

If we land up at π' , which is not better but as good as π , are π & π' optimal policies?

\Downarrow

$$V_{\pi'} = V_{\pi} \text{ for } \forall s \in \mathcal{S} \quad - (1)$$

$$\pi'(s) = \arg \max_a Q_{\pi}(s, a) \quad - (2)$$

with (1), will (2) imply the following?

$$V_{\pi'}(s) = \max_a E[R_{t+1} + \gamma V_{\pi'}(S_{t+1}) | S_t = s, A_t = a]$$

$$V_{\pi}(s) \leq Q_{\pi}(s, \pi'(s)) \\ = V_{\pi'}(s)$$

$$Q_{\pi}(s, a) = E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = a]$$

Model-based
MDP

$$P(s', r | s, a)$$

availability of
model implies
you know this
probability -
which is a 4-arg-
ument function

We first discuss policy iteration (PI)

Refer to PI pseudocode on pg. 80
from T.B.

$$A = \{1, 1, 0, -1\}$$
$$\arg \max_i A$$
$$= 0, 1 \text{ or } 1, 2$$

1. $P \in \rightarrow$ given π , you find $V_{\pi}(s) \forall s \in \mathcal{S}$ by iteratively solving Bellman equations.

2. P.I. \rightarrow given $\pi \rightarrow V_{\pi}(s) \xrightarrow{\forall s} \pi' \rightarrow V_{\pi'}(s) \xrightarrow{\forall s} \pi'' \rightarrow V_{\pi''}(s) \dots$

until you locate
 π^* & $V_{\pi^*}(s)$

Coming back to last slide —

$$V_{\pi}(s) = V_{\pi'}(s) = Q_{\pi}(s, \pi'(s))$$
$$= \max_a Q_{\pi}(s, a)$$