Policy Gradient Theorem, REINFORCE Algorithm

Reading

- Reinforcement Learning
 - http://www.incompleteideas.net/book/the-book-2nd.html
 - Chapters 13.1-13.3
- David Silver lecture on Policy Gradients
 - https://www.youtube.com/watch?v=KHZVXao4qXs&t=3s

Overview

- In Q-learning, we select actions based on their q values
- With policy gradient methods, the controller is just a function that has no notion of q values
 - Of course, during training, it will be trained to select actions that maximize q values
- Policy gradient methods are more flexible than standard Qlearning for a number of reasons
 - Can handle partially observable MDPs
 - Can handle continuous control systems
- At the same time, training with policy gradient methods is very unstable

Policy Gradients vs Q-learning

- In Q-learning, we have one function (approximation)
 - We have an estimate q(s, a) for each s, a pair
 - Those estimates define a deterministic policy
- In policy gradient methods, we have one function to estimate action values and a separate function for the policy
 - We update the two separately (using gradients)
 - Even if the Q approximations are (temporarily) wrong, the policy may not be affected much since it's slowly updated according to its learning rate

Finite-MDP Setup

- The finite MDP setup is the same as before
 - An MDP is the usual 5-tuple (S, A, P, R, η)
- The main difference is that now the policy does not depend on the q-values:
 - -recall that $\pi(a|s; \theta)$ is the probability that action a is taken from state s
 - —the parameters $oldsymbol{ heta}$ are determined during training
 - —looks the same as before except there is no explicit computation of q-values

Policy Example

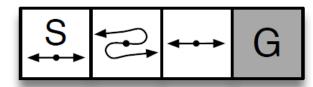
- Can encode a probabilistic policy, with parameters $oldsymbol{ heta}$, using softmax
 - -How?
 - Let the current state be s
 - Let x(s, a) be a feature vector of all states and actions, e.g., one-hot encoding
 - Then the probability of taking action *a* from state *s* is:

$$\pi(a|s;\boldsymbol{\theta}) = \frac{e^{\boldsymbol{\theta}^T x(s,a)}}{\sum_{a'} e^{\boldsymbol{\theta}^T x(s,a')}}$$

- The encoding x(s, a) can be any encoding, including non-linear functions of the states and actions
- Of course, π may also be arbitrarily complex (wink, wink)

Partial Observability

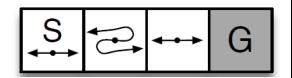
Consider this short corridor example



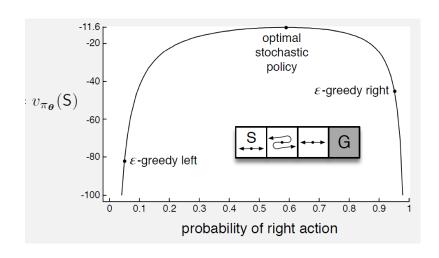
- Actions are left/right, but their effect is reversed in state 2
- -Suppose features are $x(s, right) = [1\ 0], x(s, left) = [0\ 1]$
 - Same features regardless of the value of s
 - You don't see which state you're in effectively
- Reward of -1 after each step
- What is the optimal policy (without knowing where you are)?
 - Need to make two rights and a left
 - So cannot be deterministic
 - Turns out a coin flip with a slight bias to the right is optimal
 - More next

Partial Observability, cont'd

What would an action-value method do?



- $-\operatorname{If} Q([1\ 0]) > Q([0\ 1])$, always go right
 - Or with ϵ -greedy probability
- Cannot learn different policies per state
 - At best, take correct action w.p. ϵ
- Policy gradient method will learn a better probability than ϵ greedy



Policy Gradient Setup

- Unlike Q-learning, we now have a policy that is learned separately from the q-values
- The policy π is defined in the same way as before:

$$\pi(a|s; \boldsymbol{\theta}) = \mathbb{P}_{\pi}[A_t = a|S_t = s]$$

- The state values, $v_{\pi}(s)$, and action values, $q_{\pi}(s,a)$, are defined in the same way as before
 - The main difference is that the policy is now trained separately from the value estimates
 - We are now directly training the policy to maximize the value of each state

