# RL Lecture 3:

**Dynamic      Programming      policy (for us)**

↓                    ↓                    ↓

sequential      optimize a (program)      Idea → Break into sub-problems

nature of
problem                                    ↓

(which we      Problems should have      Combine solutions of
have)          2 Properties              Sub-problems.

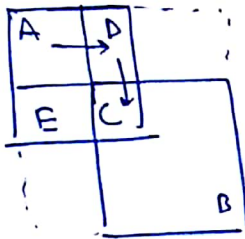① Optimal substructure → problem consists of
                                          smaller similar problems

② Overlapping subproblems          e.g.          smallest distance

The subproblems should                      similar      from A to B, can
occur multiple times.                      subproblems ←      be viewed as
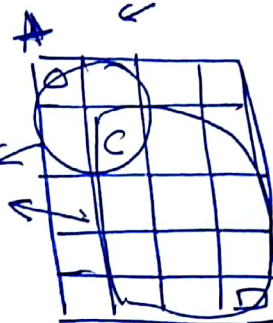                                                          smallest distance
e.g.                                                      from A to C &
                                                          then C → B.



Suppose I tried
path A → D and then
went from D → C and I need to solve the
problem of C → B now.
However, if I go from A → E & then E → C.

This is an          ← I still need to solve the same problem for C → B.
example of overlapping
subproblems.

        MDP ⇸ satisfies both → Bellmann equatz gives recursive
                                            decomposition

In DP → solutions can be cached.

In RL, Value functz → cache of the informatz already figured out
about the state → then work the way backwards to get
the result.

                              Knowledge of
Recall:   Planning → ∧MDP already given (eg. rules of a game)
          and then solve the MDP.

2 types of problems :
     ① Prediction → output the value functz $V_\pi$.
          I/p → MDP / MRP      and policy

② Control → MDP given but we need to figure out the most important / optimal policy

o/p → optimal value function $V_\pi$ ( hence, also the optimal policy)

I/p → MDP.

## Policy Evaluation (MDP & Policy given, what is the reward)

Bellman expectation equation to be used in an iterative manner.

$V_i$ → value function at each iteration

$V_1$ → $\begin{bmatrix} V_{1s_1} \\ V_{1s_2} \\ \vdots \\ V_{1s_n} \end{bmatrix}$ → initial value function → 0 is a safe default.

Then, using the 2-step lookahead → find the next new value function

$V_2$ → . . . so on.

$V_\pi$ → final optimal value function.

Synchronous Backups:

$\begin{bmatrix} V_{1s_1} \\ \vdots \\ V_{1s_n} \end{bmatrix}$ $\xrightarrow[\text{states}]{\text{Update all the}}$ $\begin{bmatrix} V_{2s_1} \\ \vdots \\ V_{2s_n} \end{bmatrix}$ . . . . →

So, at each step, we are updating all the states using the value function at the last iteration → Synchronous backup.

at each iteration $k+1$:

→ update all states from previous iteration

→ $V_{k+1}(s)$ from $V_k(s')$

→ $s'$ → successor state of $s$.

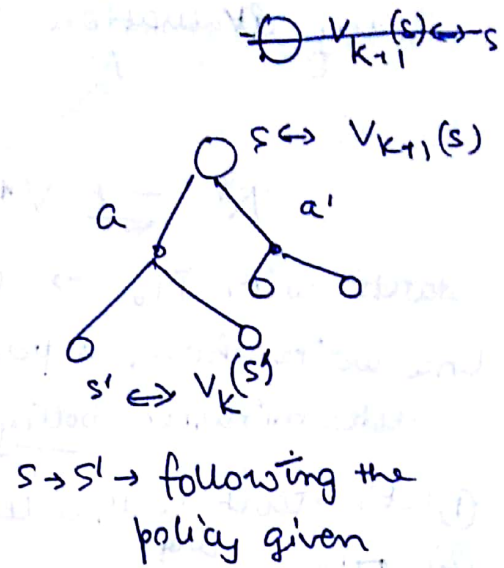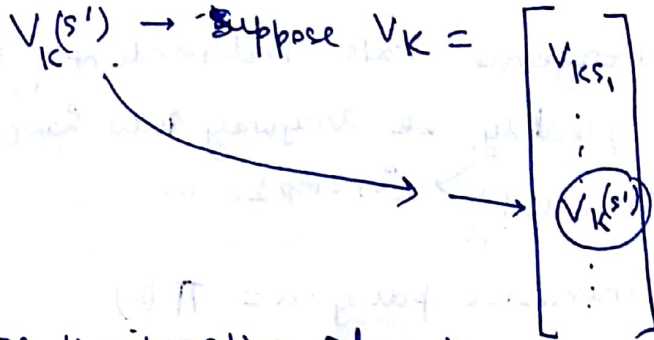This does converge!!

