# Reinforcement Learning

Designing, Visualizing and Understanding Deep Neural Networks

# CS W182/282A

Instructor: Sergey Levine UC Berkeley



# From prediction to control





- i.i.d. distributed data (each datapoint is independent)
- ground truth supervision
- objective is to predict the right label

These are not **just** issues for control: in many cases, real-world deployment of ML has these same **feedback** issues **Example:** decisions made by a traffic prediction system might affect the route that people take, which changes traffic

- each decision can change future inputs (not independent)
- supervision may be high-level (e.g., a goal)
- objective is to accomplish the task

# How do we specify what we want?





high reward



### Definitions

Markov decision process

S – state space

states  $s \in \mathcal{S}$  (discrete or continuous)

 $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{T}, r\}$ 

 $\mathcal{A}$  – action space

r – reward function







transition operator or transition probability or dynamics (all synonyms)

### Definitions

Markov decision process

 $\mathcal{S}$  – state space

states  $s \in \mathcal{S}$  (discrete or continuous)

 $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{T}, r\}$ 

- $\mathcal{A}$  action space
- r reward function





 $r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ 



Andrey Markov



**Richard Bellman** 

### Definitions

partially observed Markov decision process

 $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{E}, r\}$ 

- S state space states  $s \in S$  (discrete or continuous)
- $\mathcal{A}$  action space actions  $a \in \mathcal{A}$  (discrete or continuous)
- $\mathcal{O}$  observation space observations  $o \in \mathcal{O}$  (discrete or continuous)
- $\mathcal{T}$  transition operator (like before)
- $\mathcal{E}$  emission probability  $p(o_t|s_t)$

r - reward function  $r: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ 

For now, we'll stick to regular (fully observed) MDPs, but you should know that this exists



### The goal of reinforcement learning



$$\underbrace{p_{\theta}(\mathbf{s}_{1}, \mathbf{a}_{1}, \dots, \mathbf{s}_{T}, \mathbf{a}_{T})}_{p_{\theta}(\tau)} = p(\mathbf{s}_{1}) \prod_{t=1}^{T} \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) p(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t})$$

$$\theta^{\star} = \arg\max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

### The goal of reinforcement learning





### Wall of math time (policy gradients)

### The goal of reinforcement learning



$$\underbrace{p_{\theta}(\mathbf{s}_{1}, \mathbf{a}_{1}, \dots, \mathbf{s}_{T}, \mathbf{a}_{T})}_{p_{\theta}(\tau)} = p(\mathbf{s}_{1}) \prod_{t=1}^{T} \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) p(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t})$$

$$\theta^{\star} = \arg\max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

### Evaluating the objective

$$\theta^{\star} = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$
$$J(\theta)$$



$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right] \approx \frac{1}{N} \sum_{i} \sum_{t} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

### Direct policy differentiation

$$\theta^{\star} = \arg \max_{\theta} E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$
$$J(\theta)$$

### a convenient identity

$$\pi_{\theta}(\tau)\nabla_{\theta}\log\pi_{\theta}(\tau) = \pi_{\theta}(\tau)\frac{\nabla_{\theta}\pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} = \underline{\nabla_{\theta}\pi_{\theta}(\tau)}$$

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)] = \int \pi_{\theta}(\tau)r(\tau)d\tau$$
$$\sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t})$$

$$\nabla_{\theta} J(\theta) = \int \underline{\nabla_{\theta} \pi_{\theta}(\tau)} r(\tau) d\tau = \int \underline{\pi_{\theta}(\tau)} \nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau) d\tau = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$$

### Direct policy differentiation

$$\theta^{\star} = \arg \max_{\theta} J(\theta)$$

$$J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)]$$

$$\log \text{ of both} \quad \text{sides} \quad \pi_{\theta}(\tau)$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$$

$$\nabla_{\theta} \left[ \log r(\mathbf{s}_{1}) + \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) + \log p(\mathbf{s}_{t+1} | \mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right) \right]$$



### Evaluating the policy gradient

recall: 
$$J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right] \approx \frac{1}{N} \sum_{i} \sum_{t} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right) \right]$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$
what is this?











