# Lecture 12: Supervision via Bellman

In this lecture we consider an alternative method for <mark>supervising</mark> (i.e. finding target labels for) Q functions. First we start with a fundamental Lemma.

## 1) Performance-Difference Lemma

Goal: understand $V^{\pi}$ vs. $V^{\pi'}$ in terms of the difference between $\pi$ vs. $\pi'$.

Lemma (Performance Difference):

$$V^{\pi}(s_0) - V^{\pi'}(s_0) = \frac{1}{1-\gamma} \underset{s \sim d_{s_0}^{\pi}}{\mathbb{E}} \left[ \underset{a \sim \pi(s)}{\mathbb{E}} \overbrace{\left[ Q^{\pi'}(s,a) \right]}^{A^{\pi}(s,a)} - V^{\pi'}(s) \right]$$

For $r(s,a) \in [0,1]$,

$$\left| V^{\pi}(s_0) - V^{\pi'}(s) \right| \leq \frac{1}{(1-\gamma)^2} \underset{s \sim d_{s_0}^{\pi}}{\mathbb{E}} \left[ \sum_{a \in A} \underbrace{|\pi(a|s) - \pi'(a|s)|}_{\|\pi(\cdot|s) - \pi'(\cdot|s)\|_1} \right]$$

The first expression inspires us to define

<u>Def</u> (Advantage) $A^{\pi}(s,a) = Q^{\pi}(s,a) - V^{\pi}(s)$

The "advantage" of taking action $a$ at state $s$ rather than following $\pi$.

Notice that $A^{\pi}(s, \pi(s)) = 0$.

Also notice $\underset{a}{\text{argmax}}\, A^{\pi}(s,a) = \underset{a}{\text{argmax}}\, Q^{\pi}(s,a)$

## Proof of PDL:

$$V^{\pi}(S_0) - V^{\pi'}(S_0) = V^{\pi}(S_0) - \underset{a_0 \sim \pi(S_0)}{\mathbb{E}}\left[r(S_0, a_0) + \gamma \underset{S_1 \sim P(S_0, a_0)}{\mathbb{E}}[V^{\pi'}(S_1)]\right] + \underset{a_0 \sim \pi(S_0)}{\mathbb{E}}\left[r(S_0, a_0) + \gamma \underset{S_1 \sim P(S_0, a_0)}{\mathbb{E}}[V^{\pi'}(S_1)]\right] - V^{\pi'}(S_0)$$

$$= \gamma \underset{\substack{a_0 \sim \pi(S_0) \\ S_1 \sim P(S_0, a_0)}}{\mathbb{E}}\left[V^{\pi}(S_1) - V^{\pi'}(S_1)\right] + \underset{a_0 \sim \pi(S_0)}{\mathbb{E}}\left[Q(a_0, S_0) - V^{\pi'}(S_0)\right]$$

The first statement in the Lemma follows by _iteration_
(similar to Simulation Lemma)

$$\underset{a \sim \pi(S)}{\mathbb{E}}\left[Q^{\pi'}(S, a) - V^{\pi'}(S)\right] = \underset{a \sim \pi(S)}{\mathbb{E}}\left[Q^{\pi'}(S, a)\right] - \underset{a \sim \pi'(S)}{\mathbb{E}}\left[Q^{\pi'}(S, a)\right]$$

$$= \sum_{a \in \mathcal{A}} (\pi(a|S) - \pi'(a|S)) Q^{\pi'}(S, a)$$

Therefore,

$$|V^{\pi}(S_0) - V^{\pi'}(S_0)| \leq \frac{1}{1-\gamma} \underset{S \sim d_{S_0}^{\pi}}{\mathbb{E}}\left[\sum_{a \in \mathcal{A}} |\pi(a|S) - \pi'(a|S)| Q^{\pi'}(S, a)\right]$$

The second statement follows by noting $0 \leq Q^{\pi'}(S, a) \leq \frac{1}{1-\gamma}$ $\square$

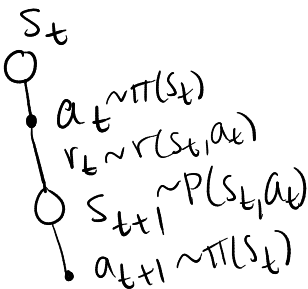We can use the PDL to prove monotonic improvement of policy iteration (HW 2).

