# ON-POLICY BELLMAN EQUATIONS

Fix $\sigma : S \to A$. How do we compute $val_\sigma$?

$$val_\sigma(s) = \mathbb{E}_\sigma\left[\sum_{i \geq 0} \gamma^i r_i\right]$$

$$val_\sigma(s) = \sum_{s', r} \delta(s', r \mid s, \sigma(s))\left[r + \gamma \, val_\sigma(s')\right]$$

$\hookrightarrow$ this is a set of linear equations with unknowns $(val_\sigma(s))_{s \in S}$

# OPERATOR POINT OF VIEW

Define $\quad F_S = \{v : S \to \mathbb{R}\} \cong \mathbb{R}^S$

$$\Phi_\sigma : F_S \to F_S$$

$$\Phi_\sigma(v)(s) = \sum_{s', r'} \delta(s', r \mid s, \sigma(s))\left[r + \gamma v(s')\right]$$

Bellman's equations: $\quad \Phi_\sigma(val_\sigma) = val_\sigma$

Theorem: $\Phi_\sigma$ is $\gamma$-contracting.

$$\| \Phi_\sigma(v) - \Phi_\sigma(v') \|_\infty \leq \gamma \| v - v' \|_\infty$$

So it has a unique fixed point: $val_\sigma$

Algorithm:     $v_0(s) = 0$

$v_{k+1} = \Phi_\sigma(v_k)$

stop when $\|v_{k+1} - v_k\| \leq \varepsilon$

Lemma: $v_k = \mathbb{E}_\sigma \left[ \sum_{i=0}^{k-1} \gamma^i r_i \right]$

# STRATEGY IMPROVEMENT $\equiv$ ITERATION

$\sigma \rightsquigarrow val_\sigma \rightsquigarrow \sigma'$

$\sigma'(s) = \underset{a \in A}{\arg\max} \sum_{s', r} \delta(s', r | s, a) \left[ r + \gamma \, val_\sigma(s') \right]$

Theorem

if $\sigma$ is not optimal then $val_\sigma < val_{\sigma'}$

Strategy improvement algorithm

$\sigma \underset{eval}{\rightsquigarrow} val_\sigma \underset{impr}{\rightsquigarrow} \sigma_1 \underset{eval}{\rightsquigarrow} val_{\sigma_1} \underset{impr}{\rightsquigarrow} \sigma_2 \cdots\cdots \sigma_b$ optimal

# OFF-POLICY BELLMAN EQUATIONS

$$Val_*(s) = \max_{a \in A} \sum_{s', r} \delta(s', r | s, a) \left[ r + \gamma Val_*(s') \right]$$

$$① \ (v)(s) = \max_{a \in A} \sum_{s', r'} \delta(s', r | s, \sigma(s)) \left[ r + \gamma v(s') \right]$$

**Theorem :** ① is $\gamma \cdot$ contracting :

$$\| ①(v) - ①(v') \|_\infty < \gamma \| v - v' \|_\infty$$

So it has a unique fixed point : $val_*$

Value iteration algorithm :

$$v_0 \rightsquigarrow v_1 = ①(v_0) \rightsquigarrow v_2 = ①(v_1) \rightsquigarrow \ldots$$

$$until \quad \| v_{k+1} - v_k \|_\infty \leq \varepsilon$$