Optimization Function

- What function should the policy optimize?
 - Maximize the value $v_{\pi}(s)$ for all s
 - What are some issues with this?
 - Don't know the real v_{π}
 - Also, v_{π} is policy-specific, so it changes every time we change π
 - If we knew $v_{\pi}(s)$ for each state s, how would we train π ?
- To maximize v_{π} , need to calculate its gradient w.r.t. $oldsymbol{ heta}$
 - Here, the policy π is parameterized by $m{ heta}$ (written $\pi_{m{ heta}}$)
 - In this case, the gradient is called a policy gradient
 - The policy can be any function, as usual
 - E.g., a neural network's parameters
 - When clear from context, we'll just write π
 - Only requirement is that it's differentiable w.r.t $oldsymbol{ heta}$

Policy Gradients

• Look at policy gradient (w.r.t. θ) in finite state case:

$$\nabla v_{\pi}(s) = \nabla \left[\sum_{a} \pi(a|s) q_{\pi}(s, a) \right]$$

$$= \sum_{a} \nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \nabla q_{\pi}(a, s)$$

$$= \sum_{a} \nabla \pi(a|s) q_{\pi}(s, a) +$$

$$+ \pi(a|s) \nabla \sum_{s'} P(s, a, s') [R(s, a, s') + \gamma v_{\pi}(s')]$$

$$= \sum_{a} \nabla \pi(a|s) q_{\pi}(s, a) + \gamma \pi(a|s) \sum_{s'} P(s, a, s') \nabla v_{\pi}(s')$$

• Notice that ∇v_{π} appears recursively

$$\nabla v_{\pi}(s) = \sum_{a} \nabla \pi(a|s) q_{\pi}(s,a) + \gamma \pi(a|s) \sum_{s'} P(s,a,s') \\ \left[\sum_{a'} \nabla \pi(a'|s') q_{\pi}(s',a') + \gamma \pi(a'|s') \sum_{s''} P(s',a',s'') \nabla v_{\pi}(s'') \right]$$

- Notation:
- $\mathbb{P}[s \to x, k, \pi]$ is the probability that state x is visited from state s after k steps (following policy π)
 - $\mathbb{P}[s \to x, 0, \pi] = 1$ if s = x and 0, otherwise
 - $\mathbb{P}[s \to x, 1, \pi] = \sum_a \pi(a|s) P(s, a, x)$
 - $\mathbb{P}[s \to x, 2, \pi] = \sum_{a} \pi(a|s) \sum_{s'} P(s, a, s') \sum_{a'} \pi(a'|s') P(s', a', x)$

$$\nabla v_{\pi}(s) = \sum_{a} \nabla \pi(a|s) q_{\pi}(s,a) + \gamma \pi(a|s) \sum_{s'} P(s,a,s')$$

$$\left[\sum_{a'} \nabla \pi(a'|s') q_{\pi}(s',a') + \gamma \pi(a'|s') \sum_{s''} P(s',a',s'') \nabla v_{\pi}(s'') \right]$$

Look at first term:

$$\sum_{a} \nabla \pi(a|s) q_{\pi}(s,a) =$$

$$= \sum_{x \in S} \mathbb{P}[s \to x, 0, \pi] \sum_{a} \nabla \pi(a|x) q_{\pi}(x,a)$$

- since $\mathbb{P}[s \to x, 0, \pi] = 1$ only when x = s

$$\nabla v_{\pi}(s) = \sum_{a} \nabla \pi(a|s) q_{\pi}(s,a) + \gamma \pi(a|s) \sum_{s'} P(s,a,s')$$

$$\left[\sum_{a'} \nabla \pi(a'|s') q_{\pi}(s',a') + \gamma \pi(a'|s') \sum_{s''} P(s',a',s'') \nabla v_{\pi}(s'') \right]$$

Look at second term (rename s' to x):

$$\gamma \sum_{a} \pi(a|s) \sum_{x} P(s, a, x) \left[\sum_{a'} \nabla \pi(a'|x) q_{\pi}(x, a') \right] =$$

$$= \gamma \sum_{x} \left[\sum_{a'} \nabla \pi(a'|x) q_{\pi}(x, a') \right] \sum_{a} \pi(a|s) P(s, a, x)$$

$$= \sum_{x \in S} \mathbb{P}[s \to x, 1, \pi] \left[\sum_{a'} \nabla \pi(a'|x) q_{\pi}(x, a') \right]$$

