# APS 1080 - LECTURE 3
## MODEL BASED RL - DYNAMIC PROGRAMMING

$$M \underset{S, R}{\overset{A}{\rightleftharpoons}} E$$

~Goal

So $A_0$ $R, S_1 A_1$ ... $S_i | A_i$ $R_{i+1}$ $S_{i+1}$ $A_{i+1}$ ... $S_T$

What Action to take if
you're in State $i$.

Return: Sum of subsequent rewards as a result of our action

$$G_i = R_{i+1} + \gamma R_{i+2} + \gamma^2 R_{i+3} + \dots$$

Goal: Select $A_i$ such that the expected value of the return is
   maximized
$\downarrow$
Action

$$\max \mathbb{E}[G_i]$$

## VALUE FUNCTION

$$V_\pi(s) = \mathbb{E}[G_t | S_t = s]$$

$$q_\pi(s) = \dots$$

Given $V_\pi$

$P$ : model of the environment $\Big\}$ We can solve Action Selection Problem.

To solve Action Selection Problem, you need either:

① $V$ and $P \implies \pi$

② $Q \implies \pi$

$\downarrow$

The Policy, which is the mechanism for action selection.

The states are drawn from $S = \left\{ S^{(1)}, S^{(2)}, ..., S^{(N)} \right\}$

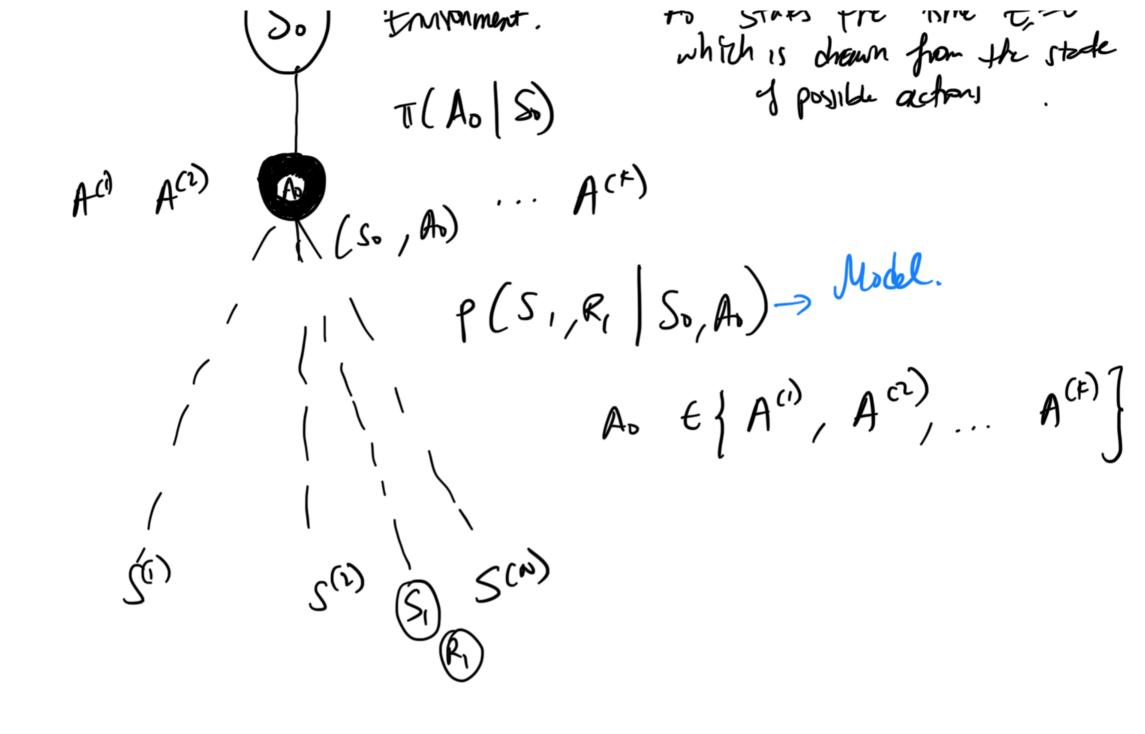| $S$ | $V_\pi(s)$ |
|---|---|
| $S^{(1)}$ | $\mathbb{E}\left[ G_t \mid S_t = S^{(1)} \right]$ |
| $S^{(2)}$ | |
| $\vdots$ | |
| $S^{(N)}$ | |

Target state: which state to go to

then consult "$P$" to see what action to take to get to that probability.

$$V_\pi(s) = \mathbb{E}\left[ G_t \mid S_t = S \right]$$

Subscript: Time sequence
Superscript: Selection from the set.

$S_0$

Environment.

to states are time $t_{i=0}$ which is drawn from the state of possible actions.

$\pi(A_0 \mid S_0)$

$A^{(1)} \quad A^{(2)}$

$A_0$

$A^{(F)}$

$(S_0, A_0)$

$p(S_1, R_1 \mid S_0, A_0) \rightarrow$ Model.

$A_0 \in \{A^{(1)}, A^{(2)}, \dots A^{(F)}\}$

$S^{(1)} \quad S^{(2)} \quad S_1 \quad S^{(N)}$

$R_1$

This is how $\pi$ and $p$ interplay in continuous state/action rewards.
$S \rightarrow A$ transition

Action $\rightarrow$ State (weighted by distribution $p$)

The model we may or may not have it.

The

with $p$ we can calculate $V_\pi, \pi$ and $\pi^*$
$\downarrow$
optimal policy.

