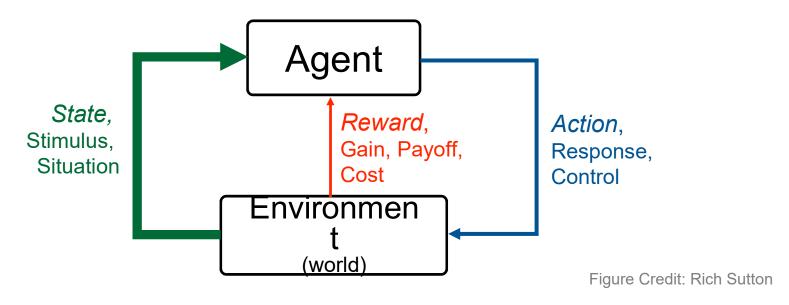
Topics:

- Reinforcement Learning Part 2
 - Q-Learning
 - Deep Q-Learning

CS 4803-DL / 7643-A ZSOLT KIRA

RL: Sequential decision making in an environment with evaluative feedback.



- **Environment** may be unknown, non-linear, stochastic and complex.
- Agent learns a policy to map states of the environments to actions.
 - Seeking to maximize cumulative reward in the long run.



- MDPs: Theoretical framework underlying RL
- lacktriangle An MDP is defined as a tuple $(\mathcal{S},\mathcal{A},\mathcal{R},\mathbb{T},\gamma)$

 ${\cal S}$: Set of possible states

 ${\cal A}\,$: Set of possible actions

 $\mathcal{R}(s,a,s')$: Distribution of reward

 $\mathbb{T}(s,a,s')$: Transition probability distribution, also written as p(s'|s,a)

 γ : Discount factor

- MDPs: Theoretical framework underlying RL
- lacktriangle An MDP is defined as a tuple $(\mathcal{S},\mathcal{A},\mathcal{R},\mathbb{T},\gamma)$

 ${\cal S}$: Set of possible states

 ${\cal A}\,$: Set of possible actions

 $\mathcal{R}(s,a,s')$: Distribution of reward

 $\mathbb{T}(s,a,s')$: Transition probability distribution, also written as p(s'|s,a)

 γ : Discount factor

Interaction trajectory: $\ldots s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}, r_{t+2}, s_{t+2}, \ldots$

What we want

e.g. A policy π State Action

$$\pi^* = rg \max_{\pi} \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | \pi
ight]$$

Definition of optimal policy

Some intermediate concepts and terms

A **Value function** (how good is a state?)

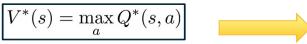
$$V: \mathcal{S}
ightarrow \mathbb{R} \quad V^{\pi}(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, \pi
ight]$$

A Q-Value function (how good is a state-action pair?)

$$Q: \mathcal{S} \times \mathcal{A} \to \mathbb{R} \quad Q^{\pi}(s, a) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t r_t | s_0 = s, a_0 = a, \pi \right]$$

$$Q^*(s,a) = \underset{\sim p(s'|s,a)}{\mathbb{E}} [r(s,a) + \gamma V^*(s')]$$
 (Math in previous lecture)

Equalities relating optimal quantities



$$\pi^*(s) = \arg\max_{a} Q^*(s, a)$$

We can then derive the Bellman Equation

$$Q^*(s, a) = \sum_{s'} p(s'|s, a) \left[r(s, a) + \gamma \max_{a} Q^*(s', a') \right]$$

This must hold true for an optimal Q-Value!

-> Leads to dynamic programming algorithm to find it

Value Iteration Update:

$$V^{i+1}(s) \leftarrow \max_{a} \sum_{s'} p(s'|s,a) \left[r(s,a) + \gamma V^{i}(s') \right]$$

Q-Iteration Update:

$$Q^{i+1}(s,a) \leftarrow \sum_{s'} p\left(s'|s,a\right) \left[r\left(s,a\right) + \gamma \max_{a'} Q^{i}(s',a')\right]$$

The algorithm is same as value iteration, but it loops over actions as well as states



For Value Iteration:

Theorem: will converge to unique optimal values
Basic idea: approximations get refined towards optimal values
Policy may converge long before values do

Time complexity per iteration $O(|\mathcal{S}|^2|\mathcal{A}|)$

Feasible for:

- 3x4 Grid world?
- Chess/Go?
- Atari Games with integer image pixel values [0, 255] of size 16x16 as state?