$$V^{\pi^{t+1}}_{(S_0)} - V^{\pi^t}_{(S_0)} = \frac{1}{1-\gamma} \underset{S \sim d_{S_0}^{\pi^{t+1}}}{\mathbb{E}}\left[A^{\pi^t}(S, \pi^{t+1}(S))\right]$$

# 2) Supervision via Bellman Equation

Recall the Bellman Expectation Equation:

$$Q^{\pi}(s,a) = r(s,a) + \gamma \underset{s' \sim P(s,a)}{\mathbb{E}}\left[V^{\pi}(s')\right]$$

$$= r(s,a) + \gamma \underset{s' \sim P(s,a)}{\mathbb{E}}\left[Q^{\pi}(s',a')\right]$$

$$\underset{\text{stochastic}}{\text{possibly}} \longrightarrow a' \sim \pi(s')$$

IDEA: We can bootstrap a label for ==supervision== with just one time step!

$s_t$
$a_t \sim \pi(s_t)$
$r_t \sim r(s_t, a_t)$
$s_{t+1} \sim P(s_t, a_t)$
$a_{t+1} \sim \pi(s_t)$

At time $t$, we are at $s_t$ and sample & take action $a_t \sim \pi(s_t)$. As a result we observe $r_t$ and $s_{t+1}$. Then we sample $a_{t+1} \sim \pi(s_{t+1})$.

Then our target/label is defined as:

$$==y_t = r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1})== \qquad (s_t, a_t, y_t)$$

$$y_t \approx Q^{\pi}(s_t, a_t)$$

This is sometimes called "Temporal Difference" target

The TD error is

$$r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1}) - \hat{Q}(s_t, a_t)$$

In the tabular setting, a basic Algorithm:

## Alg: SARSA subroutine ("state-action-reward-state-action")

initialize $Q^0$, $s_0 \sim \mu_0$, $a_0 \sim \pi(s_0)$

for $t = 0, 1, \dots$

$\quad$ Take action $a_t$, observe $s_{t+1} \sim P(s_t, a_t)$ & $r_t \sim r(s_t, a_t)$

$\quad$ Sample $a_{t+1} \sim \pi(s_{t+1})$

$\quad$ update $Q^{t+1}(s_t, a_t) = (1-\alpha) Q^t(s_t, a_t) + \alpha\left(r_t + \gamma Q^t(s_{t+1}, a_{t+1})\right)$

This subroutine can be incorporated into an approximate dynamic programming algorithm (ie as the sample based policy evaluation step)

## Policy Improvement w/ ε-greedy

SARSA requires sufficient exploration to converge (for how a formal statement & proof are out of scope)
A common strategy is ε-greedy:

$$\pi(s) = \begin{cases} \operatorname*{argmax}_a Q(s,a) & \text{w.p. } 1-\varepsilon \\ a_0 & \text{w.p. } \varepsilon/A \\ a_1 & \text{w.p. } \varepsilon/A \\ \vdots \end{cases}$$

or using slightly different notation:

$$\pi(a|s) = \begin{cases} 1 - \varepsilon + \frac{\varepsilon}{A} & a = \operatorname{argmax} Q(s,a) \\ \frac{\varepsilon}{A} & \text{o.w.} \end{cases}$$
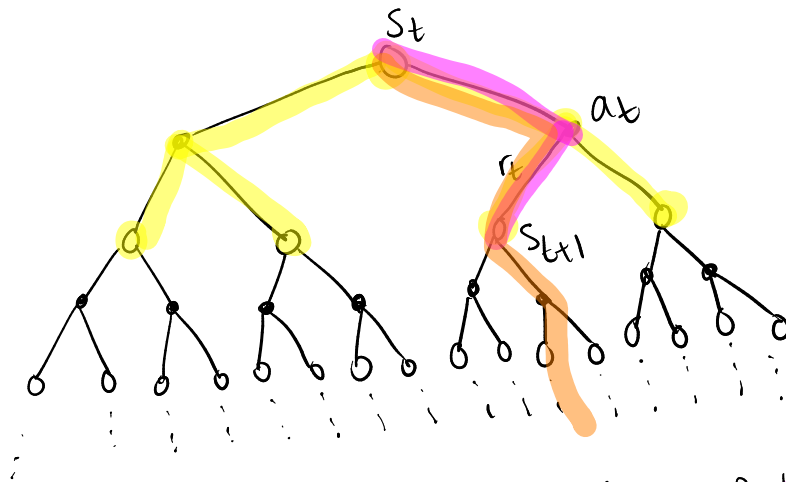
# Comparison with Rollout-based Supervision (MC):

1) TD can update $Q$ function online at every step, MC must wait untill end of rollout

2) TD is __biased__ when $\hat{Q} \neq Q^\pi$

$r_t + Q^\pi(s_{t+1}, a_{t+1})$ is unbiased, but we don't know $Q^\pi$!