[ Expectate backup → as Expectation over backups.
  Rather than max → in the case of Optimal Policy ]

use the 2-step lookahead tree. ~~lookahead~~.

vector)

$$V^{k+1} = R^\pi + \gamma P^\pi V^k.$$



$V_k(s') \rightarrow$ Suppose $V_K = \begin{bmatrix} V_{ks_1} \\ \vdots \\ \boxed{V_k(s')} \\ \vdots \end{bmatrix}$

$s \rightarrow s' \rightarrow$ following the policy given

Here, the iterative algorithm
is that, at each state $s$, to compute the $V_{k+1}$ value for it,
we look at the values of the value function for all its
successors (1 step look-ahead) from the previous iteration

This is guaranteed to converge.

Also, even though we have one policy (eg. random policy ), the
Value function essentially helps us get a better policy
just by looking at the rewards and value functn of the
next step.

Policy iteration → improve the policy using a feedback
mechanism to feed in the improved policy obtained through
the value function.

Policy given → $\pi$
① Evaluate the policy
② Improve by acting greedily wrt $V_\pi$
$$\pi' = greedy(V_\pi) \rightarrow \text{whichever gives better value function, pick that!}$$

This process always converges to $\pi^*$.
(atleast 1 deterministic optimal policy exists).

Policy evaluation $\longrightarrow$ Policy improvement

$$\pi^* \geqslant \pi$$

$$\pi^* \rightleftarrows V^*.$$

starts with $\Pi_0 \rightarrow$ convergence rate independent of $\pi_0$.
Once we're chosen a policy greadily, we anyway now have a
deterministic policy. (Think!) $\rightarrow$ in step 1.

① Let's start with a deterministic policy : $a = \pi(s)$

② $\Pi'(s) = \underset{a \in A}{arg\ max}\ q_\pi(s, a) \rightarrow$ was in state $s$, took action $a$

and then followed policy $\pi$

③ $q'_\pi(s, \pi'(s))$

$\rightarrow$ what will be the value function.

$$= \underset{a \in A}{max}\ q_\pi(s, a) \geqslant \boxed{q_\pi(s, \pi(s))} \overset{②}{=} V_\pi(s)$$

∴ Atleast improves one
step following $\pi'(s)$

as from state $s$, followed
followed $\pi$ to get next
action and followed
$\pi$ from there
onwards

④ $V_{\pi'}(s) \geqslant V_\pi(s)$

$V_\pi(s) \leq \underset{\underbrace{\qquad}}{q_\pi(s, \pi'(s))}$

$\rightarrow$ unroll it

$$= E_{\Pi'}(R_{t+1} + \gamma V_\pi(s_{t+1}) \mid s_t = s)$$

$$\leq E_{\pi'}[R_{t+1} + \gamma q_\pi(s_{t+1}, \Pi'(s_{t+1})) \mid s_t = s]$$

$$\leq E_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \cdots \mid s_t = s]$$

$$= V_{\pi'}(s).$$

RL. lecture 3.

So, proved that following the greedy policy, we would, atleast
atleast equal value, if not, more. Still, haven't told that-
it'll keep increasing. / go to optimal.

