# Linear Programming (LP) for MDPs.

Primal form.

$$\min_{V \in \mathbb{R}^S} d_0^T V$$

$\to$ "init state distribution"
$d_0 > 0, \sum_s d_0(s) = 1.$

$$\text{s.t.} \quad V \geq \mathcal{T} V$$

$\to$ Bellman optimality op.

① Why LP? Constraints: $\forall s \in S.$

$$V(s) \geq \boxed{\max_{a \in A}} \left( R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s,a)} [V(s')] \right).$$

Nonlinear!     Can convert to $|A|$ linear constraints:

$$\underline{V(s) \geq R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s,a)} [V(s')]}, \quad \forall a \in A.$$

linear!

∴ Primal form has $|S|$ decision varibles, & $|S| \times |A|$ constraints.

② Why it solves planning problem?

Claim: if $d_0(s) > 0, \forall s \in S.$ then opt. sol. to primal
is $V = V^*.$

Proof: constraint $V \geq \mathcal{T} V.$

---
Lemma. monotone property of $\mathcal{T}$. $\forall V \geq V'. \mathcal{T}V \geq \mathcal{T}V'.$

---

We will show: $\underline{V \geq \mathcal{T}V} \Rightarrow V \geq V^*.$

Invoke monotonicity on $V$ & $V' = \mathcal{T}V. \Rightarrow \underline{\mathcal{T}V \geq \mathcal{T}(\mathcal{T}V)}$

$$V \geq \mathcal{T}V \geq \mathcal{T}^2 V \Rightarrow V \geq \mathcal{T}^2 V.$$

Now let $V' = \mathcal{T}^2 V \Rightarrow V \geq \mathcal{T}V \geq \mathcal{T}^3 V \Rightarrow V \geq \mathcal{T}^3 V.$
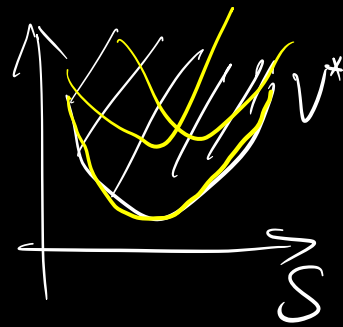
$\cdots \quad \forall n, \quad V \geq \mathcal{T}^n V.$

take limit $n \to \infty$. $V \geq T^\infty V = V^*$.

Therefore ① any feasible $V \geq V^*$.

② $V^*$ is feasible: $V^* = TV^*$

$\Rightarrow \underline{\min\limits_{V} d_0^T V \text{ is achieved w/ } V = V^*.}$

---

$\boxed{\text{Dual form}}$ $\quad \max\limits_{d \in \mathbb{R}^{S \times A}, \, d \geq 0} \quad d^T R \quad \longrightarrow$ reward function.

s.t. $\forall s \in S, \quad \sum\limits_a d(s,a) = \gamma \sum\limits_{\hat{s}, \tilde{a}} d(\hat{s}, \tilde{a}) P(s | \hat{s}, \tilde{a}) + (1-\gamma) d_0(s).$

Interpretation: $d$ plays the role of occupancy.

In fact. dual constraint characterizes all possible occupancies in the MDP.

if $d = d^\pi$. $\Rightarrow d^T R = (d^\pi)^T R = \mathbb{E}_{(s,a) \sim d^\pi} [R(s,a)]$.

$$= \mathbb{E}_{s \sim d_0} [V^\pi(s)].$$

Verify $d = d^\pi$ is feasible for any $\pi$:

$d^\pi(s,a) = (1-\gamma) \sum\limits_{t=1}^{\infty} \gamma^{t-1} d_t^\pi(s,a) \quad \longrightarrow$ distribution of $s_t, a_t$ under $\pi, d_0$.

$\underline{\underline{d^\pi(s)}} = (1-\gamma) \sum\limits_{t=1}^{\infty} \gamma^{t-1} \underline{\underline{d_t^\pi(s)}}$

Dual constraint: $\sum\limits_a d(s,a) = \gamma \sum\limits_{\hat{s}, \tilde{a}} d(\hat{s}, \tilde{a}) P(s | \hat{s}, \tilde{a}) + (1-\gamma) d_0(s)$

$$\text{LHS} = \sum_{a \in A} d^\pi(s,a) = d^\pi(s) = (1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} d_t^\pi(s).$$

$$\text{RHS (first term)} = \gamma(1-\gamma) \sum_{t=1}^{\infty} \sum_{\hat{s}, \hat{a}} \gamma^{t-1} \underbrace{d_t^\pi(\hat{s}, \hat{a})}_{\triangle} \underbrace{P(s|\hat{s}, \hat{a})}_{\triangle}$$

a function $s$.

a distribution of state.

$\hookrightarrow$ $s$ drawn from the dist $\Longleftrightarrow (\tilde{s}, \tilde{a}) \sim \underbrace{d_t^\pi}_{\triangle}, \; s \sim \underbrace{P(\cdot|\hat{s}, \hat{a})}_{\triangle}$

$\Rightarrow$ the distribution is $d_{t+1}^\pi(s)$.

$$\Rightarrow \text{RHS (first term)} = \gamma(1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} d_{t+1}^\pi(s)$$

$$= (1-\gamma) \sum_{t=1}^{\infty} \gamma^{t} d_{t+1}^\pi(s)$$

$$= (1-\gamma) \sum_{t=2}^{\infty} \gamma^{t-1} d_{t}^\pi(s) \quad \Longleftarrow$$

$$\text{RHS (second term)} = (1-\gamma) d_0(s) = (1-\gamma) d_1^\pi(s)$$
$$\forall \pi.$$



$$d^\pi(s,a) = d^\pi(s) \cdot \pi(a|s)$$

Remark: solution to dual is $d = \underbrace{d^{\pi^*}}_{\triangle} \xrightarrow{\text{back out}} \pi^*.$

$\forall \pi$ (stat. possibly stochastic). given $\underbrace{d^\pi(s,a)}_{\triangle} \Rightarrow \pi(a|s) = \dfrac{d^\pi(s,a)}{\sum_{a'} d^\pi(s,a')}$