MC is __unbiased__

3) Variance of TD estimate due to one stochastic transition:
$$a_t \sim \pi(s_t)$$
$$s_{t+1} \sim P(s_t, a_t)$$
$$a_{t+1} \sim \pi(s_{t+1})$$

Variance of MC due to __many__ transitions Therefore higher.

4) BOTH methods supervise $Q^\pi$ using data collected from rollouts with $\pi$, ie. they are both __on policy__



Dynamic Programing Bellman Expectation:
1 time step, all possible outcomes

Rollout-based (MC):
Many timestep, sampled outcome

Bellman-based (TD):
One timestep, sampled outcome

# 3) Supervision with Bellman Optimality

So far, we focus on estimating $Q^\pi$. But we ultimately only care about $Q^*$. Can we focus on this directly?

## recall: Bellman optimality:

$$Q^*(s,a) = r(s,a) + \gamma \underset{s' \sim P(s,a)}{\mathbb{E}}\left[\max_{a'} Q^*(s',a')\right]$$
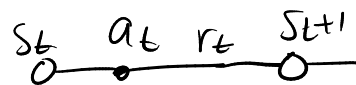
## recall: Value Iteration:

An algorithm for finding an optimal policy that focused on $Q^*$ directly

Init. $Q^0$
for $t=0,1,\dots$
$\qquad Q^{t+1} = \text{Bellman Op}(Q^t)$

BellmanOperator($Q$):
$$Q^{t+1}(s,a) = r(s,a) + \gamma \underset{s' \sim P(s,a)}{\mathbb{E}}\left[\max_{a'} Q^t(s',a')\right]$$

$$\underset{\circ}{s_t} \quad \underset{\bullet}{a_t} \quad r_t \quad \underset{\circ}{s_{t+1}}$$

## Sample-based supervision:

$$y_t = r_t + \gamma \max_a \hat{Q}(s_{t+1}, a) \qquad (s_t, a_t, y_t)$$

$$y_t \approx Q^*(s_t, a_t)$$

## Alg: Q-learning in the tabular setting

initialize $Q$
for $t=0,1,\dots$
$\qquad$ take action $a_t$ (e.g. $\varepsilon$ greedy) & observe $s_{t+1} \sim P(s_t, a_t), \ r_t \sim r(s_t, a_t)$
$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha\left(r_t + \gamma \max_{a'} Q(s_{t+1}, a')\right)$$

# Some properties of Bellman optimality based supervision

1) updates at every timestep
2) biased label when $Q \neq Q^*$
3) variance depends on randomness from one timestep
4) Not specific to a policy, so can use <u>off policy</u> data.

## 4) Function approximation

Bellman-based supervision (like rollout based) gives us labels that we can use to train models: $\{(s_i, a_i, y_i)\}_{i=1}^{N}$

$$\underline{ERM}: \quad \min_{Q \in \mathcal{Q}} \sum_{i=1}^{N} (Q(s_i, a_i) - y_i)^2$$

suppose parametrized model class
$$\mathcal{Q} = \{Q_\theta \mid \theta \in \mathbb{R}^d\}$$

Bellman-based supervision is online & incremental. So rather than full ERM minimization, it is common to do gradient descent updates to $\theta$ using incoming data.

$$\nabla_\theta (Q_\theta(s_i, a_i) - y_i)^2 = 2(Q_\theta(s_i, a_i) - y_i) \nabla_\theta Q_\theta(s_i, a_i)$$

update looks like
$$\theta \leftarrow \theta + \alpha (Q_\theta(s_i, a_i) - y_i) \nabla_\theta Q_\theta(s_i, a_i))$$

could be Bellman-exp (SARSA) or Bellman-opt (Q-learning)