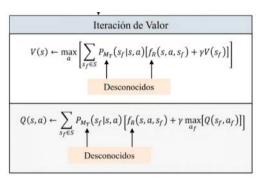# Week 10: V(s) and Q(s,a) for Bellman's Optimality

Bellman's Optimality Equations are defined in two ways: $V(s)$ and $Q(s,a)$ functions.
Both give us a set of equations for $V(s)$ and $Q$, respectively. In essence, they're the same.
- $\rightarrow V(s)$ is a vector (one variable function)
  - $\rightarrow Q(s,a)$ is a matrix (two variable function)

|  |  | Q(s,a) |  |
|---|---|---|---|
|  | a1 | a2 | a(m=3) |
| s1 | Q(s1,a1) | Q(s1,a2) | Q(s1,a3) |
| s2 | Q(s2,a1) | Q(s2,a2) | Q(s2,a3) |
| s3 | Q(s3,a1) | Q(s3,a2) | Q(s3,a3) |
| s4 | Q(s4,a1) | Q(s4,a2) | Q(s4,a3) |
| s(n=5) | Q(s5,a1) | Q(s5,a2) | Q(s5,a3) |

the max Q of each row, goes to $V(s)$ vector and the $Q(s,a)$ action is the action of $V(s_i)$

$\Rightarrow$ to know the action is to know Q matrix

$Q(s3,a1)$ is the accumulated reward of all trayectories that start in $s3$ and execute $a1$ on that $s3$ initial state

**Iteración de Valor**

$$V(s) \leftarrow \max_{a} \left[ \sum_{s_f \in S} P_{M_T}(s_f | s, a) [f_R(s, a, s_f) + \gamma V(s_f)] \right]$$

Desconocidos

$$Q(s,a) \leftarrow \sum_{s_f \in S} P_{M_T}(s_f | s, a) \left[ f_R(s, a, s_f) + \gamma \max_{a_f} [Q(s_f, a_f)] \right]$$

Desconocidos

The bellman equations will discover the politic: set of actions.
  $\rightarrow$ the optimal politic
    $\nearrow$ and thus $Q(s,a)$ too

$\rightarrow$ Since $V(s)$ involves $\max()$ function, the solution cannot be analytical: we will use a numerical method, Value Iteration.
For this, both the Transition Model and reward function must be known.

Taking the example world:

$\rightarrow$ The world contains the set of states $S = \{ S_1, s_2, s_3, sF_1, sF_2 \}$ where $s_1$ = initial state and, $sF_1$ and $sF_2$ are final states:

| SF₁ | S₁ | S₂ | S₃ | SF₂ |
|---|---|---|---|---|

$\rightarrow$ The world has the following set of Actions $A = \{ \rightarrow, \leftarrow \}$, where:
  - $\rightarrow$ = agent moves to left, one cell
  - $\leftarrow$ = agent moves to right, one cell  $\Big\}$ deterministic example

$\rightarrow$ The reward function $f_R(s, a, s_f) = f_R(s_f)$ only depends on the state that the agent arrives to.

$$f_R(s_f) = \begin{matrix} sF_1 \\ s_1 \\ s_2 \\ s_3 \\ sF_2 \end{matrix} \begin{bmatrix} -10 \\ 0 \\ -0.4 \\ -0.4 \\ 10 \end{bmatrix}$$

In the excel file, choose the sheet WITHOUT $F_\pi(s)$:
we have the calculation with $V(s)$ Value Iteration, under Bellman's Optimal Politic

## Calculation with V(s) Value Iteration, under Bellman's Optimal Politic

| | a1=← | a2=→ | | |
|---|---|---|---|---|
| | sF1 | s1 | s2 | s3 | sF2 |

| | sF1 | s1 | s2 | s3 | sF2 |
|---|---|---|---|---|---|
| fR(s)= | -10 | 0 | -0.04 | -0.04 | 10 |
| V(s)= | 0 | 0 | 0 | 0 | 0 |
| V(s)= | 0 | -0.04 | 0 | 10 | 10 |

→ input

γ= 0.9

**Bellman's Optimality Eq. for V(s):**

$$V(s) \leftarrow \max_a \left[ \sum_{s_f \in S} P_{M_T}(s_f | s, a)\left[ f_R(s,a,s_f) + \gamma V(s_f) \right] \right] \quad (1)$$

→ Reward func.

→ V(s) values: since it's a vector, we have a value per state

→ Init with zeroes.

$P_{MT} = 1$ for $s_1$, ∅ for rest

here, we write for each state the Bellman's Optimality Equation:

Formally
$$V(s) \leftarrow \max_a \left[ \sum_{s_f \in S} P_{M_T}(s_f | s, a)\left[ f_R(s,a,s_f) + \gamma V(s_f) \right] \right]$$

⟹ =MAX(E4+$E$10*E6, F4+$E$10*F6)

1 for current state, ∅ for the rest

maximum — deterministic

max of two terms: one for each action ($a_1 = ←$, $a_2 = →$)

1st iteration result: copy it and paste nums in above V(s) (input) → w/o format

Since it's deterministic, the sum Σ disappears, and thus we only got the $f_R(s_f) + \gamma V(s_f)$ term for each action.

for $a_1 = ←$ in the cell:

$(-10_1 + 0.9 * ∅_1, \quad ∅_2 + 0.9 * ∅_2)$

$a_1 = ←$ , $a_2 = →$

→ immediate reward after transition with $a_i$

All 5 V(s) must converge at the same iteration

we say Convergence when ΔV(s) is very small → we decide how small = 0.1 or 0.001 or 0.0001

v(s), $v_F$

Here we are looking for the Optimal Politic, and we will do so from these V(s) values when converged: V(s) values' approximation, but how do we know the action?

