

MDP with finite state & action:

setting: $a \in A, s \in S, P_{st}^a$ ($\sum_t P_{st}^a = 1 \forall s, a$)
 (P^a is a transition matrix)

$\gamma \in (0, 1), r_s^a, \pi_s^a$ ($\sum_a \pi_s^a = 1 \forall s$)

Define $V_s^\pi = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_{s_t}^a \mid S_0 = s \right]$
 $a_0 \sim \pi_s^a$
 $S_{t+1} \sim P_{S_t, \cdot}^{a_t}$

$$V_s^* = \max_{\pi} V_s^\pi \quad \& \quad V^* = \begin{pmatrix} V_1^* \\ \vdots \\ V_s^* \end{pmatrix}$$

Bellman Eq: $V_s^\pi = r_s^\pi + \gamma \sum_t P_{st}^\pi V_t^\pi$ $r_s^\pi = \sum_a r_s^a \pi_s^a, P_{st}^\pi = \sum_a P_{st}^a \pi_s^a$

pf: $V_s^\pi = \sum_a r_s^a \pi_s^a + \gamma \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_{s_{t+1}}^a \mid S_0 = s \right]$
 $= \sum_a \mathbb{E} [V_{s_1} \mid S_0 = s, a_0 = a] \pi_s^a = \sum_{s_1, a} P_{s_1, a}^a V_{s_1}^\pi \pi_s^a$

optimal Bellman Eq: $V_s^* = \max_a \left(r_s^a + \gamma \sum_t P_{st}^a V_t^* \right)$
 \Leftrightarrow (1)

(P) $\begin{cases} \min_{V \in \mathbb{R}^{|S|}} \sum P_s V_s \\ \text{s.t.} \quad r_s^a + \gamma \sum_t P_{st}^a V_t - V_s \leq 0 \quad \forall a, s. \end{cases}$

(D) $\begin{cases} \max_{\lambda \in \mathbb{R}^{|A| \times |S|}} \sum_a (\lambda^a)^T r^a \\ \text{s.t.} \quad \lambda_s^a \geq 0 \quad \forall a, s \\ \sum_a (I - \gamma P^a)^T \lambda^a = e \end{cases} \rightarrow$ $|S| \times |A|$ -dim max prob
 with $|A| \times |S| + |S|$ constraints.

$$e^T V^* = \max_{\pi \in \Delta^{|S|}} e^T V^\pi$$

$$\text{s.t. } \begin{cases} V^\pi = r^\pi + \gamma P^\pi V \\ \sum_a \pi_s^a = 1, \pi_s^a \geq 0 \end{cases} \quad (P_{st}^\pi = \sum_a P_{st}^a \pi_s^a)$$

$$(P) V^* = \min e^T V$$

$|S|$ -dimensional min prob.

$$\begin{cases} r^{a_1} + \gamma P^{a_1} V - V \leq 0 \\ \vdots \\ r^{a_m} + \gamma P^{a_m} V - V \leq 0 \end{cases}$$

$|A| \times |S|$ ineq constraints.

MDP with 2 actions, 2 states $A = \{1, 2\}$, $S = \{1, 2\}$

$$P^1 = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad P^2 = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

$$r^1 = \begin{pmatrix} 1 \\ \frac{1}{2} \end{pmatrix}, \quad r^2 = \begin{pmatrix} \frac{3}{4} \\ \frac{1}{4} \end{pmatrix}, \quad \gamma = \frac{1}{2}$$

(P) :

$$\min \frac{1}{2}(V_1 + V_2)$$

$$\text{s.t. } 1 + \frac{1}{2} V_2 \leq V_1$$

$$\frac{1}{2} + \frac{1}{2} V_1 \leq V_2$$

$$\frac{3}{4} + \frac{1}{2} \left(\frac{1}{2} V_1 + \frac{1}{2} V_2 \right) \leq V_1$$