If we have equality $\rightarrow$ improvements stop. !

$$q_\pi(s, \pi'(s)) = \max_{a \in A} q_\pi(s, a) \geq q_\pi(s, \pi(s))$$

$$= v_\pi(s)$$

then Bellman optimality

eq² satisfied

$$V_\pi(s) = \max_{a \in A} q_\pi(s, a) \Rightarrow V_\pi(s) = V_*(s) \quad \forall s \in S$$

$$\pi \rightarrow \text{optimal}$$
$$\text{policy.}$$

∴. Policy iteration solves MDP.

Question $\rightarrow$ If in the 3rd iterate itself, we can get the optimal
values, iterations 4 to ∞ are wasted. $\rightarrow$ Any way to
truncate the iterative procedure / do an approximation?

① 
Modified Policy Iteration $\rightarrow$ add an early stopping criterion
eg. if the Bellman eq² value difference
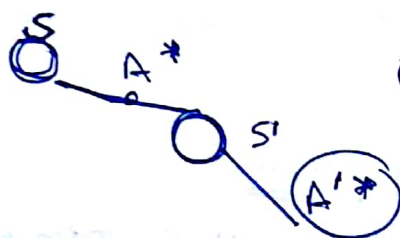in consequent iterations is < $\varepsilon$.

② Train / Run for only K iterations. eg (K=3 above).
Extreme case $\rightarrow$ K=1 $\rightarrow$ value iteration ⟵ ( most famous
use g DP).

Optimal policy

$\downarrow$ ① Take optimal action 1st $\rightarrow A_*$

② wherever you land $\rightarrow$ the next state, because of this
action, it should be optimal from here too.

$\pi(a|s) \rightarrow$ optimal when :

① $V_\pi(s) = V_*(s)$

② for all reachable s',
$V_\pi(s') = V_*(s')$.

S ⊙ $A^*$ ⟍ ⊙ S' ⊙ $A'^*$

# Deterministic Value Iteration:

Assumption → already know the best solution for the next step. → $s' → V_*(s')$. → use it to find the optimal sol$^n$ for the previous step.

$$V_*(s) = \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s') \right)$$

↓

Can be found by one-step look ahead.

| Difference bet$^n$ policy iterat$^n$ & value iteration?

Policy iteration → Start with random policy, find the value funct$^n$ and improve policy based on this value function.

Value Iteration → Start with random value function and keep improving that value function

# Policy Iteration

① Initialization → $\pi(s) \in A(s)$ (arbitrary)

② Policy evaluation:

  Repeat       as becomes deterministic on greedy choice ↓

    $\Delta \Leftarrow 0$
    for all $s \in S$:   $a = \pi(s)$
     $V_{k+1}(s) = \cancel{\max}$ ,

$$R_s^a + \gamma \sum_{s'} P_{ss'}^a V_k(s')$$

    $\Delta = \max(\Delta, |V_{k+1}(s) - V_k(s)|)$
   if $\Delta < \theta →$ stop

③ Policy improvement:
  for all $s \in S$
   $\pi(s) = \max_a q_\pi(s,a)$.
  if no change in $\pi(s)$ for all $s$, then policy stable and stop. Else go to ②.

# Value Iterat?

→ ✦ Initialization : $V_*(s) = 0 \; \forall \; s \in S$  $\qquad$ eg.

→ Repeat
$\qquad$ for all $s \in S$ :

$$V_*(s) = \max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'} V_*(s') \right).$$

$\qquad$ if change for all $s$, very small ($< \varepsilon$). → stop.

→ for all $s$,

$$\Pi(s) = \arg\max_a \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a V_*(s') \right). \quad \rceil \; \substack{\text{policy} \\ \text{updated} \\ \downarrow \text{Just} \\ \text{once}} \quad \text{not}$$

We are not solving the full RL problem as MDP's knowledge already given.

Value Iteration → any intermediate value (before converging) may not actually be the value funct? for any policy, ie..
may not be $V_\pi$ for any $\pi$.
$\qquad$ But at the end, we get the value funct? for the optimal policy.

