

Policy Iteration

Init $\pi_0: S \rightarrow A$.

$$Tf(s) = \underset{a}{\operatorname{argmax}} f(s, a)$$

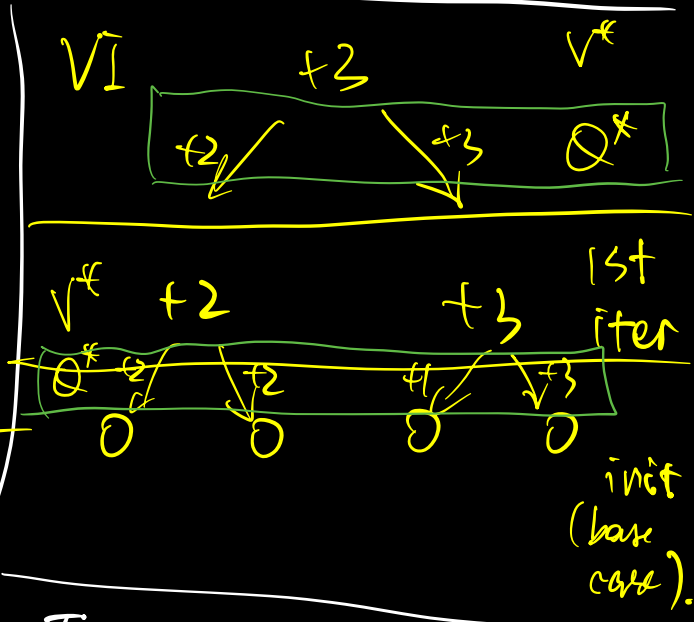
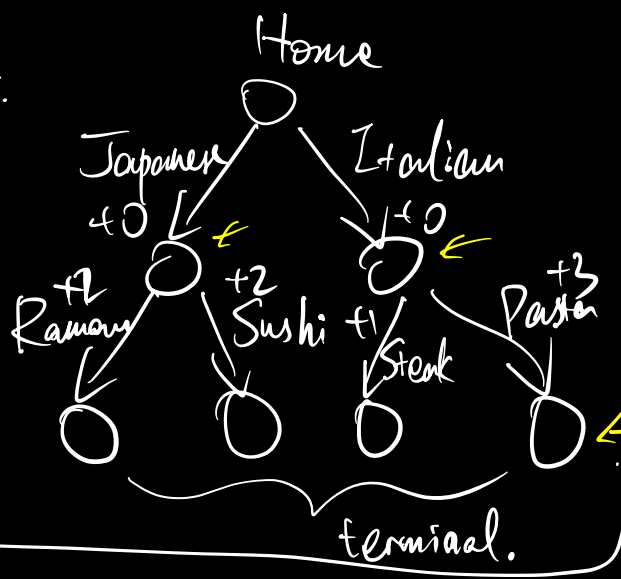
For $k=1, 2, 3, \dots$

Policy evaluation: Compute $Q^{\pi_{k-1}} \in \mathbb{R}^{S \times A}$

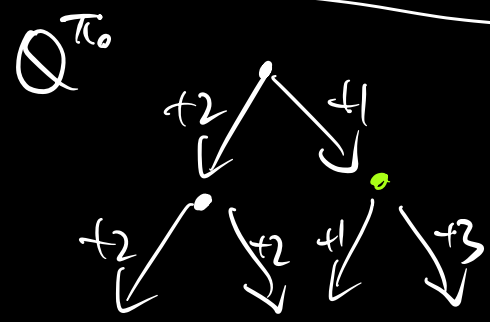
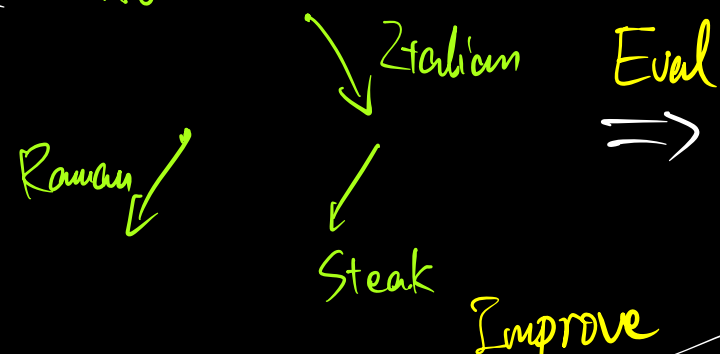
Policy Improvement: $\pi_k \leftarrow \pi_{Q^{\pi_{k-1}}}$