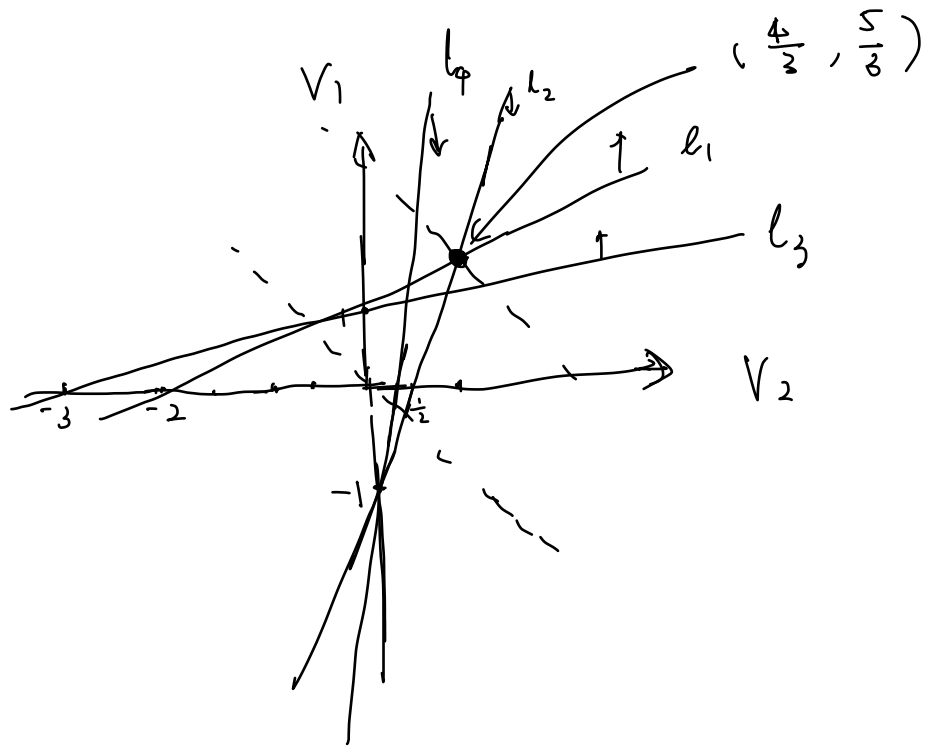
$$\frac{1}{4} + \frac{1}{2} \left(\frac{1}{2} V_1 + \frac{1}{2} V_2 \right) \leq V_2$$

$$V_1 \geq \frac{1}{2} V_2 + 1 \leftarrow \ell_1$$

$$V_1 \leq 2V_2 - 1 \leftarrow \ell_2$$

$$V_1 \geq \frac{1}{3} V_2 + 1 \leftarrow \ell_3$$

$$V_1 \leq 3V_2 - 1$$



$$(D) : V^* = \max \quad \lambda_1^1 + \frac{1}{2} \lambda_2^1 + \frac{3}{4} \lambda_1^2 + \frac{1}{4} \lambda_2^2$$

$$\text{s.t.} \quad \lambda_1^1, \lambda_2^1, \lambda_1^2, \lambda_2^2 \geq 0$$

$$a - \lambda_1^1 + \lambda_1^2 - \frac{1}{2} (\lambda_2^1 + \frac{1}{2} \lambda_1^1 + \frac{1}{2} \lambda_2^2) = \frac{1}{2}$$

$$\lambda_2^1 + \lambda_2^2 - \frac{1}{2} (\lambda_1^1 + \frac{1}{2} \lambda_1^1 + \frac{1}{2} \lambda_2^2) = \frac{1}{2}$$