Rewriting the policy gradient:

$$\nabla v_{\pi}(s) = \sum_{x \in S} \mathbb{P}[s \to x, 0, \pi] \sum_{a} \nabla \pi(a|x) q_{\pi}(x, a) +$$

$$+ \gamma \sum_{x} \mathbb{P}[s \to x, 1, \pi] \sum_{a} \nabla \pi(a|x) q_{\pi}(x, a)$$

$$+ \sum_{x} \gamma \pi(a|s) \sum_{s'} P(s, a, s') \left[\sum_{a'} \gamma \pi(a'|s') \sum_{s''} P(s', a', s'') \nabla v_{\pi}(s'') \right]$$

 We can continue the expansion in the same fashion for future steps

$$\nabla v_{\pi}(s_0) = \sum_{s \in \mathcal{S}} \sum_{k=0}^{\infty} \gamma^k \mathbb{P}[s_0 \to s, k, \pi] \sum_{a} \nabla \pi(a|s) q_{\pi}(s, a)$$

- We can treat the sum of probabilities as the discounted aggregate state visitation "probability"
 - -Call it d_{π}
 - Similar to the stationary distribution μ_{π} but not the same $\mu_{\pi}P=\mu_{\pi}$
- What probability does $oldsymbol{\mu}_{\pi}$ capture?

$$\lim_{k\to\infty} \mathbb{P}[s_0\to s,k,\pi]$$

• If you want to treat d_{π} as a real probability distribution, need to normalize it so that it sums up to 1

$$\nabla v_{\pi}(s_0) = \sum_{s \in \mathcal{S}} \sum_{k=0}^{\infty} \gamma^k \mathbb{P}[s_0 \to s, k, \pi] \sum_{a} \nabla \pi(a|s) q_{\pi}(s, a)$$

- We can treat the sum of probabilities as the discounted aggregate state visitation probability
 - Call it d_{π}
- So, finally

$$\nabla v_{\pi}(s_0) = \sum_{s} d_{\pi}(s) \sum_{a} \nabla \pi(a|s) q_{\pi}(s,a)$$

• This is the policy gradient theorem!

Using the policy gradient theorem

- To improve a given a policy π_{θ} , we observe the next state-action-reward pair, and compute the gradient
- Note that we can think of the gradient as an expectation

$$\nabla v_{\pi}(s_0) = \sum_{s} d_{\pi}(s) \sum_{a} \nabla \pi(a|s) q_{\pi}(s,a)$$
$$= \mathbb{E}_{d_{\pi}} \left[\sum_{a} \nabla \pi(a|S_t) q_{\pi}(S_t,a) \right]$$

- Technically need to normalize d_π
 - That's just a constant which would be multiplied by the learning rate anyway
- How do we approximate the expectation using real data?
 - Average over real data

Using the policy gradient theorem, cont'd

Note that we can think of the gradient as an expectation

$$\nabla v_{\pi}(s_0) = \mathbb{E}_{d_{\pi}} \left| \sum_{a} \nabla \pi(a|S_t) q_{\pi}(S_t, a) \right|$$

- How do we approximate the expectation using real data?
 - Average over real data
 - For each state s, compute gradient over all actions:

$$\sum_{a} \nabla \pi(a|s) q_{\pi}(s,a)$$

- Any issues with this?
- Need to know all $q_{\pi}(s, a)$

REINFORCE algorithm

- The benefit of the policy gradient theorem is that we can compute gradients w.r.t. θ and improve the policy
 - As long as we have good estimates \hat{q} of the real q function
 - We'll discuss several ways to get \hat{q}
- We could directly instantiate a gradient-descent algorithm:

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha \sum_{a} \hat{q}(S_t, a) \nabla \pi(a|S_t; \boldsymbol{\theta})$$

- Of course, this requires good estimates \hat{q} for all state-action pairs
 - -Also a θ update updates the entire function, not just the current action parameters
 - Could be very unstable

