 84x84x4 stack of last 4 frames

(after RGB->grayscale conversion, downsampling, and cropping)

Action selection \rightarrow $\nabla_{\theta_i} Q$ (output)

Q-network Architecture

Q, ∇_{θ_i} , FC layer \rightarrow $\nabla_{\theta_i} Q$.

Output of FC: Q values.



Current state s_t : 84x84x4 stack of last 4 frames

(after RGB->grayscale conversion, downsampling, and cropping)

Action selection \rightarrow $\nabla_{\theta_i} Q$ (output)

Mnih et al. NIPS Workshop 2013; Nature 2015

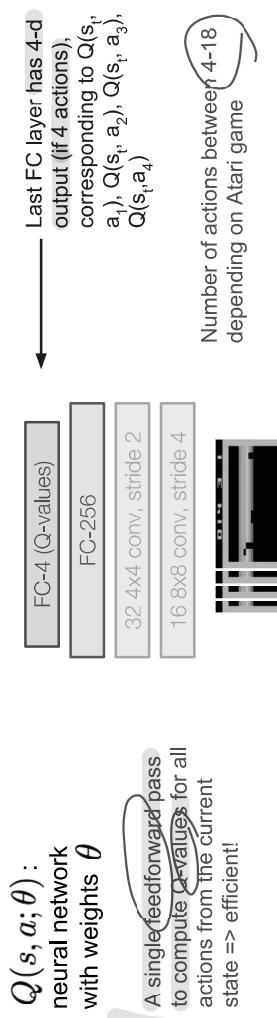
Mnih et al. NIPS Workshop 2013; Nature 2015

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 44 May 23, 2017 Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 14 - 45 May 23, 2017 Fei-Fei Li & Justin Johnson & Serena Yeung

Q-network Architecture



Continual update of transitions (s_t, a_t, r_t, s_{t+1}) as game episodes are played

Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)] \nabla_{\theta_i} Q(s, a; \theta_i)$$

Training the Q-network: Loss function (from before)

Remember: want to find a Q-function that satisfies the Bellman Equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q^*(s', a') | s, a]$$

Forward Pass

$$\text{Loss function: } L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a; \theta_i))^2]$$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$

Backward Pass ~~backward gradient~~ \rightarrow NW ~~gradient~~ \rightarrow Gradient update (with respect to Q-function parameters θ):

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)] \nabla_{\theta_i} Q(s, a; \theta_i)$$

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 48 May 23, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 49 May 23, 2017

Mnih et al. NIPS Workshop 2013; Nature 2015

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

Shade, action, reward, next state

Address these problems using **experience replay**

- Continually update a **replay memory** table of transitions (s_t, a_t, r_t, s_{t+1}) as game (experience) episodes are played
- Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples

Mnih et al. NIPS Workshop 2013; Nature 2015

Training the Q-network: Experience Replay

Learning from batches of consecutive samples is problematic:

- Samples are correlated => inefficient learning
- Current Q-network parameters determines next training samples (e.g. if maximizing action is to move left, training samples will be dominated by samples from left-hand size) => can lead to bad feedback loops

- Address these problems using **experience replay**
 - Continually update a **replay memory** table of transitions (s_t, a_t, r_t, s_{t+1}) as game (experience) episodes are played
 - Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples
- Each transition can also contribute to ~~multiple~~ weight updates => greater data efficiency

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 51 May 23, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 14 - 52 May 23, 2017

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1,  $M$  do
    Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ 
    for  $t = 1, T$  do
        With probability  $\epsilon$  select a random action  $a_t$ 
        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
    end for
end for

```

Putting it together: Deep Q-Learning with Experience Replay

Algorithm 1 Deep Q-learning with Experience Replay

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1,  $M$  do
    Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ 
    for  $t = 1, T$  do
        With probability  $\epsilon$  select a random action  $a_t$ 
        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Take the action  $(a_t,$ 
        and observe the
        reward  $r_t$ 
        and next
        state  $s_{t+1}$ 
        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
        for terminal  $\phi_{j+1}$ 
        for non-terminal  $\phi_{j+1}$ 
        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
    end for
end for

```

Policy Gradients

What is a problem with Q-learning?
 The Q-function can be very complicated! \curvearrowleft

Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair

Policy Gradients

What is a problem with Q-learning?
 The Q-function can be very complicated!

Example: a robot grasping an object has a very high-dimensional state => hard to learn exact value of every (state, action) pair

But the policy can be much simpler: just close your hand
 Can we learn a policy directly, e.g. finding the best policy from a collection of policies?
 \oplus value X . Policy directly

Policy Gradients

Formally, let's define a class of parametrized policies: $\Pi = \{\pi_\theta, \theta \in \mathbb{R}^m\}$

For each policy, define its value:

$$J(\theta) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi_\theta \right]$$

Expected value of rewards that would be received.

We want to find the optimal policy $\theta^* = \arg \max_{\theta} J(\theta)$

How can we do this?

Gradient ascent on policy parameters!