$\begin{matrix} | & | \\ b & d \end{matrix}$

$$\Leftrightarrow \max \quad 2a + \frac{3}{2}c - \frac{1}{2} = k$$

$$\text{s.t.} \quad a, b, c \geq 0$$

$$4a + 3c - 2b \geq 2$$

$$5a + 4c - b = 4$$

$$\Rightarrow a = b = 1, \quad c = d = 0.$$

$$\Rightarrow \begin{cases} \lambda_1^1 = 1 & \lambda_1^2 = 0 \\ \lambda_2^1 = 1 & \lambda_2^2 = 0 \end{cases} \Rightarrow \text{always take action 1 for both states}$$

$$\Rightarrow \left\{ \begin{array}{l} v_1^* = \frac{5}{3}, \quad v_2^* = \frac{4}{3} \\ \text{at both states, optimal action are 1} \end{array} \right.$$

$$\textcircled{3} \quad \pi_s^a = \frac{1}{w_s} \lambda_s^a \leftarrow \text{normalized } \lambda, \quad w_s = \sum_a \lambda_s^a$$

$$1) \text{ then } \lambda_s^a \geq 0 \quad \forall a, s \quad (\Leftrightarrow) \quad \left\{ \begin{array}{l} \pi_s^a \geq 0 \quad \forall a, s \quad (\Leftrightarrow) \pi \in \Delta^{S^1} \\ \sum_a \pi_s^a = 1 \end{array} \right.$$

$$2) \quad \sum_a (I - \gamma P^a)^T \lambda^a = e \quad (\Leftrightarrow) \quad w = (I - \gamma P^\pi)^{-T} e$$

$$3) \quad \sum_a (\lambda^a)^T r^a \quad (\Leftrightarrow) \quad w^T r^\pi$$

$$4) \quad \max_{\pi \in \Delta^{S^1}} e^T (I - \gamma P^\pi)^{-1} r^\pi \quad (\Leftrightarrow) \quad \max_{\pi \in \Delta^{S^1}} e^T v^\pi$$

s.t. $v^\pi = (I - \gamma P^\pi)^{-1} r^\pi$
 $v^\pi = \textcircled{II} r^\pi - \gamma P^\pi v^\pi$

$$\textcircled{2} \quad \text{Lagrangian } L(v, \lambda) = \sum_s e_s v_s + \sum_{a, s} \lambda_s^a (r_s^a + \gamma \sum_t P_{st}^a v_t - v_s)$$

$$\min_v L(v, \lambda) = e^T v + \sum_a (r^a + \gamma P^a v - v)^T \lambda^a = f(\lambda)$$

$$\text{Dual: } \max_{\lambda_s^a} f(\lambda)$$

$$\text{s.t. } \lambda_s^a \geq 0$$

$$\min_v L(v, \lambda) = \underbrace{e + \gamma \sum_a (P^a - I)^T \lambda^a}_{g(\lambda)}^T v + \underbrace{\sum_a r^a \lambda^a}_{h(\lambda)} \quad \left\{ \begin{array}{l} -\infty \text{ if } g(\lambda) \neq 0 \\ h(\lambda) \text{ if } g(\lambda) = 0 \end{array} \right.$$

$$\rightarrow w^* = \max_{\lambda_s^a} \sum_a r^a \lambda^a \quad \sum_a \Gamma^a \lambda^a$$

$$\text{s.t. } \sum_a (I - \gamma P^a)^T \lambda^a = e$$

"Strong duality": If the obj fun is a convex & the constraints are linear

$$e^T v^* = f(x^*)$$

" \Rightarrow "

The solution to the optimal BZ solves the (P) optimization prob.

$$V_s^* = \max_{a \in A} (r_s^a + \gamma \sum_t P_{st}^a V_t^*)$$

this means for $\forall s, \exists a_s^*$ s.t. $V_s^* = r_s^{a_s^*} + \gamma \sum_t P_{st}^{a_s^*} V_t^*$
 $\& \forall a \neq a_s^* \Rightarrow V_s^* \geq r_s^a + \gamma \sum_t P_{st}^a V_t^*$