↳ each V(s) was calculated as the max($a_1$, $a_2$) thus if we write each of the two:

| Cell | =MAX(E4+$E$10*E6, F4+$E$10*F6) | |
|---|---|---|
| | $a_1$ ← | $a_2$ → |

| | | | |
|---|---|---|---|
| 55.542887 | x→ $a_2$ | 72.826047 | |
| 55.542887 | x→ $a_1$ | 80.918447 | |
| 72.826047 | x→ $a_2$ | 89.954447 | |
| 80.918447 | x→ $a_2$ | 99.994447 | |
| 89.954447 | x→ $a_2$ | 99.994447 | |

we need to write how the cell decided so that we can deduce the action.

Are the same from Q table

gives the maximum possible reward on each state.

↳ Optimal Politic: if agent started on $s_1, s_2, s_3, \dots$ it needs to move according to this politic: $a_2$ →

The solution to Bellman's Opt Eq are the 5 V(s) values, from which we can deduce the Optimal Politic.

## Calculation with Q(s,a) Value Iteration, under Bellman's Optimal Politic

move ← , move →, move → but frontier (sF)

| | a1=← | a2=→ | | |
|---|---|---|---|---|
| | sF1 | s1 | s2 | s3 | sF2 |

| | sF1 | s1 | s2 | s3 | sF2 | |
|---|---|---|---|---|---|---|
| fR(s)= | -10 | 0 | -0.04 | -0.04 | 10 | |
| input | a2=→ | a2=→ | a2=→ | a2=→ | a2=→ | |
| Q(s,a1=←)= | 0 | 0 | 0 | 0 | 0 | a1=← |
| Q(s,a2=→)= | 0 | 0 | 0 | 0 | 0 | a2=→ |
| formulas | | | | | | |
| Q(s,a1=←)= | -10 | -10 | 0 | -0.04 | -0.04 | a1=← → Eqs for $a_1 = ←$ |
| Q(s,a2=→)= | 0 | -0.04 | -0.04 | 10 | 10 | a2=→ → Eqs for $a_2 = →$ |

γ= 0.9

→ reward func.

This politic is discovered at the end ↳ Bellman's Equations outcome!

Whereas here, it is a matrix: two rows (actions) and 5 columns (states)

→ Init with zeroes

In this matrix we write the Bellman's Optimality Equations (10)

this max will tell us which action for the final Optimal Politic

For Q(s,a) this is the Bellman's Opt Equation System

$$Q(s,a) \leftarrow \sum_{s_f \in S} P_{M_T}(s_f | s, a)\left[ f_R(s,a,s_f) + \gamma \max_{a_f}[Q(s_f, a_f)] \right]$$

we will have 10 equations looking like this

→ Since it's the same deterministic problem, the sum and $P_{MT}$ disappear.

For this cell we have:

$$Q(s,a) \leftarrow \sum_{s_f \in S} P_{M_T}(s_f | s, a)\left[ f_R(s,a,s_f) + \gamma \max_{a_f}[Q(s_f, a_f)] \right]$$

⟹ =L4+$L$12*MAX(L6,L7)

Cell row is $a_1 = ←$

gamma, $Q(sF_1, ←)$ $Q(sF_1, ←)$

immediate reward

$-10 + 0.9 * \max_{a_f} ( ∅ \quad ∅ )$

Cell → represents the value of the average accumulated reward given that we started on state s, and took (on that state), action a: the average acc. reward of all trajectories that start on $s_1$ (it's column name) and execute $a_1 = ←$ (row name) on that $s_1$ state they started with, and from then on, they do whatever.

→ Iterate as before: put the formula matrix's results into the input matrix values and so on.

→ Not all systems of Equations converge under Value Iteration, but Bellman's Optimality Equation Systems always converge (either $V(s)$ or $Q(s,a)$ versions).

→ What determines the Optimal Politic is the Reward function $f_R(s,a,s_f)$

→ Both $V(s)$ and $Q(s,a)$ give the same output:

| $V(s)=$ | 71.8618263 | 79.9542263 | 88.9902263 | 99.0302263 | 99.0302263 | |
|---|---|---|---|---|---|---|
| $Q(s,a1=\leftarrow)=$ | 54.5786663 | 54.5786663 | 71.8618263 | 79.9542263 | 88.9902263 | a1=← |
| $Q(s,a2=\rightarrow)=$ | 71.8618263 | 79.9542263 | 88.9902263 | 99.0302263 | 99.0302263 | a2=→ |

$V(s)$ will take the max Q value, so in order examples $V(s)$ will have other Q row value.

→ $Q$ needs more memory : $V(s)$ calculates the same but when needed.
→ For $V(s)$ and $Q(s,a)$ we need $P_{MT}$ and $f_R$ to be known.