### Comparison to maximum likelihood

policy gradient: 
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

maximum likelihood: 
$$\nabla_{\theta} J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right)$$



### What did we just do?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau_{i}) r(\tau_{i})}_{\sum_{t=1}^{T} \nabla_{\theta} \log_{\theta} \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}$$

bad stuff is made less likely

simply formalizes the notion of "trial and error"!

REINFORCE algorithm:

1. sample 
$$\{\tau^i\}$$
 from  $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$  (run it on the robot)  
2.  $\nabla_{\theta} J(\theta) \approx \sum_i \left( \sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$   
3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ 

maximum likelihood:  $\nabla_{\theta} J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} \log \pi_{\theta}(\tau_i)$ 

### Partial observability



 $\mathbf{a}_t$ 

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{o}_{i,t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

Markov property is not actually used!

Can use policy gradient in partially observed MDPs without modification

### Making policy gradient actually work

### A better estimator

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left( \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

Causality: policy at time t' cannot affect reward at time t when t < t'

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left( \sum_{\substack{i=1 \ t' \in \mathbf{I}}}^{T} r(\mathbf{s}_{i,t''}, \mathbf{a}_{i,t''}) \right)$$
  
"reward to go"  
 $\hat{Q}_{i,t}$ 

### Baselines

a convenient identity  $\pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) = \nabla_{\theta} \pi_{\theta}(\tau)$ 



$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} \log \pi_{\theta}(\tau) [n(\pi)) - b]$$
$$b = \frac{1}{N} \sum_{i=1}^{N} r(\tau) \qquad \text{but... are we allowed to do that??}$$

$$E[\nabla_{\theta} \log \pi_{\theta}(\tau)b] = \int \pi_{\theta}(\tau)\nabla_{\theta} \log \pi_{\theta}(\tau)b \, d\tau = \int \nabla_{\theta}\pi_{\theta}(\tau)b \, d\tau = b\nabla_{\theta} \int \pi_{\theta}(\tau)d\tau = b\nabla_{\theta}1 = 0$$

subtracting a baseline is *unbiased* in expectation!

average reward is *not* the best baseline, but it's pretty good!

### Policy gradient is on-policy

 $\theta^{\star} = \arg \max_{\theta} J(\theta)$ 

 $J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)]$ 

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$$
this is trouble...

- Neural networks change only a little bit with each gradient step
- On-policy learning can be extremely inefficient!

REINFORCE algorithm: 1. sample  $\{\tau^i\}$  from  $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$  (run it on the robot) 2.  $\nabla_{\theta} J(\theta) \approx \sum_i \left( \sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$ 3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ 

## Off-policy learning & importance sampling

 $\theta^{\star} = \arg \max_{\theta} J(\theta)$ 

 $J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)]$ 

what if we don't have samples from  $\pi_{\theta}(\tau)$ ? (we have samples from some  $\bar{\pi}(\tau)$  instead)

$$J(\theta) = E_{\tau \sim \bar{\pi}(\tau)} \left[ \frac{\pi_{\theta}(\tau)}{\bar{\pi}(\tau)} r(\tau) \right]$$
$$\pi_{\theta}(\tau) = p(\mathbf{s}_1) \prod_{t=1}^{T} \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$$

$$\frac{\pi_{\theta}(\tau)}{\bar{\pi}(\tau)} = \frac{p(\mathbf{s}_1) \prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)}{p(\mathbf{s}_1) \prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t) p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)} = \frac{\prod_{t=1}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=1}^T \bar{\pi}(\mathbf{a}_t | \mathbf{s}_t)}$$

importance sampling  $E_{x \sim p(x)}[f(x)] = \int p(x)f(x)dx$   $= \int \frac{q(x)}{q(x)}p(x)f(x)dx$   $= \int q(x)\frac{p(x)}{q(x)}f(x)dx$   $= E_{x \sim q(x)}\left[\frac{p(x)}{q(x)}f(x)\right]$ 

### Deriving the policy gradient with IS

 $\theta^{\star} = \arg\max_{\theta} J(\theta)$ 

 $J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)]$ 

can we estimate the value of some *new* parameters  $\theta'$ ?