① feasibility: $\#$

② For \forall feasible V , $V_s - V_s^* \geq 0$ for $\forall s$

pf: subtracting $V_s^* = r_s^{a_s^*} + \gamma \sum_t P_{st}^{a_s^*} V_t^*$ from the V_s

$$V_s^* + \gamma \sum_t P_{st}^{a_s^*} V_t \leq V_s \quad \gamma \sum_t P_{st}^{a_s^*} (V_t - V_t^*) \leq V_s - V_s^*$$

$$r_s^{a_s^*} + \gamma \sum_t P_{st}^{a_s^*} V_t^* = V_s^*$$

$$\gamma P^* (V - V^*) \leq V - V^*, \quad P^* = \begin{pmatrix} P_{11}^{a_1^*} \\ \vdots \\ P_{|S|, |S|}^{a_{|S|}^*} \end{pmatrix}$$

Claim 1: If $(P - I)b \leq 0$ with $0 < \gamma < 1$ & $\sum_t P_{st} = 1 \Rightarrow b \geq 0$

$$V - V^* \leq 0 \Rightarrow e^T V^* = \min_{\forall \text{ feasible } V} e^T V$$

" \Leftarrow " The minimizer of (P) solves the optimal BZ.

① KKT condition: The minimizer v^* of (P) & the maximizer λ^* satisfies.

$$\lambda_i^* g_i(v^*) = 0 \quad \text{for } \forall i$$

(g_i 's are the constraints in (P))

$$\Rightarrow (\lambda_s^a)^* (r_s^a - \gamma \sum_t P_{st}^a V_t^* - V_s^*) = 0 \quad \forall a, s \quad (*)$$

② For fixed s , we want to show $\exists \lambda_s^*$ s.t. $r_s^a - \gamma \sum_t P_{st}^a V_t^* - V_s^* = 0$

(*) implies: $\lambda_s^a = 0$ or $r_s^a - \gamma \sum_t P_{st}^a V_t^* - V_s^* = 0$

proof by contradiction: if $\forall a, (\lambda_s^a)^* = 0$ & $r_s^a - \gamma \sum_t P_{st}^a V_t^* - V_s^* \neq 0$

then $\sum_a (I - \gamma P^a)^T \lambda^a = e \Rightarrow$ for this $s \Rightarrow \sum_a \lambda_s^a - \gamma \sum_t P_{ts}^a \lambda_t^a = e_s$

$\Rightarrow -\gamma \sum_t P_{ts}^a \lambda_t^a = e_s \Rightarrow e_s \leq 0 \quad \otimes$

proof of claim 1: pf by contradiction: if $\exists b_s < 0$

$$\sum_t P_{st} b_t \leq b_s$$

$$\geq \gamma \sum_{\{b_t \geq 0\}} P_{st} b_t - \gamma \sum_{\{b_t < 0\}} P_{st} |b_t|$$

because $\sum_t P_{st} = 1$

$$\geq \gamma \sum_{\{b_t \geq 0\}} P_{st} b_t - \gamma \max_{\{b_t < 0\}} |b_t|$$

$$\Rightarrow \underbrace{\gamma \sum_{\{b_t \geq 0\}} P_{st} b_t}_{\geq 0} \leq b_s + \gamma \max_{\{b_t < 0\}} |b_t| = b_s - \gamma \max_{\{b_t < 0\}} b_t = (1-\gamma)b_s < 0$$

for $s = \arg \max_{\{b_t < 0\}} b_t \quad \otimes$

Regularized reward.

