# Deriving the Policy Gradient

A policy $\pi_\theta(a|x)$ specifies (as a function of parameters $\theta$) the probability of selecting action $a$ when in state $x$. For the moment, we'll be taking $x$ to be a discrete state and actions $a$ discrete actions, but a lot of this is going to generalize to continuous spaces as well.

Let $p_0(x)$ be the probability of *starting* in state $x$ at time 0.

We have that the *value* of policy $\pi_\theta$, or expected (infinite, discounted) reward associated with this policy is given by

$$V(\theta) = \sum_x p_0(x) J_\theta(x). \tag{1}$$

where

$$J_\theta(x) = \sum_a \pi_\theta(a|x) \left[ r_{x,a} + \beta \sum_y p_{x,y}^a J_\theta(y) \right]. \tag{2}$$

We'd like to find $\theta^*$ that maximizes the value, i.e., the policy that produces the largest possible expected value in application. In order to do this, we want to effectively do a sort of gradient ascent -

$$\theta_{\text{new}} \leftarrow \theta_{\text{old}} + \alpha \nabla_\theta V(\theta_{\text{old}}) \tag{3}$$

In order to do such a thing, we need to be able to calculate or at least estimate $\nabla V(\theta)$. To do so, we need to compute or estimate the $\nabla J_\theta(x)$. But the recursive nature of $J_\theta(x)$ makes that difficult.

$$
\begin{aligned}
\nabla J_\theta(x) &= \sum_a \left[ \nabla \pi_\theta(a|x) \right] \left[ r_{x,a} + \beta \sum_y p_{x,y}^a J_\theta(y) \right] + \sum_a \pi_\theta(a|x) \nabla \left[ r_{x,a} + \beta \sum_y p_{x,y}^a J_\theta(y) \right] \\
&= \sum_a \left[ \nabla \pi_\theta(a|x) \right] J_\theta(x,a) + \sum_a \pi_\theta(a|x) \nabla \left[ r_{x,a} + \beta \sum_y p_{x,y}^a J_\theta(y) \right].
\end{aligned}
\tag{4}
$$

Note the use of $J_\theta(x,a)$ notation to indicate the 'expected utility of taking action $a$ in state $x$, following policy $\pi_\theta$ beyond that point). Simplifying...

$$\nabla J_\theta(x) = \sum_a \left[ \nabla \pi_\theta(a|x) \right] J_\theta(x,a) + \beta \sum_a \pi_\theta(a|x) \sum_y p_{x,y}^a \nabla \left[ J_\theta(y) \right]. \tag{5}$$

It's convenient to re-arrange the second sum in the following way:

$$\nabla J_\theta(x) = \sum_a \left[ \nabla \pi_\theta(a|x) \right] J_\theta(x,a) + \beta \sum_y \nabla \left[ J_\theta(y) \right] \sum_a \pi_\theta(a|x) p_{x,y}^a. \tag{6}$$

Note that the summation on the left can be interpreted nicely, as the probability from transitioning from $x$ to $y$ in one step *under the policy* $\pi_\theta$. Hence,

$$\nabla J_\theta(x) = \sum_a \left[ \nabla \pi_\theta(a|x) \right] J_\theta(x,a) + \beta \sum_y \nabla \left[ J_\theta(y) \right] P_1^\theta(x \to y). \tag{7}$$

Note in the above - we want to compute $\nabla J_\theta(x)$. We can compute $\nabla \pi_\theta$ based on the structure of the policy. We can *estimate* $J_\theta(x,a)$ from data. Similarly, we can estimate $P_1^\theta(x \to y)$ from data. The only unknown term in the above

(relatively speaking) is $\nabla [J_\theta(y)]$. We can start to nest these, however. Note that for any $x_0, x_1, x_2, \ldots$:

$$\nabla J_\theta(x_0) = \sum_{a_0} [\nabla \pi_\theta(a_0|x_0)] J_\theta(x_0, a_0) + \beta \sum_{x_1} \nabla [J_\theta(x_1)] P_1^\theta(x_0 \to x_1)$$

$$\nabla J_\theta(x_1) = \sum_{a_1} [\nabla \pi_\theta(a_1|x_1)] J_\theta(x_1, a_1) + \beta \sum_{x_2} \nabla [J_\theta(x_2)] P_1^\theta(x_1 \to x_2) \tag{8}$$

$$\nabla J_\theta(x_2) = \sum_{a_2} [\nabla \pi_\theta(a_2|x_2)] J_\theta(x_2, a_2) + \beta \sum_{x_3} \nabla [J_\theta(x_3)] P_1^\theta(x_2 \to x_3)$$

$$\ldots$$

Substituting in, we can push off the unknown gradient-

$$\begin{aligned}
\nabla J_\theta(x_0) &= \sum_{a_0} [\nabla \pi_\theta(a_0|x_0)] J_\theta(x_0, a_0) + \beta \sum_{x_1} P_1^\theta(x_0 \to x_1) \nabla [J_\theta(x_1)] \\
&= \sum_{a_0} [\nabla \pi_\theta(a_0|x_0)] J_\theta(x_0, a_0) \\
&\quad + \beta \sum_{x_1} P_1^\theta(x_0 \to x_1) \left[ \sum_{a_1} [\nabla \pi_\theta(a_1|x_1)] J_\theta(x_1, a_1) + \beta \sum_{x_2} \nabla [J_\theta(x_2)] P_1^\theta(x_1 \to x_2) \right] \\
&= \sum_{a_0} [\nabla \pi_\theta(a_0|x_0)] J_\theta(x_0, a_0) \\
&\quad + \beta \sum_{x_1} P_1^\theta(x_0 \to x_1) \sum_{a_1} [\nabla \pi_\theta(a_1|x_1)] J_\theta(x_1, a_1) \\
&\quad + \beta^2 \sum_{x_1} \sum_{x_2} P_1^\theta(x_0 \to x_1) P_1^\theta(x_1 \to x_2) \nabla [J_\theta(x_2)]
\end{aligned} \tag{9}$$