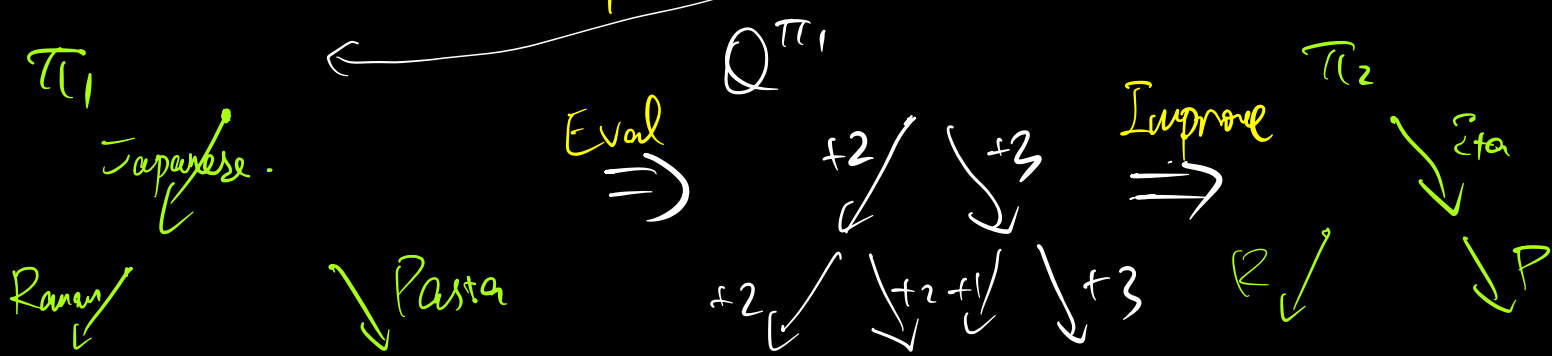
$$\pi_{k-1} = \pi^* \quad Q^{\pi_{k-1}} = Q^* \quad \pi_k = \pi_{Q^*} = \pi^*$$

Example:



PI: π_0 :





$$Q_{t+1} \quad \gamma = 0.9$$

$$f_0 = 0, \quad f_1 = 1, \quad f_2 = 1 + \gamma \cdot 1 = 1.9$$

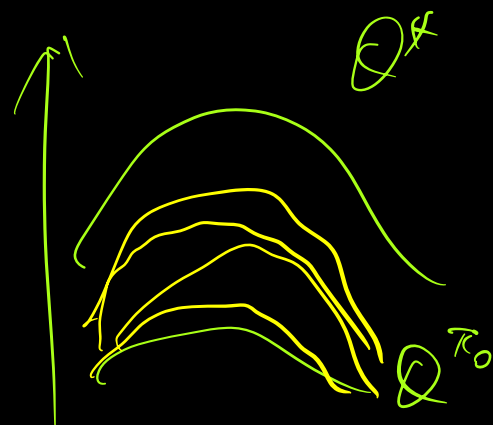
$$f_3 = 1 + \gamma + \gamma^2, \quad \dots$$

Monotone Policy Improvement: $\forall k$.

$$V^{\pi_k}(s) \geq V^{\pi_{k-1}}(s) \quad \forall s$$

if $\pi_{k-1} \neq \pi^*$, $\exists s. V^{\pi_k}(s) > V^{\pi_{k-1}}(s)$ δ

$$\Rightarrow \# \text{iter} \leq |A|^{|S|}$$



Lemma: T^π is monotone. that is, SXA

$$\forall f \leq f' \quad T^\pi f \leq T^\pi f'$$

$$\begin{aligned} & (Tf)(s,a) - (Tf')(s,a) \\ &= (R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [f(s',\pi)]) \\ &\quad - (R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [f'(s',\pi)]) \\ &= \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [f(s',\pi) - f'(s',\pi)] \\ &\leq 0. \end{aligned}$$

Want to show: $Q^{\pi_k} \leq Q^{\pi_{k+1}}$ alg.

$$Q^{\pi_k} = T^{\pi_k} Q^{\pi_k} \leq T Q^{\pi_k} = T^{\pi_{k+1}} Q^{\pi_k}$$

$$(Tf)(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [\max_{a'} f(s',a')]$$

$$(T^\pi f)(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [f(s',\pi)]$$

$$\forall f, Tf = T^\pi f.$$

$$Q^{\pi_k} \subseteq \mathcal{T}^{\pi_{k+1}} Q^{\pi_k}$$

$$\mathcal{T}^{\pi_{k+1}} Q^{\pi_k} \subseteq \mathcal{T}^{\pi_{k+1}} (\mathcal{T}^{\pi_{k+1}} Q^{\pi_k})$$

$$Q^{\pi_k} \subseteq \underbrace{\left(\mathcal{T}^{\pi_{k+1}} \right)^\infty}_{\text{"VI"}} Q^{\pi_k} = Q^{\pi_{k+1}}.$$

Recall: VI for Q^π : $\forall f_0 \in \mathbb{R}^{S \times A}$

$$f_k \leftarrow \mathcal{T}^\pi f_{k-1}.$$

$$f_k = (\mathcal{T}^\pi)^k f_0$$

as $k \rightarrow \infty$, $f_k \rightarrow Q^\pi$.