Summary: MDP Algorithms

Value Iteration

 Bellman update to state value estimates

Q-Value Iteration

Bellman update to (state, action) value estimates



Reinforcement Learning, Deep RL



- Recall RL assumptions:
 - $\mathbb{T}(s, a, s')$ unknown, how actions affect the environment.
 - ullet $\mathcal{R}(s,a,s')$ unknown, what/when are the good actions?
- But, we can learn by trial and error.
 - Gather experience (data) by performing actions.
 - Approximate unknown quantities from data.

Reinforcement Learning



- Old Dynamic Programming Demo
 - https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_dp.html
- RL Demo
 - https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_td.html

Slide credit: Dhruv Batra



Sample-Based Policy Evaluation?

We want to improve our estimate of V by computing these averages:

$$V_{k+1}^{\pi}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_k^{\pi}(s')]$$

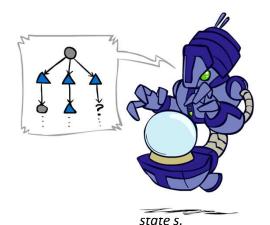
• Idea: Take samples of outcomes s' (by doing the action!) and average

$$sample_{1} = R(s, \pi(s), s'_{1}) + \gamma V_{k}^{\pi}(s'_{1})$$

$$sample_{2} = R(s, \pi(s), s'_{2}) + \gamma V_{k}^{\pi}(s'_{2})$$
...
$$sample_{n} = R(s, \pi(s), s'_{n}) + \gamma V_{k}^{\pi}(s'_{n})$$

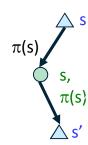
$$V_{k+1}^{\pi}(s) \leftarrow \frac{1}{n} \sum_{i} sample_{i}$$

What's the difficulty of this algorithm?



Temporal Difference Learning

- Big idea: learn from every experience!
 - Update V(s) each time we experience a transition (s, a, s', r)
 - Likely outcomes s' will contribute updates more often
- Temporal difference learning of values
 - Policy still fixed, still doing evaluation!
 - Move values toward value of whatever successor occurs: running average



Sample of V(s):
$$sample = R(s, \pi(s), s') + \gamma V^{\pi}(s')$$

Update to V(s):
$$V^{\pi}(s) \leftarrow (1-\alpha)V^{\pi}(s) + (\alpha)sample$$

Same update:
$$V^{\pi}(s) \leftarrow V^{\pi}(s) + \alpha(sample - V^{\pi}(s))$$

Q-Learning

• We'd like to do Q-value updates to each Q-state:

$$Q_{k+1}(s,a) \leftarrow \sum_{s'} T(s,a,s') \left[R(s,a,s') + \gamma \max_{a'} Q_k(s',a') \right]$$

- But can't compute this update without knowing T, R
- Instead, compute average as we go
 - Receive a sample transition (s,a,r,s')
 - This sample suggests

$$Q(s, a) \approx r + \gamma \max_{a'} Q(s', a')$$

- But we want to average over results from (s,a)
- So keep a running average

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[r + \gamma \max_{a'} Q(s', a')\right]$$

Q-Learning Properties

- Amazing result: Q-learning converges to optimal policy -- even if you're acting suboptimally!
- This is called off-policy learning
- Caveats:
 - You have to explore enough
 - You have to eventually make the learning rate small enough
 - ... but not decrease it too quickly
 - Basically, in the limit, it doesn't matter how you select action



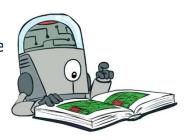


Deep Q-Learning



Generalizing Across States

- Basic Q-Learning keeps a table of all q-values
- In realistic situations, we cannot possibly learn about every single state
 - Too many states to visit them all in training
 - Too many states to hold the q-tables in memory
- Instead, we want to generalize:
 - Learn about some small number of training states from experience
 - Generalize that experience to new, similar situations
 - This is the fundamental idea in machine learning!



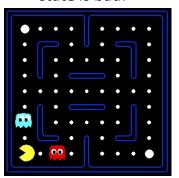


[demo - RL pacman]



Example: Pacman

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!





