

Review:

Shortest Path LP's

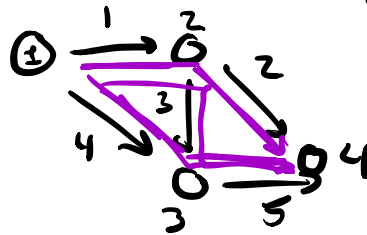
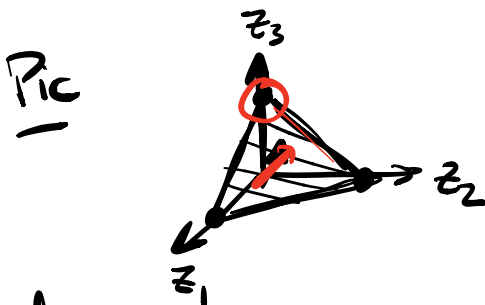
1) Enumerate Paths \mathcal{R} (set of routes)

$R \in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{R}|}$ $r \in \mathcal{R}$ particular route

Indicator matrix for which edges are in ea route

$$\begin{cases} \min_{z \in \mathbb{R}^{|\mathcal{R}|}} l^T z \\ \text{s.t. } \mathbf{1}^T z = 1, z \geq 0 \end{cases} \quad l^T = C^T R \quad l_r \text{ cost of taking route } r$$

$$l_r = \sum_{e \in \mathcal{E}} c_e$$

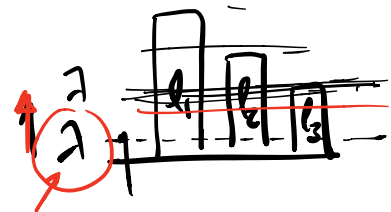


Dual

$$\begin{cases} \max_{\lambda, u \in \mathbb{R}^{|\mathcal{R}|}} \lambda \\ \text{s.t. } \lambda \mathbf{1}^T = l^T - u^T, u \geq 0 \end{cases}$$

$\Rightarrow \lambda \mathbf{1}^T \leq l^T \quad \lambda \mathbf{1}^T = l^T - u^T$

$$\lambda = l_r - u_r, u_r \geq 0 \iff \lambda \leq l_r$$



2) Edge Formulation

$$x^* = R z^*$$

$$\min_{x \in \mathbb{R}^{|E|}} C^T x$$

$$\text{s.t. } E x = b \quad x \geq 0$$

incidence matrix

source sink vector

Dual

$$\begin{aligned} & \rightarrow \max_{v \in \mathbb{R}^{|I|}} -v^T b \\ & \rightarrow \mu \in \mathbb{R}^{|E|} \end{aligned}$$

$$\text{s.t. } -v^T E = C^T - \mu^T, \mu \geq 0$$

$$v^T E \leq C^T$$

$$b = \begin{bmatrix} -1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

← origin

← dest

$$-v^T b = v_0 - v_d$$

cost to go at origin cost to go at dest.

v : value function
"cost-to-go" from ea. node to destination


μ_e : inefficiency of taking edge e

$$E = E_i - E_0$$

$$v^T E = v^T E_i - v^T E_0$$

$$v^T E = C^T - \mu^T$$

for ea. edge $e \dots$



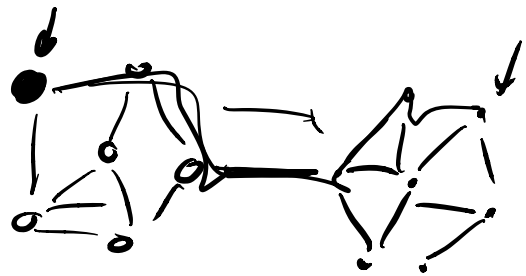
$$-(v_j - v_i) = c_e - \mu_e$$

$$\rightarrow v_i = \underline{c_e} - \underline{\mu_e} + v_j \Rightarrow \underline{v_i \leq c_e + v_j}$$

$$\left(\begin{array}{l} \min_x \max_{v, \mu \geq 0} c^T x + v^T (Ex - b) - \mu^T x \quad \left\{ \begin{array}{l} x \geq 0 \\ x < 0 \end{array} \right. \leftarrow -v^T b \\ \max_v \min_x c^T + v^T E - \mu^T = 0 \end{array} \right)$$

v_i should be a lower bound on the cost to go at node i

$Ex = b$
 ↪ always an extra eqn.



$$\mathbb{1}^T (Ex = b)$$

$$\mathbb{1}^T E x = \mathbb{1}^T b$$

$$0x = 0$$

$[E|b]$ not full row rank ← extra row.

