How good is a policy? → look at value / utility

$$V(s_0, s_1, \dots s_T) = \sum_{t=0}^{T} R(s_t)$$

$$V(s_0, s_1, \dots s_T) = \sum_{t=0}^{T} \gamma^t R(s_t) \quad \gamma \in [0,1]$$

$$V(s)^\pi = E_{Pr}(s_0, s_1, \dots | s_0 = s, \pi) \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$

$$\pi^*_{(s)} = \arg\max_\pi V(s)^\pi \longrightarrow \pi^*_{(s)} = \arg\max_{a \in A(s)} \sum_{s'} P(s'|s,a) \, Value(s')$$
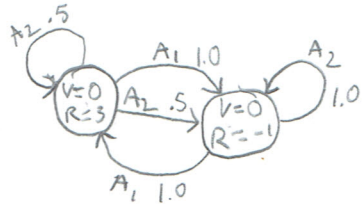
$$V(s) = R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) V(s') \qquad \text{Bellman equation}$$

For $i = 1$ to <stopping criteria>

$$V_{i+1}(s) \leftarrow R(s) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s,a) V(s')$$

$$i \leftarrow i + 1 \qquad \text{Value iteration}$$

## Ex 1



$$S = \{s_1, s_2\}$$
$$A = \{a_1, a_2\}$$
$$R(s_1) = 3$$
$$R(s_2) = -1$$
$$\gamma = .5$$

$$T = \begin{bmatrix} [P(s_1|s_1,a_1), \; P(s_1|s_1,a_2)], \\ [P(s_1|s_2,a_1), \; P(s_1|s_2,a_2)], \\ [P(s_2|s_1,a_1), \; P(s_2|s_1,a_2)], \\ [P(s_2|s_2,a_1), \; P(s_2|s_2,a_2)] \end{bmatrix}$$

$i = 1, \quad V_0(s_1) = 0, \quad V_0(s_2) = 0$

$$V_1(s_1) = R(s_1) + \gamma \max_{a \in A(s)} \sum_{s'} P(s'|s_1, a) V(s')$$

$$= 3 + .5 \times \max \begin{bmatrix} \begin{pmatrix} P(s'=s_1|s_1, a=a_1)V(s_1) \\ + \\ P(s'=s_2|s_1, a=a_1)V(s_2) \end{pmatrix} \\ \begin{pmatrix} P(s'=s_1|s_1, a=a_2)V(s_1) \\ + \\ P(s'=s_2|s_1, a=a_2)V(s_2) \end{pmatrix} \end{bmatrix}$$

$\rightarrow 0 \cdot V(s_1) \rightarrow 0 \cdot 0 = 0$ } $0 + 0 = 0$
$\rightarrow 1.0 \cdot V(s_2) \rightarrow 1.0 \cdot 0 = 0$
$\rightarrow .5 \cdot V(s_1) \rightarrow .5 \cdot 0 = 0$ } $0 + 0 = 0$
$\rightarrow .5 \cdot V(s_2) \rightarrow .5 \cdot 0 = 0$

$$= 3$$

$$\boxed{\begin{array}{l} V_1(s_1) = 3 \\ V_1(s_2) = -1 \end{array}} \longrightarrow * \text{ when } V_0(s) = 0, \; V_1(s) = R(s), \text{ therefore } V_1(s_2) = R(s_2)$$

$$= -1$$

$$T = \left[ \begin{array}{c} [\overset{111}{0}, \overset{112}{.5}], [\overset{121}{1.0}, \overset{122}{0}] \\ [\underset{211}{1.0}, \underset{212}{.5}], [\underset{221}{0}, \underset{222}{1.0}] \end{array} \right]$$

$i = 2, \; V_1(s_1) = 3, \; V_1(s_2) = -1$

$$V_2(s_1) = R(s_1) + \gamma \max_{a \in A(s)} \left[ \overset{s' \in S}{\underset{s'}{\sum}} P(s'|s, a) V(s') \right]$$

$$= 3 + .5 \cdot \max \left[ \begin{array}{l} P(s'=s_1|s_1, a=a_1)V(s_1) + P(s'=s_2|s_1, a=a_1)V(s_2), \\ P(s'=s_1|s_1, a=a_2)V(s_1) + P(s'=s_2|s_1, a=a_2)V(s_2) \end{array} \right]$$

