



- The linear-program interpretation is even more powerful when we deal with constrained MDP
- We will restrict ourselves to the average cost problem.
- Suppose that in addition to a per-stage cost  $g(i,u)$ , there is a per-stage penalty of  $y(i,u)$ .
- We want to minimize the average cost

$$\lim_{t \rightarrow +\infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} g(X_k, u_k) \right]$$

subject to a constraint of the average penalty

$$\lim_{t \rightarrow +\infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} y(X_k, u_k) \right] \leq V$$

### Linear program

- Let  $\lambda_{iu} = \lambda_i \delta_{iu}$

= The probability of being at state  $i$  & use action  $u$

- The average penalty is

$$\sum_{i,u} \lambda_{iu} y(i,u)$$

- Thus, the constrained DP can be written as

$$\min \sum_{i,u} \lambda_{iu} g(i,u)$$

$$\text{sub to } \sum_n \lambda_{jn} = \sum_{in} \lambda_{in} p_{ij}(n) \quad \forall j$$

$$\sum_{in} \lambda_{in} y(i, n) \leq V \quad (*)$$

$$\sum_{in} \lambda_{in} = 1$$

- In general, the solution may have multiple non-zero  $\lambda_{in}$  for a state  $i$ ,

In other words, a probabilistic policy is needed.

## Duality

- Since the Linear program is also a convex program, duality holds
- Associate a Lagrange multiplier  $\lambda$  to (\*)
- The Lagrangian is

$$L(\vec{\lambda}, \lambda) = \sum_{in} \lambda_{in} g(i, n) + \lambda \sum_{in} \lambda_{in} y(i, n) - \lambda V$$

The dual objective is

$$D(\lambda) = \min L(\vec{\lambda}, \lambda)$$

$$\text{sub to } \sum_j \lambda_{jn} = \sum_{in} \lambda_{in} p_{ij}(n)$$

$$\sum_{in} \lambda_{in} = 1$$

- But this is simply an average-cost problem with per-stage cost

$$g(i, n) + \lambda y(i, n).$$

- In a more symbolic way, we can rewrite the constrained DP as

$$i \rightarrow i \stackrel{t-1}{=} 0 \text{ or } \dots$$

constrained DP as U:

$$\min \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} g(x_k, u_k) \right]$$

$$\text{Sub to } \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} y(x_k, u_k) \right] \leq V$$

Associate a Lagrange multiplier  $\lambda$  to the constraint.

$$L = \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} g(x_k, u_k) + \lambda y(x_k, u_k) \right] - \lambda V$$

- We should then minimize

$$\min \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} g(x_k, u_k) + \lambda y(x_k, u_k) \right]$$

- This is just an average-cost MDP!

---

We can then apply all results from duality.

- There exists a policy  $\mu$  (possibly probabilistic) &  $\lambda$  such that

$$\text{KKT} \left\{ \begin{array}{l} \mu \text{ minimizes } \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} g(x_k, u_k) + \lambda y(x_k, u_k) \right] \quad (\text{an MDP}) \\ \lambda \geq 0 \\ \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} y(x_k, u_k) \right] \leq V \\ \lambda \cdot \left[ \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} y(x_k, u_k) \right] - V \right] = 0 \end{array} \right.$$

- Any pairs of  $\mu, \lambda$  that satisfy that KKT condition are also optimal.

- The following iterative algo will converge for  $\lambda$

At step  $t$ :

- Solve the average cost MDP with stage cost

$$g(x_k, u_k) + \lambda y(x_k, u_k)$$

— Update

$$\lambda(t+1) = \left[ \lambda(t) + \alpha \left[ \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{k=0}^{T-1} y(x_k, u_k) \right] - v \right] \right]^+$$

# Decomposition

Friday, April 14, 2023 11:41 AM

- As in convex optimization, this type of duality can be very helpful for decomposing a large problem into smaller problems!

- Consider  $M$  copies of the MDP

- The  $m$ -th MDP has a per-stage cost  $g^m(x_k^m, u_k^m)$  and per-stage penalty  $y^m(x_k^m, u_k^m)$

- If the penalties are not coupled, i.e., each MDP has a separate constraint on

$$\lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} y^m(x_k^m, u_k^m) \right]$$

Then of course each MDP can be solved independently.

- What if the penalty constraints are coupled

$$\begin{aligned} \min \quad & \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{m=1}^M \sum_{k=0}^{t-1} g^m(x_k^m, u_k^m) \right] \\ \text{sub to} \quad & \lim_{t \rightarrow \infty} \frac{1}{t} \mathbb{E} \left[ \sum_{m=1}^M \sum_{k=0}^{t-1} y^m(x_k^m, u_k^m) \right] \leq V \quad (\star) \end{aligned}$$

- We can no longer solve each MDP separately!

- Further, solving this global MDP will likely run into the curse-of-dimensionality!

... instead, solving MDPs given a policy will likely run into the curse-of-dimensionality!

- Instead, associate a Lagrange multiplier  $\lambda$  to (\*).

- The Lagrangian (precise form can be written through the corresponding LP)

$$\begin{aligned} L(\vec{\lambda}, \lambda) &= \lim_{t \rightarrow +\infty} \frac{1}{t} \mathbb{E} \left[ \sum_{m=1}^M \sum_{k=0}^{t-1} g^m(x_k^m, u_k^m) \right] \\ &\quad + \lambda \left[ \lim_{t \rightarrow +\infty} \frac{1}{t} \mathbb{E} \left[ \sum_{m=1}^M \sum_{k=0}^{t-1} y^m(x_k^m, u_k^m) \right] - V \right] \\ &= \sum_{m=1}^M \left\{ \lim_{t \rightarrow +\infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} g^m(x_k^m, u_k^m) + \lambda y^m(x_k^m, u_k^m) \right] \right\} \\ &\quad - \lambda V \end{aligned}$$

- Therefore, the  $m$ -th MDP can optimize

$$\min_{\lambda_m} \lim_{t \rightarrow +\infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} g^m(x_k^m, u_k^m) + \lambda y^m(x_k^m, u_k^m) \right]$$

- This is a much smaller problem!

- Finally, the Lagrange multiplier can be updated by

$$\begin{aligned} \lambda^{(l+1)} &= \left[ \lambda^{(l)} \right. \\ &\quad \left. + \alpha \left( \sum_{m=1}^M \lim_{t \rightarrow +\infty} \frac{1}{t} \mathbb{E} \left[ \sum_{k=0}^{t-1} y^m(x_k^m, u_k^m) \right] - V \right) \right]^+ \end{aligned}$$

using the optimal policy based on  $\lambda^{(l)}$ .

using the optimal policy  
based on  $\lambda(V)$ .

can be replaced by

$$E[y^m(x_k^m, u_k^m)]$$

↑  
assume  $k$  is the  
steady-state