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 / (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:

$$V(s) = w_1 f_1(s) + w_2 f_2(s) + \ldots + w_n f_n(s)$$

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \dots + w_n f_n(s,a)$$

- Advantage: our experience is summed up in a few powerful numbers
- Disadvantage: states may share features but can actually be very different in value!

- State space is too large and complicated for feature engineering though!
- Recall: Value iteration not scalable (chess, RGB images as state space, etc)
- Solution: Deep Learning! ... more precisely, function approximation.
 - Use deep neural networks to learn state representations
 - Useful for continuous action spaces as well

Deep Reinforcement Learning

Georgia Tech

Value-based RL

(Deep) Q-Learning, approximating $Q^*(s,a)$ with a deep Q-network

Policy-based RL

ullet Directly approximate optimal policy π^* with a parametrized policy $\pi^*_{ heta}$

Model-based RL

- lacktriangle Approximate transition function T(s',a,s) and reward function $\mathcal{R}(s,a)$
- Plan by looking ahead in the (approx.) future!



Q-Learning with linear function approximators

$$Q(s, a; w, b) = w_a^{\top} s + b_a$$

- Has some theoretical guarantees
- Deep Q-Learning: Fit a deep Q-Network $\,Q(s,a; heta)\,$
 - Works well in practice
 - Q-Network can take RGB images

FC-4 (Q-values)

FC-256

32 4x4 conv, stride 2

16 8x8 conv, stride 4

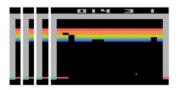


Image Credits: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



Assume we have collected a dataset:

$$\{(s, a, s', r)_i\}_{i=1}^N$$

We want a Q-function that satisfies bellman optimality (Q-value)

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

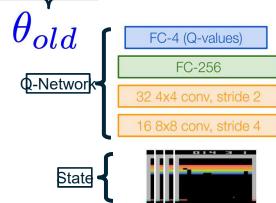
Loss for a single data point:

$$\text{MSE Loss} := \left(\frac{Q_{new}(s, a) - (r + \gamma \max_{a} Q_{old}(s', a))}{\text{Predicted Q-Value}} \right)^2$$



Compute loss: $\frac{\left(Q_{new}(s,a) - (r + \gamma \max_{a} Q_{old}(s',a))\right)^2}{\theta_{new}}$

- Backward pass: $\dfrac{\partial Loss}{\partial heta_{new}}$



MSE Loss :=
$$\left(Q_{new}(s, a) - (r + \max_{a} Q_{old}(s', a))\right)^2$$

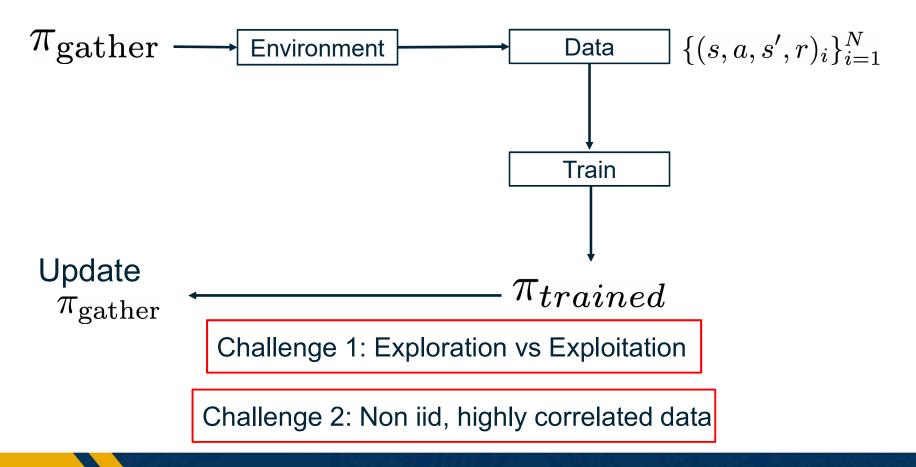
- In practice, for stability:
 - Freeze Q_{old} and update Q_{new} parameters
 - lacksquare Set $Q_{old} \leftarrow Q_{new}$ at regular intervals