REINFORCE algorithm, cont'd

- To avoid needing an estimate for each action value, one could modify the policy gradient theorem
 - Could use the return G_t directly
 - To simplify the math, focus on finite-horizon case

$$\begin{array}{ll} \nabla v_{\pi}(s) = \nabla \mathbb{E}_{\pi}[G_{1}|S_{1}=s] \\ \text{(law of total probability)} & = \displaystyle \sum_{tr=(S_{1},A_{1},R_{1}\dots)} \nabla \mathbb{P}_{\pi}[tr|S_{1}=s]\mathbb{E}_{\pi}[G_{1}|tr] \\ & = \displaystyle \sum_{tr=(S_{1},A_{1},R_{1}\dots)} \mathbb{P}_{\pi}[tr|S_{1}=s] \nabla \log (\mathbb{P}_{\pi}[tr|S_{1}=s])\mathbb{E}_{\pi}[G_{1}|tr] \\ & = \displaystyle \sum_{tr=(S_{1},A_{1},R_{1}\dots)} \mathbb{P}_{\pi}[tr|S_{1}=s] \nabla \log \left(\displaystyle \prod_{t=1}^{T} \pi(A_{t}|S_{t})P(S_{t},A_{t},S_{t+1}) \right) \mathbb{E}_{\pi}[G_{1}|tr] \\ & = \displaystyle \sum_{tr} \mathbb{P}_{\pi}[tr|S_{1}=s] \left(\displaystyle \sum_{t} \nabla \log (\pi(A_{t}|S_{t})) + \nabla \log (P(S_{t},A_{t},S_{t+1})) \right) \mathbb{E}_{\pi}[G_{1}|tr] \\ \end{array}$$

 $= \mathbb{E}_{\pi} \left[\sum_{t=1}^{T} \nabla \log \left(\pi(A_t | S_t) \right) G_1 \middle| S_1 = s \right]$

REINFORCE algorithm, cont'd

Final form for the gradient is

$$\nabla v_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=k}^{T} \nabla \log \left(\pi(A_{t}|S_{t}) \right) G_{k} \middle| S_{k} = s \right]$$

- Book proves a slightly different result
 - Expected value over next action only
- Once we have the gradient, update weights as usual

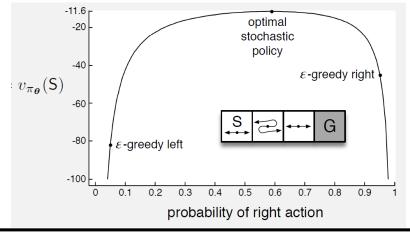
$$\boldsymbol{\theta}' = \boldsymbol{\theta} + \alpha \nabla_{\boldsymbol{\theta}} v_{\pi_{\boldsymbol{\theta}}}(s)$$

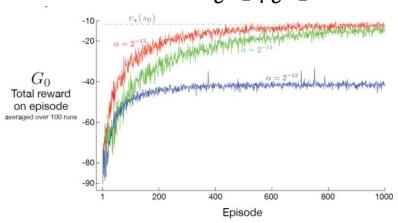
-Similar to the Monte Carlo methods where we wait until the end of the episode to observe G_t

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for π_* Input: a differentiable policy parameterization $\pi(a|s,\theta)$ Algorithm parameter: step size $\alpha > 0$ Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ (e.g., to $\mathbf{0}$) Loop forever (for each episode): Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot,\theta)$ Loop for each step of the episode $t = 0, 1, \ldots, T - 1$: $G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k \qquad (G_t)$ $\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t,\theta)$

Partial Observability, cont'd

- What would an action-value method do?
- $S \longleftrightarrow G$
- $\text{If } Q([1\ 0]) > Q([0\ 1]), \text{ always go right}$
 - Or with ϵ -greedy probability
- Cannot learn different policies per state
 - At best, take correct action w.p. ϵ
- Policy gradient method learns a better policy than ϵ -greedy
 - Suppose we use softmax policy $\pi([1\ 0]) = \frac{e^{\theta_1 \cdot 1 + \theta_2 \cdot 0}}{e^{\theta_1 + e^{\theta_2}}}$





Issues with REINFORCE

Can you spot any issues with this iteration?

$$\nabla v_{\pi}(s) = \mathbb{E}_{\pi} \left[\sum_{t=k}^{T} \nabla \log \left(\pi(A_{t}|S_{t}) \right) G_{k} \middle| S_{k} = s \right]$$

- How important is the magnitude of G_k ?
- Turns out quite a bit tasks have greatly varying returns
- Especially problematic if *good* runs have zero returns
 - Gradient is 0!
- Vanilla REINFORCE has very large variance depending on G_k
- Next time we'll discuss how to address this issue