REINFORCE algorithm

Mathematically, we can write:

$$\begin{aligned} J(\theta) &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)] \\ &= \int_{\tau} r(\tau) p(\tau; \theta) d\tau \quad (\text{Intractable! } \mathbb{E}_{\tau \sim p(\tau; \theta)} \rightarrow \text{tricky}) \end{aligned}$$

Where $r(\tau)$ is the reward of a trajectory $\tau = (s_0, a_0, r_0, s_1, \dots)$

REINFORCE algorithm

Expected reward: $J(\theta) = \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau)]$

$$= \int_{\tau} r(\tau) p(\tau; \theta) d\tau$$

Now let's differentiate this: $\nabla_{\theta} J(\theta) = \int_{\tau} r(\tau) \nabla_{\theta} p(\tau; \theta) d\tau$

However, we can use a nice trick: $\nabla_{\theta} p(\tau; \theta) = p(\tau; \theta) \frac{\nabla_{\theta} p(\tau; \theta)}{p(\tau; \theta)} = p(\tau; \theta) \nabla_{\theta} \log p(\tau; \theta)$
If we inject this back:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau \\ &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)] \end{aligned}$$

Can estimate with Monte Carlo sampling

REINFORCE algorithm

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \int_{\tau} (r(\tau) \nabla_{\theta} \log p(\tau; \theta)) p(\tau; \theta) d\tau \\ &= \mathbb{E}_{\tau \sim p(\tau; \theta)} [r(\tau) \nabla_{\theta} \log p(\tau; \theta)]\end{aligned}$$

Can we compute those quantities without knowing the transition probabilities?

We have: $p(\tau; \theta) = \prod_{t \geq 0} p(s_{t+1} | s_t, a_t) \pi_{\theta}(a_t | s_t)$

$$\text{Thus: } \log p(\tau; \theta) = \sum_{t \geq 0} \log p(s_{t+1} | s_t, a_t) + \log \pi_{\theta}(a_t | s_t)$$

$$\text{And when differentiating: } \nabla_{\theta} \log p(\tau; \theta) = \sum_{t \geq 0} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Therefore when sampling a trajectory τ , we can estimate $J(\theta)$ with

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Intuition

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Interpretation:

- If $r(\tau)$ is high, push up the probabilities of the actions seen
- If $r(\tau)$ is low, push down the probabilities of the actions seen

Might seem simplistic to say that if a trajectory is good then all its actions were good. But in expectation, it averages out!

However, this also suffers from high variance because credit assignment is really hard. Can we help the estimator?

Variance reduction

$$\text{Gradient estimator: } \nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

First idea: Push up probabilities of an action seen, only by the cumulative future reward from that state

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Second idea: Use discount factor γ to ignore delayed effects

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t' - t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

Variance reduction: Baseline

Problem: The raw value of a trajectory isn't necessarily meaningful. For example, if rewards are all positive, you keep pushing up probabilities of actions.

What is important then? Whether a reward is better or worse than what you expect to get

Idea: Introduce a baseline function dependent on the state. Concretely, estimator is now:

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t' - t} r_{t'} - b(s_t) \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

How to choose the baseline?

$$\nabla_\theta J(\theta) \approx \sum_{t \geq 0} \left(\sum_{t' \geq t} \gamma^{t'-t} r_t \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \nabla_\theta \log \pi_\theta(a_t | s_t)$$

A simple baseline: constant moving average of rewards experienced so far from all trajectories

Variance reduction techniques seen so far are typically used in “Vanilla REINFORCE”

How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

Intuitively, we are happy with an action a_t in a state s_t if $Q^\pi(s_t, a_t) - V^\pi(s_t)$ is large. On the contrary, we are unhappy with an action if it's small.

Using this, we get the estimator: $\nabla_\theta J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t)) \nabla_\theta \log \pi_\theta(a_t | s_t)$

Actor-Critic Algorithm

Problem: we don't know Q and V. Can we learn them?

Yes, using Q-learning! We can combine Policy Gradients and Q-learning by training both an **actor** (the policy) and a **critic** (the Q-function).

- The actor decides which action to take, and the critic tells the actor how good its action was and how it should adjust
- Also alleviates the task of the critic as it only has to learn the values of (state, action) pairs generated by the policy
- Can also incorporate Q-learning tricks e.g. experience replay
- **Remark:** we can define by the **advantage function** how much an action was better than expected

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$$

Actor-Critic Algorithm

Initialize policy parameters θ , critic parameters ϕ
For iteration=1, 2 ... **do**
 Sample m trajectories under the current policy

$\Delta\theta \leftarrow 0$
 For i=1, ..., m **do**
 $A_t = \sum_{t' \geq t} \gamma^{t'-t} r_t^i - V_\phi(s_t^i)$
 $\Delta\theta \leftarrow \Delta\theta + A_t \nabla_\theta \log(a_t^i | s_t^i)$
 End for
 $\Delta\phi \leftarrow \sum_t \sum \nabla_\phi ||A_t^i||^2$
 $\theta \leftarrow \alpha \Delta\theta$
 $\phi \leftarrow \beta \Delta\phi$
End for

Summary

- **Policy gradients:** very general but suffer from high variance so requires a lot of samples. **Challenge:** sample-efficiency
- **Q-learning:** does not always work but when it works, usually more sample-efficient. **Challenge:** exploration
- Guarantees:
 - **Policy Gradients:** Converges to a local minima of $J(\theta)$, often good enough!
 - **Q-learning:** Zero guarantees since you are approximating Bellman equation with a complicated function approximator

Next Time

- Guest Lecture: Song Han**
 - Energy-efficient deep learning
 - Deep learning hardware
 - Model compression
 - Embedded systems
 - And more...