How to gather experience?

$$\{(s, a, s', r)_i\}_{i=1}^N$$

This is why RL is hard





Georgia Tech

- What should π_{gather} be?
 - Greedy? -> Local minimas, no exploration $\arg\max_a Q(s,a;\theta)$
- An exploration strategy:
 - \bullet ϵ -greedy

$$a_t = \begin{cases} \arg\max_{a} Q(s, a) & \text{with probability } 1 - \epsilon \\ \text{random action} & \text{with probability } \epsilon \end{cases}$$

- Samples are correlated => high variance gradients => inefficient learning
- Current Q-network parameters determines next training samples => can lead to bad feedback loops
 - e.g. if maximizing action is to move right, training samples will be dominated by samples going right, may fall into local minima

start

R=10		R=1



- Correlated data: addressed by using experience replay
 - ightharpoonup A replay buffer stores transitions (s,a,s^{\prime},r)
 - Continually update replay buffer as game (experience) episodes are played, older samples discarded
 - Train Q-network on random minibatches of transitions from the replay memory, instead of consecutive samples
- Larger the buffer, lower the correlation

```
Algorithm 1 Deep Q-learning with Experience Replay
   Initialize replay memory \mathcal{D} to capacity N
                                                                            Experience Replay
   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
                                                                     Epsilon-greedy
            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}
                                                                                                   Q Update
           Perform a gradient descent step on (y_i - Q(\phi_i, a_i; \theta))^2 according to equation 3
       end for
  end for
```



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

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

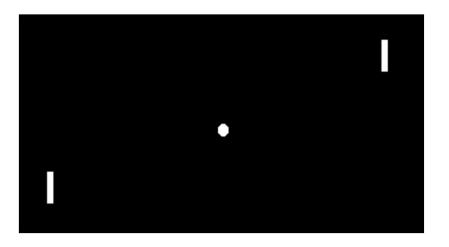
Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Case study: Playing Atari Games



Atari Games





https://www.youtube.com/watch?v=V1eYniJ0Rnk

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

Case study: Playing Atari Games



In today's class, we looked at

- Dynamic Programming
 - Value, Q-Value Iteration
- Reinforcement Learning (RL)
 - The challenges of (deep) learning based methods
 - Value-based RL algorithms
 - Deep Q-Learning

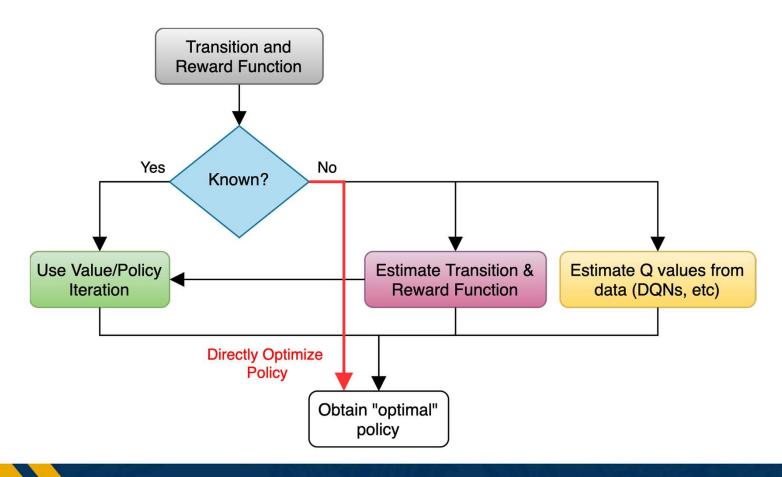
Now:

Policy-based RL algorithms (policy gradients)



Policy Gradients, Actor-Critic





Overview



ullet Class of policies defined by parameters heta

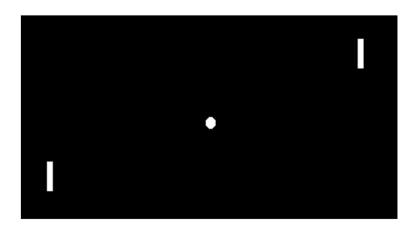
$$\pi_{\theta}(a|s): \mathcal{S} \to \mathcal{A}$$