$$\hat{r}_s^\pi = \mathbb{E}_{a \sim \pi} [r_s^a - \lambda \log \pi_s^a] = r_s^\pi - \sum_a \lambda \pi_s^a \log \pi_s^a$$

$$\hat{V}_s^\pi = \hat{r}_s^\pi + \gamma \sum_t P_{st}^\pi V_t^\pi$$

policy ∇ :

$$\max_{\pi \in \Delta^{|S|}} e^T \hat{V}_s^\pi$$

s.t. $\hat{V}^\pi = \hat{r}^\pi + \gamma P^\pi \hat{V}^\pi$

\Leftrightarrow ③

$$(D) : \quad \max_{\mu \in \mathbb{R}^{(S+1) \times 1}} \sum_{a,s} \mu_s^a \left[r_s^a - \lambda \log \left(\frac{\mu_s^a}{\sum_a \mu_s^a} \right) \right] = \sum_a (u^a)^T v^a - \sum h(\mu_s)$$

$$s.t. \quad \mu \geq 0$$

$$\sum_a (I - \gamma P^a)^T \mu^a = e$$

$$h(\mu_s) = \sum_a \mu_s^a \log \frac{\mu_s^a}{\sum_a \mu_s^a}$$

$$\textcircled{2} \quad \pi_s^a = \frac{1}{w_s} \mu_s^a \quad w_s = \sum_a \mu_s^a$$

$$\cdot \mu_s^a \geq 0 \Leftrightarrow \pi_s^a \geq 0 \quad \& \quad \sum \pi_s^a = 1$$

$$\cdot \sum_a \mu_s^a - \gamma \sum_{t,a} P_{t,s}^a \mu_t^a = e_s \Leftrightarrow \sum_a w_s \pi_s^a - \gamma \sum_{t,a} P_{t,s}^a w_t \pi_t^a = e_s$$

$$\Leftrightarrow w_s - \gamma \sum_t P_{t,s}^T w_t = e_s$$

$$\Leftrightarrow (I - \gamma P^T) w = e \quad \Leftrightarrow w = (I - \gamma P^T)^{-1} e$$

$$\cdot \sum_a (u^a)^T (r^a - \lambda \log \mu^a) \Leftrightarrow \sum_{a,s} w_s \pi_s^a (r_s^a - \lambda \log \pi_s^a)$$

$$\Leftrightarrow \sum_s w_s \sum_a \pi_s^a (r_s^a - \lambda \log \pi_s^a)$$

$$\Leftrightarrow \sum_s w_s \hat{r}_s^T$$

$$\Leftrightarrow w^T \hat{r}^T$$

$$\Leftrightarrow e^T \underbrace{(I - \gamma P^T)^{-1}}_{V^T} \hat{r}^T$$

$$(D) \stackrel{\textcircled{2}}{\Leftrightarrow} (P) \quad \min e^T v$$

$$s.t. \quad \max_{\pi \in \Delta^{(S)}} (\hat{r}^T + \gamma P^T v - v) \leq 0$$

$$(D) \quad L(v, w) = e^T v + w^T \left(\max_{\pi \in \Delta^{(S)}} \hat{r}^T + \gamma P^T v - v - \lambda h(\pi) \right)$$

$$\min_v L(v, w) = e^T v + \sum_s w_s \left(\max_{\substack{\pi_s^a \geq 0 \\ \sum_a \pi_s^a = 1}} \sum_a \pi_s^a \left(\hat{r}_s^a + \gamma \sum_{t,a} P_{st}^a v_t - v_s - \lambda \log \pi_s^a \right) \right)$$

$$\text{for fixed } s : \quad \max_{\pi} \sum \pi_s^a \left(r_s^a + \gamma \sum_{t,a} P_{st}^a v_t - v_s - \lambda \log \pi_s^a \right)$$

$$\max_{\pi_s^a} \sum w_s \pi_s^a \left(r_s^a + \gamma \sum_{t,a} P_{st}^a v_t - v_s - \lambda \log \frac{w_s \pi_s^a}{w_s} \right)$$

$$\max_{\mu_s^a} \sum \mu_s^a \left(r_s^a + \gamma \sum_{t,a} P_{st}^a v_t - v_s - \lambda \log \frac{\mu_s^a}{\sum_a \mu_s^a} \right)$$

$$\max_{u_s \geq 0} \sum_a (u^a)^T (r^a + \gamma P^a V - V) - \lambda \sum_s u_s h(u_s)$$