 $J(\theta') = E_{\tau \sim \pi_{\theta}(\tau)} \begin{bmatrix} \pi_{\theta'}(\tau) \\ \pi_{\theta}(\tau) \end{bmatrix}$  the only bit that depends on  $\theta'$ 

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \frac{\nabla_{\theta'} \pi_{\theta'}(\tau)}{\pi_{\theta}(\tau)} r(\tau) \right] = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \frac{\pi_{\theta'}(\tau)}{\pi_{\theta}(\tau)} \nabla_{\theta'} \log \pi_{\theta'}(\tau) r(\tau) \right]$$

now estimate locally, at  $\theta = \theta'$ :  $\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) r(\tau)]$ 

a convenient identity

$$\pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) = \nabla_{\theta} \pi_{\theta}(\tau)$$

## The off-policy policy gradient

 $\theta^{\star} = \arg \max_{\theta} J(\theta)$ 

 $J(\theta) = E_{\tau \sim \pi_{\theta}(\tau)}[r(\tau)]$ 

$$\frac{\pi_{\theta'}(\tau)}{\pi_{\theta}(\tau)} = \frac{\prod_{t=1}^{T} \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\prod_{t=1}^{T} \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}$$

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \frac{\pi_{\theta'}(\tau)}{\pi_{\theta}(\tau)} \nabla_{\theta'} \log \pi_{\theta'}(\tau) r(\tau) \right] \quad \text{when } \theta \neq \theta'$$

$$= E_{\tau \sim \pi_{\theta}(\tau)} \left[ \left( \prod_{t=1}^{T} \frac{\pi_{\theta'}(\mathbf{a}_{t} | \mathbf{s}_{t})}{\pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t})} \right) \left( \sum_{t=1}^{T} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{t} | \mathbf{s}_{t}) \right) \left( \sum_{t=1}^{T} r(\mathbf{s}_{t}, \mathbf{a}_{t}) \right) \right] \text{ what about causality?}$$

$$= E_{\tau \sim \pi_{\theta}(\tau)} \left[ \sum_{t=1}^{T} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{t} | \mathbf{s}_{t}) \left( \prod_{t'=1}^{t} \frac{\pi_{\theta'}(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{\pi_{\theta}(\mathbf{a}_{t'} | \mathbf{s}_{t'})} \right) \left( \sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \left( \prod_{t''=t}^{t'} \frac{\pi_{\theta'}(\mathbf{a}_{t''} | \mathbf{s}_{t''})}{\pi_{\theta}(\mathbf{a}_{t''} | \mathbf{s}_{t''})} \right) \right) \right]$$

$$\text{future actions don't affect current weight} \quad \textbf{future actions don't affect current weight} \quad \textbf{future actions algorithm (more on this in a later lecture)}$$

### A first-order approximation for IS (preview)

$$\nabla_{\theta'} J(\theta') = E_{\tau \sim \pi_{\theta}(\tau)} \left[ \sum_{t=1}^{T} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t) \left( \prod_{t'=1}^{t} \frac{\pi_{\theta'}(\mathbf{a}_{t'} | \mathbf{s}_{t'})}{\pi_{\theta}(\mathbf{a}_{t'} | \mathbf{s}_{t'})} \right) \left( \sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right) \right]$$
exponential in T...

let's write the objective a bit differently...

on-policy policy gradient: 
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}$$
  
off-policy policy gradient:  $\nabla_{\theta'} J(\theta') \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \frac{\pi_{\theta'}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})}{\pi_{\theta}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}$   
This simplification is very commonly used in practice!  $= \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \frac{\pi_{\theta'}(\mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{s}_{i,t})} \frac{\pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})} \nabla_{\theta'} \log \pi_{\theta'}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}$ 

### Policy gradient with automatic differentiation

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \frac{\nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})}{\mathbf{Q}_{i,t}}$$
pretty inefficient to compute these explicitly

How can we compute policy gradients with automatic differentiation?