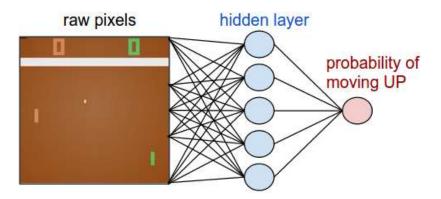
- ullet Eg: heta can be parameters of linear transformation, deep network, etc.
- Want to maximize:

$$J(\pi) = \mathbb{E}\left[\left|\sum_{t=1}^{T} \mathcal{R}(s_t, a_t)\right|\right]$$

In other words,

$$\pi^* = \arg \max_{\pi: \mathcal{S} \to \mathcal{A}} \mathbb{E} \left[\sum_{t=1}^T \mathcal{R}(s_t, a_t) \right] \longrightarrow \theta^* = \arg \max_{\theta} \mathbb{E} \left[\sum_{t=1}^T \mathcal{R}(s_t, a_t) \right]$$





Georgia Tech

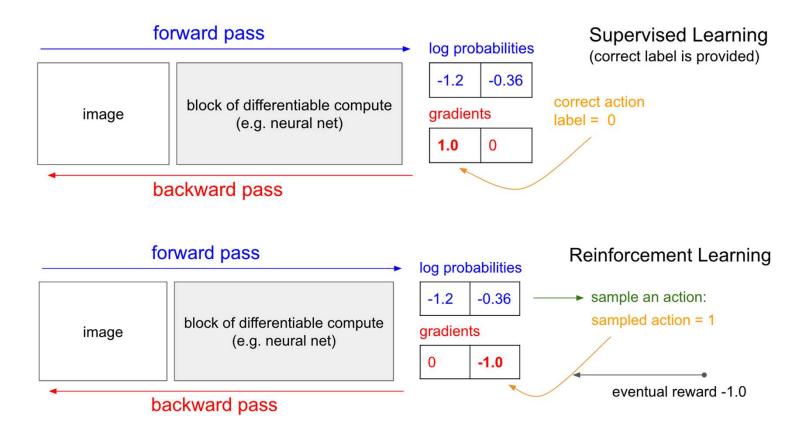


Image Source: http://karpathy.github.io/2016/05/31/rl/



Slightly re-writing the notation

Let
$$au = (s_0, a_0, \dots s_T, a_T)$$
 denote a trajectory

$$\pi_{\theta}(\tau) = p_{\theta}(\tau) = p_{\theta}(s_0, a_0, \dots s_T, a_T)$$

$$= p(s_0) \prod_{t=0}^{T} p_{\theta}(a_t \mid s_t) \cdot p(s_{t+1} \mid s_t, a_t)$$

$$\arg\max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\mathcal{R}(\tau) \right]$$

$$J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\mathcal{R}(\tau) \right]$$

$$= \mathbb{E}_{a_{t} \sim \pi(\cdot | s_{t}), s_{t+1} \sim p(\cdot | s_{t}, a_{t})} \left[\sum_{t=0}^{T} \mathcal{R}(s_{t}, a_{t}) \right]$$
How to gather data?

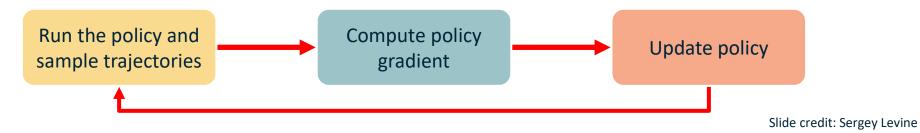
- How to gather data?
 - We already have a policy: π_{θ}
 - Sample N trajectories $\{\tau_i\}_{i=1}^N$ by acting according to π_{θ}

$$\approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} r(s_t^i, a_t^i)$$

- ullet Sample trajectories $\, au_i = \{s_1, a_1, \dots s_T, a_T\}_i$ by acting according to $\,\pi_{ heta}$
- Compute policy gradient as

$$\nabla_{\theta}J(\theta) \approx$$
 ?