$\overset{\text{s' number}}{\downarrow} \quad \overset{\text{state s number}}{\downarrow} \quad \overset{\text{action number}}{\downarrow}$

$$= 3 + .5 \cdot \max \left[ \begin{array}{l} T(1,1,1)V(s_1) + T(2,1,1)V(s_2), \\ T(1,1,2)V(s_1) + T(2,1,2)V(s_2) \end{array} \right]$$

$$= 3 + .5 \cdot \max \left( 0 \cdot 3 + 1.0 \cdot -1, \; .5 \cdot 3 + .5 \cdot -1 \right)$$

$$= (-1, \quad 1)$$

$$= 3 + .5(1) = 3.5 \qquad \boxed{V_2(s_1) = 3.5}$$

$$V_2(s_2) = R(s_2) + \gamma \max_{a \in A(s)} \left[ \overset{s' \in S}{\underset{s'}{\sum}} P(s'|s_2, a) V(s') \right]$$

$$= -1 + .5 \max \left[ \begin{array}{l} P(s'=s_1|s_2, a=a_1)V(s_1) + P(s'=s_2|s_2, a=a_1)V(s_2), \\ P(s'=s_1|s_2, a=a_2)V(s_1) + P(s'=s_2|s_2, a=a_2)V(s_2) \end{array} \right]$$

$$= -1 + .5 \max \left[ \begin{array}{l} T(1,2,1)V(s_1) + T(2,2,1)V(s_2), \\ T(1,2,2)V(s_1) + T(2,2,2)V(s_2) \end{array} \right]$$

$$= -1 + .5 \cdot \max \left( 1.0 \cdot 3.5 + 0 \cdot -1, \; 0 \cdot 3.5 + 1.0 \cdot -1 \right)$$

$$= (3.5, \quad -1)$$

$$= -1 + .5 \cdot 3.5 = .75 \qquad \boxed{V_2(s_2) = .75}$$

$i \in$ epochs (iteration)

$j \in |S|$ (current state)

$k \in |S|$ (next state)

$n \in |A|$ (current action)

$T(1,1,1) = 0$  $T(2,1,1) = 1.0$

$T(1,1,2) = .5$  $T(2,1,2) = .5$

$T(1,2,1) = 1.0$  $T(2,2,1) = 0$

$T(1,2,2) = 0$  $T(2,2,2) = 1.0$

$V_2(s_1) = 3.5$

$V_2(s_2) = .75$

$$V_i(s_j) = R(s_j) + \gamma \max_{a_n}^{a_n \in A(s)} \left[ \sum_{s_k} P(s_k | s_j, a_n) V(s_k) \right]$$

$$V_i(s_j) = R(s_j) + \gamma \max \left[ \underbrace{P(s_k = s_1 | s_j, a_n = a_1) V(s_k)}_{k=1, n=1} + \underbrace{P(s_k = s_2 | s_j, a_n = a_1) V(s_k)}_{k=2, n=1} \right.$$

for each action

$$\left. P(s_k = s_1 | s_j, a_n = a_2) V(s_k) + P(s_k = s_2 | s_j, a_n = a_2) V(s_k) \right]$$

Summation for each possible next state

$$V_i(s_j) = R(s_j) + \gamma \max \left[ \begin{array}{c} T(1,j,1)^{k,n} V(s_1^k) + T(2,j,1) V(s_2^k), \\ T(1,j,2) V(s_1) + T(2,j,2) V(s_2) \end{array} \right]$$

$$V_3(s_1) = 3 + .5 \cdot \max \left[ \begin{array}{c} \overbrace{T(1,1,1) V(s_1) + T(2,1,1) V(s_2)}^{k=1,j=1,n=1 \quad k=2,j=1,n=1}, \\ \underbrace{T(1,1,2) V(s_1) + T(2,1,2) V(s_2)}_{k=1,j=1,n=2 \quad k=2,j=1,n=2} \end{array} \right]$$

$$V_3(s_1) = 3 + .5 \cdot \max \left[ \underset{.75}{0 \cdot 3.5 + 1.0 \cdot .75}, \underset{2.125}{.5 \cdot 3.5 + .5 \cdot .75} \right]$$

$$\boxed{V_3(s_1) = 4.0625}$$

$$V_3(s_2) = -1 + .5 \cdot \max \left[ \begin{array}{c} T(1,2,1) V(s_1) + T(2,2,1) V(s_2), \\ T(1,2,2) V(s_1) + T(2,2,2) V(s_2) \end{array} \right]$$

$$= -1 + .5 \cdot \max \left[ \underset{4.0625}{1.0 \cdot 4.0625 + 0 \cdot .75}, \underset{.75}{0 \cdot 4.0625 + 1.0 \cdot .75} \right]$$

$$\boxed{V_3(s_2) = 1.03125}$$

$V_0(s_1 ... s_{|S|}) = 0$

for $i = 1$ to $\langle$stopping criteria$\rangle$

     for $j = 1$ to $|S|$     // for all states

         vals $\leftarrow \{\}$

         for $n = 1$ to $|A|$   // for all valid actions in $S_j$

             $vals_n \leftarrow 0$

             for $k = 1$ to $|S|$   // for all reachable states from $S_j$

                 $vals_n \leftarrow vals_n + T(k, j, n) V_{i-1}(s_k)$

         $V_i(s_j) \leftarrow R(s_j) + \gamma \cdot \max(vals)$

## Value iteration

now what?   How to use values to generate policy?