Note that the last sum can be re-arranged, so that $\sum_{x_1} P_1^\theta(x_0 \to x_1) P_1^\theta(x_1 \to x_2) = P_2^\theta(x_0 \to x_2)$, the probability of going from $x_0$ to $x_2$ in two steps, under policy $\pi_\theta$.

$$\begin{aligned}
\nabla J_\theta(x_0) &= \sum_{a_0} [\nabla \pi_\theta(a_0|x_0)] J_\theta(x_0, a_0) \\
&\quad + \beta \sum_{x_1} P_1^\theta(x_0 \to x_1) \sum_{a_1} [\nabla \pi_\theta(a_1|x_1)] J_\theta(x_1, a_1) \\
&\quad + \beta^2 \sum_{x_2} P_2^\theta(x_0 \to x_2) \nabla [J_\theta(x_2)]
\end{aligned} \tag{10}$$

We can continue on this way, pushing hte unknown $\nabla J_\theta(x_i)$ further and further into the future:

$$\begin{aligned}
\nabla J_\theta(x_0) &= \sum_{a_0} [\nabla \pi_\theta(a_0|x_0)] J_\theta(x_0, a_0) \\
&\quad + \beta \sum_{x_1} P_1^\theta(x_0 \to x_1) \sum_{a_1} [\nabla \pi_\theta(a_1|x_1)] J_\theta(x_1, a_1) \\
&\quad + \beta^2 \sum_{x_2} P_2^\theta(x_0 \to x_2) \sum_{a_2} [\nabla \pi_\theta(a_2|x_2)] J_\theta(x_2, a_2) \\
&\quad + \beta^3 \sum_{x_3} P_3^\theta(x_0 \to x_3) \sum_{a_3} [\nabla \pi_\theta(a_3|x_3)] J_\theta(x_3, a_3) \\
&\quad + \ldots
\end{aligned} \tag{11}$$

or

$$\nabla J_\theta(x_0) = \sum_{k=0}^{\infty} \beta^k \sum_{x_k} P_k^\theta(x_0 \to x_k) \sum_{a_k} [\nabla \pi_\theta(a_k|x_k)] J_\theta(x_k, a_k) \tag{12}$$

We can tie this all together ultimately in the following way:

$$\nabla V(\theta) = \sum_{x_0} p_0(x_0) \nabla J_\theta(x_0) = \sum_{k=0}^{\infty} \beta^k \sum_{x_k} \sum_{x_0} p_0(x_0) P_k^\theta(x_0 \to x_k) \sum_{a_k} [\nabla \pi_\theta(a_k|x_k)] J_\theta(x_k, a_k) \tag{13}$$

Noting that $\sum_{x_0} p_0(x_0) P_k^\theta(x_0 \to x_k)$ is simply the probability that the $k$-th state $X_k$ is $x_k$, when the steps are executed under policy $\pi_\theta$,

$$\nabla V(\theta) = \sum_{k=0}^{\infty} \beta^k \sum_{x_k} \mathbb{P}_\theta(X_k = x_k) \sum_{a_k} [\nabla \pi_\theta(a_k|x_k)] J_\theta(x_k, a_k) \tag{14}$$

It's worth noting in the above that at this point, everything is technically known, or can be estimated from data collected over simulation. (We would imagine truncating the summation for instance if $\beta^k$ becomes sufficiently small, or if the game 'end'.) The issue to some extent is the sum over the states $x_k$ and actions $a_k$. To that end, the above is frequently rewritten in the following way:

$$\begin{aligned}
\nabla V(\theta) &= \sum_{k=0}^{\infty} \beta^k \sum_{x_k} \mathbb{P}_\theta(X_k = x_k) \sum_{a_k} \pi_\theta(a_k|x_k) \left[ \frac{\nabla \pi_\theta(a_k|x_k)}{\pi_\theta(a_k|x_k)} \right] J_\theta(x_k, a_k) \\
&= \sum_{k=0}^{\infty} \beta^k \sum_{x_k} \mathbb{P}_\theta(X_k = x_k) \sum_{a_k} \pi_\theta(a_k|x_k) \nabla \ln[\pi_\theta(a_k|x_k)] J_\theta(x_k, a_k) \\
&= \sum_{k=0}^{\infty} \beta^k \sum_{x_k, a_k} \mathbb{P}_\theta(X_k = x_k, A_k = a_k) \nabla \ln[\pi_\theta(a_k|x_k)] J_\theta(x_k, a_k) \\
&= \sum_{k=0}^{\infty} \beta^k \mathbb{E}[\nabla \ln[\pi_\theta(A_k|X_k)] J_\theta(X_k, A_k)]
\end{aligned} \tag{15}$$

Note in the above, we now have an expression for $\nabla V$ in which everything is either computable ($\nabla \pi_\theta$), or estimable from data ($J_\theta(X_k, A_k)$). Implementation of Policy Gradient strategies all take the approach of various estimation schemes of $J_\theta$ in order to try to accurately estimate the above from the data and update the policy accordingly.