• Update policy parameters: $heta \leftarrow heta + lpha
abla_{ heta} J(heta)$



The REINFORCE Algorithm



$$\begin{split} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)}[\mathcal{R}(\tau)] \\ &= \nabla_{\theta} \int \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau & \text{Expectation as integral} \\ &= \int \nabla_{\theta} \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau & \text{Exchange integral and gradient} \\ &= \int \nabla_{\theta} \pi_{\theta}(\tau) \cdot \frac{\pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} \cdot \mathcal{R}(\tau) d\tau & \\ &= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau & \nabla_{\theta} \log \pi(\tau) = \frac{\nabla_{\theta} \pi(\tau)}{\pi(\tau)} \\ &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)}[\nabla_{\theta} \log \pi_{\theta}(\tau) \mathcal{R}(\tau)] \end{split}$$

$$\pi_{\theta}(\tau) = p(s_0) \prod_{t=0}^{T} p_{\theta}(a_t \mid s_t) \cdot p(s_{t+1} \mid s_t, a_t)$$

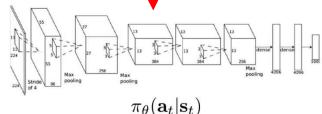
$$egin{aligned}
abla_{ heta} J(heta) &= \mathbb{E}_{ au \sim p_{ heta}(au)} [
abla_{ heta} \log \pi_{ heta}(au) \mathcal{R}(au)] \
abla_{ heta} \left[rac{\log p(s_0)}{\sum_{t=1}^T \log \pi_{ heta}(a_t|s_t)} + \sum_{t=1}^T rac{\log p(s_{t+1} + s_t, a_t)}{\sum_{t=1}^T \log p(s_{t+1} + s_t, a_t)}
ight] \end{aligned}$$

$$\nabla_{\theta} \left[\log p(s_0) + \sum_{t=1}^{T} \log \pi_{\theta}(a_t|s_t) + \sum_{t=1}^{T} \log p(s_{t+1} \mid s_t, a_t) \right]$$

Doesn't depend on Transition probabilities!

$$= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(a_{t}|s_{t}) \cdot \sum_{t=1}^{T} \mathcal{R}(s_{t}, a_{t}) \right]$$







 \mathbf{a}_t

Continuous Action Space?

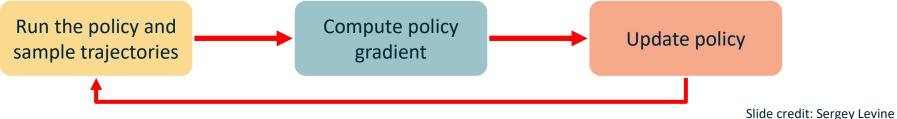
Deriving The Policy Gradient



- ullet Sample trajectories $au_i = \{s_1, a_1, \dots s_T, a_T\}_i$ by acting according to $\pi_{oldsymbol{ heta}}$
- Compute policy gradient as

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i}^{N} \left[\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta} \left(a_{t}^{i} \mid s_{t}^{i} \right) \cdot \sum_{t=1}^{T} \mathcal{R} \left(s_{t}^{i} \mid a_{t}^{i} \right) \right]$$

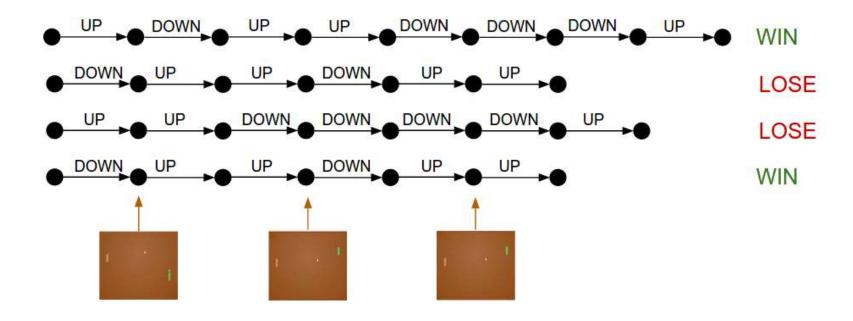
• Update policy parameters: $heta \leftarrow heta + lpha
abla_{ heta} J(heta)$



Slide Credit. Sergey Levill

The REINFORCE Algorithm





Slide credit: Dhruv Batra



Issues with Policy Gradients

- Credit assignment is hard!
 - Which specific action led to increase in reward
 - Suffers from high variance → leading to unstable training