Bellman expectat? equation

$$\to \quad V_{\Pi}(s) = \sum_{a \in A} \Pi(a/s) \left[ R_s^a + \gamma \sum_{s'} P_{ss'} V_*(s') \right]$$

$\qquad$ Since $\pi(s)$ becomes deterministic, $\pi(a/s) = 1$ for one action and $0$ for the rest $a \in A$. (leading to this)

# Extensions to DP

Till now synchronous.

Asynchronous → pick any state and update its value

funct² → $V_{k+1}(s)$ (for eg)

now immediately use this value for the calculation of values
function for other states instead of $V_k(s)$.

⇒ Breaks iteration.

→ Better computation

① In-place

Basically, as soon as we update a value
iteration ⟶ use it directly

② Prioritized sweeping → use mag² g Bellman error → to
choose which states to update 1st.
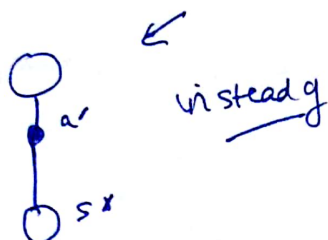
Bellman error $= \left| \max_{a \in A} \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'} V(s') \right) - V(s) \right|$

③ Real time DP → Use only states that are actually visited
by the agent → use agent's experience → real-world
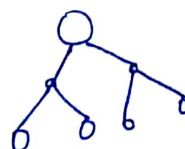
$$S_t = \max_{a \in A} \left( R_{s_t}^a + \gamma \sum_{s'} P_{ss'}^a V(s'_t) \right)$$

DP → does full width backups → i.e. goes through all the actions
and subsequently all the states. to compute its predictions

→ not feasible and we need to know the full model
dynamics

Instead → we'll do → sample backups → sample an
action, sample subsequent state $s$... etc and
then do the backup instead of full



vs instead of

④ helps overcome an
dimensionality

eg.g sampling with DP

Approximate DP → Approx. value function → function approximate $\hat{V}(s, \omega)$

### Fitted Value Iteration

**Repeat**

$s$ states sampled

① Sample $\mathcal{S} \subseteq S$

② for each $s \in \mathcal{S}$,

$$\tilde{V}_k(s) = \max_{a \in A} \left( R_s^a + \gamma \sum_{s'} P_{ss'}^a \hat{V}(s', \omega_k) \right)$$

Use $\{ s, \tilde{V}_k(s) \}$ → to **train** $\hat{V}(s', \omega_{k+1})$

Contraction mapping → confirms the convergence of value & policy iteration.

Vector space →



→ dimension = # states = $|S|$

Distance betz state value functions $u$ & $v$. → $\infty$ norm

Want to show that Bellman backup brings value functions closer in subsequent iteration and hence, it converges.

$$\Rightarrow \|U - V\|_\infty = \max_{s \in S} |u(s) - v(s)|$$

$\gamma$ - contraction

Bellman expectats operator → $T^\pi$ → brings value functions closer by atleast $\gamma$. times

$$T^\pi_{(v)} = R^\pi + \gamma P^\pi V$$

$$\| T^\pi(U) - T^\pi(V) \|_\infty \leq \| \gamma P^\pi \|V - U\|_\infty \|$$

$$\leq \gamma \|U - V\|_\infty$$

now, if U, V are not already same, then ↗ earlier dist was only $\|U-V\|_\infty$.

following contraction mapping theorem

iterative policy evaluats converges to $V_\pi$.

∴ Policy iterats → $V_*$.

for vector space V, closed under operator T(v), T → $\gamma$-contracts

• T converges to a unique fixed point

• Linear convergence rate → $\gamma$..

1  Bellmann optimality backup operator $\rightarrow T^*$

$$T^*(v) = \max_a \underline{R^a + \gamma P^a v} \rightarrow$$ similarly a $\gamma$-contraction

$\rightarrow$ hence, ~~policy~~ value iterate also converges.