We need a graph such that its gradient is the policy gradient!

maximum likelihood: 
$$\nabla_{\theta} J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \qquad J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$$

Just implement "pseudo-loss" as a weighted maximum likelihood:

$$\tilde{J}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}$$
cross entropy (discrete) or squared error (Gaussian)

### Policy gradient with automatic differentiation

Pseudocode example (with discrete actions):

Maximum likelihood:

```
# Given:
# actions - (N*T) x Da tensor of actions
# states - (N*T) x Ds tensor of states
# Build the graph:
logits = policy.predictions(states) # This should return (N*T) x Da tensor of action logits
negative_likelihoods = tf.nn.softmax_cross_entropy_with_logits(labels=actions, logits=logits)
loss = tf.reduce_mean(negative_likelihoods)
gradients = loss.gradients(loss, variables)
```

## Policy gradient with automatic differentiation

Pseudocode example (with discrete actions):

Policy gradient:

```
# Given:
# actions - (N*T) x Da tensor of actions
# states - (N*T) x Ds tensor of states
# q_values - (N*T) x 1 tensor of estimated state-action values
# Build the graph:
logits = policy.predictions(states) # This should return (N*T) x Da tensor of action logits
negative_likelihoods = tf.nn.softmax_cross_entropy_with_logits(labels=actions, logits=logits)
weighted_negative_likelihoods = tf.multiply(negative_likelihoods, q_values)
loss = tf.reduce_mean(weighted_negative_likelihoods)
gradients = loss.gradients(loss, variables)
```

$$\tilde{J}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t} | \hat{Q}_{i,t}) \mathbf{q_values}$$

### Policy gradient in practice

- Remember that the gradient has high variance
  - This isn't the same as supervised learning!
  - Gradients will be really noisy!
- Consider using much larger batches
- Tweaking learning rates is very hard
  - Adaptive step size rules like ADAM can be OK-ish
  - We'll learn about policy gradient-specific learning rate adjustment methods later!

### What are policy gradients **used** for?

#### Iteration 0







learning for control

"hard" attention: network selects where to look (Mnih et al. Recurrent Models of Visual Attention)

Discrete latent variable models (we'll learn about these later!)

**Long story short:** policy gradient (REINFORCE) can be used in any setting where we have to **differentiate** through a stochastic but non-differentiable operation

### Review

- Evaluating the RL objective
  - Generate samples
- Evaluating the policy gradient
  - Log-gradient trick
  - Generate samples
- Policy gradient is on-policy
- Can derive off-policy variant
  - Use importance sampling
  - Exponential scaling in T
  - Can ignore state portion (approximation)
- Can implement with automatic differentiation – need to know what to backpropagate
- Practical considerations: batch size, learning rates, optimizers

### REINFORCE algorithm:

1. sample  $\{\tau^i\}$  from  $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$  (run it on the robot) 2.  $\nabla_{\theta} J(\theta) \approx \sum_i \left( \sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left( \sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$ 3.  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ 

## Example: trust region policy optimization

- Natural gradient with automatic step adjustment
- Discrete and continuous actions



## Policy gradients suggested readings

• Classic papers

- Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning: introduces REINFORCE algorithm
- Baxter & Bartlett (2001). Infinite-horizon policy-gradient estimation: temporally decomposed policy gradient (not the first paper on this! see actor-critic section later)
- Peters & Schaal (2008). Reinforcement learning of motor skills with policy gradients: very accessible overview of optimal baselines and natural gradient
- Deep reinforcement learning policy gradient papers
  - Levine & Koltun (2013). Guided policy search: deep RL with importance sampled policy gradient (unrelated to later discussion of guided policy search)
  - Schulman, L., Moritz, Jordan, Abbeel (2015). Trust region policy optimization: deep RL with natural policy gradient and adaptive step size
  - Schulman, Wolski, Dhariwal, Radford, Klimov (2017). Proximal policy optimization algorithms: deep RL with importance sampled policy gradient