Dynamic programming does NOT use experience

How we determine $V_\pi$ and $\pi$ when we don't have $P$?
$\rightarrow$ you need experience + learning
to be surrogates of

$$V_\pi(S) = \sum_{a \in A} \pi(a|s) \sum_{s' \in S} \sum_{r \in \mathbb{R}} \left[ p(s',r|s,a) \left[ r + \gamma V_\pi(s') \right] \right]$$

Probability of selecting an action given you're in state s.

Indices of a for loop

for loop

consider all possible states

over all possible rewards

Expected value → Weighted Average over the entire tree.

You're building the Agent.

No knowledge about the problem → $\pi$ would be an equiprobable distribution.
→ Eventually, you'd get to a more granular $\pi$ with a less naive distribution.

We have an equation that is recursively defined.

$$V_\pi(S) = \int^\Gamma V_\pi(\cdot)$$

gamma

We can calculate $V_\pi(\cdot)$ via linear system methods but it's not very scalable because you end up with too many variables and it's very computationally expensive. Also, it would only work if you have "p".

Instead of doing that, we calculate $V_\pi(\cdot)$ via a successive approximation.

$$f(\cdot) = \Gamma^* (f(\cdot))$$

$\uparrow$

Contraction

$f(\cdot) \longleftarrow$ initiate $\dots$ function $\forall X$

loop

    loop $\forall X$

        prev $\leftarrow f(x)$

        next $\leftarrow \Gamma^* (prev)$

        $f(x) \leftarrow$ next

        $\Delta(x) \leftarrow |$ next $-$ prev $|$

    if $\|\Delta\| < 0$

        $\longrightarrow$ break.

Initialize $V_\pi(x)$ for all $x$ arbitrarily.
We set the terminal states to zero and all the other values at an arbitrary number.

$V_\pi$

Stales | 

$S^{(1)}$
$S^{(2)}$
$\vdots$
$S^{(N)}$

$$V_\pi(S_T) = 0$$

| $\pi$ | $a^{(1)}$ $a^{(2)}$ ... $a^{(k)}$ |
|---|---|
| $S^{(1)}$ | |
| $S^{(2)}$ | |
| $\vdots$ | |
| $S^{(N)}$ | |

$\Sigma = 1$
$\Sigma = 1$

Input $\diamondsuit \pi$

Init $V_\pi(s)$ arbitrarily $\forall s$

loop $\Delta \longleftarrow 0$

       loop $\forall S$

            prev $\longleftarrow V_\pi(s)$

            Next $\longleftarrow \Gamma\left(V_\pi(\cdot)\right)$

$$\sum_\pi \ldots \Sigma\Sigma \left[\not{p} + \gamma V_\pi(s')\right]$$

      $\Delta \longleftarrow \max\left[\Delta, (\text{prev} - \text{next})\right]$

if $\Delta < \theta$:

    Break

Iterate over the states, successively estimating the values of the value functions

when the deltas are low enough then you stop.

Given an arbitrary model and policy, we can determine a value function

$$P, \pi \longrightarrow V_\pi(\cdot) \qquad \longrightarrow \text{Evaluation of the policy.}$$

    Evaluation of $\pi$ in $E \sim$ "Prediction"

$$\left\{ \begin{array}{l} \pi \xrightarrow{\text{eval}} V_\pi(\cdot) \\[2ex] \pi \longrightarrow \pi' \end{array} \right.$$

eval

$$\begin{cases} \pi' \longrightarrow V_{\pi}'(\cdot) \\ \pi' \xrightarrow{V_{\pi}'} \pi'' \\ \quad\quad Improvement \end{cases}$$

$$\begin{cases} \pi'' \longrightarrow V_{\pi}''(\cdot) \\ \pi'' \longrightarrow \pi''' \end{cases}$$

$$\left. \begin{array}{l} until \quad V_{\pi}(i) \simeq V_{\pi}(i+1) \\ \quad\quad \pi(i) \simeq \pi(i+1) \end{array} \right] \begin{array}{l} Fixed\ Point \\ \quad \pi_* \end{array}$$

You modify the policy to help you improve the likelihood of going to the most desirable states.

Summarize Dynamic Programming

① $p, \pi_0 \longrightarrow V_{\pi_0}$      (Prediction)

② $\pi_0, V_0 \longrightarrow \pi_1$    $(\pi_1 \geq \pi_0)$    (Improvement)

③ Generalized Policy Iteration (G.P.I)

           Prediction
           Improvement

           Until Fixed Point is reached $\longrightarrow \pi_*$

④ GPI $\longrightarrow$ Evaluation   ①
                Improvement   ②

                Evaluation   ①
                Improvement   ②

GPI: Eval Improv. Eval Improv. $\Big|$ $\Big|$ Eval Improv. Eval

$\pi_0$ ① ② ① ② $\Big|$ ... $\Big|$ $V_{\pi_n}$ $V_{\pi_n}$ $V_{\pi_n}$

Fixed Point

$V_{\pi*}$ , $\pi*$

This is an asymptotic Process

Other approach

Run it a few times and truncate the evaluation.

Select the best value amongst bag of truncated evaluations.

Step ④ helps us move away from the need of having a model.

What are other ways to approximate a value function that would not require the use of a "$\rho$"?

Next week:

How to approximate $V_\pi$ without $P$ by leveraging experience?

→ How can we approximate $q_\pi$ and a policy from experience?
└→ $\pi$