$$\min_V \mathcal{L}(V, u) = \min_V \left(e^T V + \max_u [\dots] \right)$$

$$= \max_u \min_V \underbrace{e^T V + \dots}$$

$$\Rightarrow \max_u \sum_a (u^a)^T r^a - \lambda \sum_s u_s h(u_s)$$

$$e^T + \sum_a (u^a)^T (\gamma P^a - I) = 0$$

$$\Leftrightarrow (I - \gamma P^a)^T u^a = e$$

$$(p) \Leftrightarrow \max_{\pi \in \Delta^a} (\sum_a \pi^a + \gamma P^a V - V) \leq 0$$

$$\max_{\pi_s} \left[\sum_a r_s^a \pi_s^a + \gamma \sum_t P_{st}^a V_t \pi_s^a - V_s - \lambda \pi_s^a \log \pi_s^a \right]$$

$$f(\pi_s) = \sum_a (r_s^a + \gamma P_{st}^a V_t - V_s) \pi_s^a - \lambda \pi_s^a \log \pi_s^a$$

g_s^a Lemma: $\max_{\pi} \sum_a g^a \pi^a - \lambda \sum_a \pi^a \log \pi^a$
 \downarrow
 $\pi^a \propto e^{\frac{1}{\lambda} g^a}$
 $\lambda \log \left(\sum_a e^{\frac{1}{\lambda} g^a} \right)$

$$\frac{\partial f(\pi_s^a)}{\partial \pi_s^a} = g_s^a - \lambda \log \pi_s^a - \lambda = 0$$

$$\Rightarrow \frac{g_s^a - \lambda}{\lambda} = \log \pi_s^a$$

$$\Rightarrow \pi_s^a \propto e^{\frac{1}{\lambda} (g_s^a - \lambda)}$$

$$\Rightarrow \pi_s^a = \frac{e^{\frac{1}{\lambda} (g_s^a - \lambda)}}{\sum_a e^{\frac{1}{\lambda} (g_s^a - \lambda)}} = \frac{e^{\frac{1}{\lambda} g_s^a}}{\sum_a e^{\frac{1}{\lambda} g_s^a}} \underbrace{\qquad}_{\omega_s}$$

$$f(\pi_s^*) = \sum_a g_s^a \frac{e^{\frac{1}{\lambda} g_s^a}}{\omega_s} - \lambda \sum_a \frac{e^{\frac{1}{\lambda} g_s^a}}{\omega_s} \log \left(\frac{e^{\frac{1}{\lambda} g_s^a}}{\omega_s} \right)$$

$$\begin{aligned}
&= \frac{1}{w_s} \sum_a g_s^a e^{\frac{1}{\lambda} g_s^a} - \frac{\lambda}{w_s} \sum_a e^{\frac{1}{\lambda} g_s^a} \left(\frac{1}{\lambda} g_s^a - \log(w_s) \right) \\
&= \frac{\lambda}{w_s} \sum_a e^{\frac{1}{\lambda} g_s^a} \log(w_s) \\
&= \lambda \log w_s = \lambda \log \left(\sum_a e^{\frac{1}{\lambda} g_s^a} \right)
\end{aligned}$$

$$(p) \Leftrightarrow \min e^T v$$

$$\text{s.t. } \lambda \log \left(\sum_a e^{\frac{1}{\lambda} (r_s^a - \sum_t p_{st}^a v_t - v_s)} \right) \leq 0$$

$$(p) \Leftrightarrow v_s^* = \max_{\pi_s \in \Delta} \sum_a \left(r_s^a + \nu \sum_t p_{st}^a v_t^* - \lambda \log \pi_s^a \right) \pi_s^a$$

$$\Leftrightarrow v_s^* = \lambda \log \left(\sum_a e^{\frac{1}{\lambda} (r_s^a + \nu \sum_t p_{st}^a v_t^*)} \right)$$