Fix problem:

$$\underline{u}^T \underline{u} (Ex = b)$$

$$u = \begin{bmatrix} \mathbb{1} & 0 & 0 \\ 0 & \mathbb{I} & 0 \\ \mathbb{1} & \dots & \mathbb{1} \end{bmatrix} = \begin{bmatrix} \mathbb{I} & 0 \\ \mathbb{1}^T & \mathbb{1} \end{bmatrix} \quad u^{-1} = \begin{bmatrix} \mathbb{I} & 0 \\ -\mathbb{1}^T & \mathbb{1} \end{bmatrix}$$

$$\underline{u}^T (\underline{u} Ex = \underline{u} b)$$

$$\begin{bmatrix} I & 0 \\ \underline{1}^T & 1 \end{bmatrix} E x = \begin{bmatrix} I & 0 \\ \underline{1}^T & 1 \end{bmatrix} b$$

$$\begin{bmatrix} [I \ 0] E \\ \underline{1}^T E \end{bmatrix} x = \begin{bmatrix} [I \ 0] E \\ 0 \end{bmatrix} x = \begin{bmatrix} [I \ 0] b \\ \underline{1}^T b \end{bmatrix} = \begin{bmatrix} [I \ 0] b \\ 0 \end{bmatrix}$$

isolated the redundant constraint

$$\rightarrow \underbrace{[I \ 0] E}_{\substack{\text{all but} \\ \text{last row of } E}} x = \underbrace{[I \ 0] b}_{\substack{\text{all but} \\ \text{last row} \\ \text{of } b}}$$

$$v^T (E x - b)$$

$$v^T u^T (u E x - u b)$$

$$v^T u^T \left(\begin{bmatrix} [I \ 0] E \\ 0 \end{bmatrix} x - \begin{bmatrix} [I \ 0] b \\ 0 \end{bmatrix} \right)$$

$$v^T \begin{bmatrix} [I \ 0] \\ -\underline{1}^T \ 1 \end{bmatrix} \left(\begin{bmatrix} [I \ 0] E \\ 0 \end{bmatrix} x - \begin{bmatrix} [I \ 0] b \\ 0 \end{bmatrix} \right) \leftarrow$$

performing a coord transform on dual variable v



$$v^{iT} = v^T \begin{bmatrix} [I \ 0] \\ -\underline{1}^T \ 1 \end{bmatrix} =$$

$$\underline{v^{iT}} = \underline{[v_0 \ -v_s \ -v_d]} \begin{bmatrix} [I \ 0] \\ -\underline{1}^T \ 1 \end{bmatrix} = \begin{bmatrix} v_0 - v_d & v_s - v_d & v_d \end{bmatrix}$$

$$\underline{[v_0 - v_d \ \dots \ v_s - v_d \ \dots \ v_d]} \left(\begin{bmatrix} [I \ 0] E \\ 0 \end{bmatrix} x - \begin{bmatrix} [I \ 0] b \\ 0 \end{bmatrix} \right)$$

$$\underline{[v_0 - v_d \ \dots \ v_{k-1} - v_d]} \left([I \ 0] E x - [I \ 0] b \right)$$

$$\max -v^T b$$

$$\text{s.t. } -v^T E = c^T - \mu^T, \mu \geq 0$$

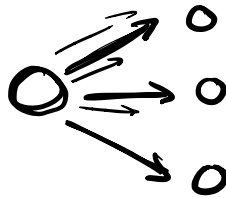
$$\left(\begin{array}{l} \max -v^T u b \Rightarrow v'_0 \rightarrow \text{maximize the cost-to-go from the origin} \\ \text{s.t. } -v'^T u E = c^T - \mu^T, \mu \geq 0 \end{array} \right.$$

MARKOV DECISIONS PROCESSES (MDP)

flow problems w/ stochastic transitions

$$G = (S, \mathcal{E}) \leftarrow$$

Action s:



→ A_s : actions at state s

for $a \in A_s \Rightarrow \text{Prob}(s'|s,a) = \text{Prob}(s'|a)$ Transition kernel

$$A = \bigcup_s A_s$$

↳ all actions

ea. action $a \in A$ implies being in a particular state s

→ actually represents a state action pair

- ea. action is only available from 1 state

Incidence Matrices

• $A \in \mathbb{R}^{|S| \times |A|}$

$$A_{sa} = \begin{cases} 1 & \text{if } a \in A_s \\ 0 & \text{otherwise} \end{cases}$$

indicator matrix for which actions are available from which states

• $P \in \mathbb{R}^{|S| \times |A|}$

Transitional kernel matrix

$P_{s'a} = \text{Prob}(s'|a)$

• $W \in \mathbb{R}^{|E| \times |A|}$

Trans kernel ...

if edge e
 $s \rightarrow s'$

$\text{Prob}(e|a) = \text{Prob}(s'|a,s)$

$W_{ea} = \text{Prob}(e|a,s)$
 $= \text{Prob}(e|a)$

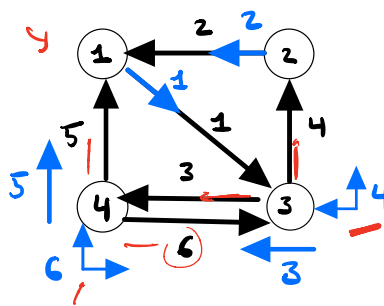
Ex. $E_i = \begin{bmatrix} 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$

$E_o = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$

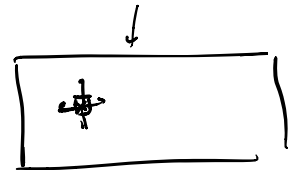
$A = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}$

$P = \begin{bmatrix} 0 & 1 & 0 & 0 & 1 & 0.5 \\ 0 & 0 & 0 & 0.5 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0.5 \\ 0 & 0 & 1 & 0.5 & 0 & 0 \end{bmatrix}$

$W = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0.5 & 0 & 0 \\ 0 & 0 & 0 & 0.5 & 1 & 0.5 \\ 0 & 0 & 0 & 0 & 0 & 0.5 \end{bmatrix}$



A: actions
E: edges



transition information

Relationships

$$|S| \int_{|S|} P = \int_{|S|} E_i |I| \int_{|I|} W \quad A = E_0 W$$

Before:

non stochastic mass conservation

$$\rightarrow \underline{x} \in \mathbb{R}^{|I|} \rightarrow \boxed{E_i x = E_0 x} \quad E x = 0$$

Now: $y \in \mathbb{R}^{|A|}$ stochastic flow vector

y_a : mass taking action a over the actions

$\varphi \in \mathbb{R}^{|S|}$ probability distribution over states

φ_s : mass in state s

for a particular distribution y over actions we can compute the corresponding state dist

$$\rightarrow \underline{\varphi} = \underline{A} y \leftarrow \text{summing probabilities on all actions in ea state}$$

steady stochastic flow.

$$\boxed{\underline{P} y = \underline{\varphi} = \underline{A} y} \leftarrow$$

where \underline{P} is mass distribution over states that mass transitions to

$x = Wy$
 \downarrow corresponding dist on edges
 \downarrow dist on actions

$$Py = Ay \quad P = E_i W \quad A = E_o W$$

$$\rightarrow \underline{E_i W y} = \underline{E_o W y} \Rightarrow \underline{(E_i - E_o) W y} = 0$$

\underline{x}

Mass conservation:

everything here is "column stochastic"
 entries in ea. column sum to 1

$$\mathbb{1}_{|S|}^T E_i = \mathbb{1}_{|S|}^T E_o = \mathbb{1}_{|E|}^T$$

$$\mathbb{1}_{|S|}^T P = \mathbb{1}_{|S|}^T A = \mathbb{1}_{|A|}^T$$

$$\mathbb{1}_{|E|}^T W = \mathbb{1}_{|A|}^T$$

$$\underline{\mathbb{1}_{|A|}^T y} = 1, \quad \mathbb{1}_{|S|}^T \varphi = 1, \quad \mathbb{1}_{|E|}^T x = 1$$

$$\varphi = Ay \quad \underline{\mathbb{1}_{|S|}^T \varphi} = \mathbb{1}_{|S|}^T A y = \underline{\mathbb{1}_{|A|}^T y} = 1$$

Rewards $r \in \mathbb{R}^{|A|}$ r_a : reward for taking action a

FINITE HORIZON MDP LP: $x^T Q x + u^T R u$

PRIMAL
max

$\rightarrow y(t)$
 $t=0, \dots, T$

$\sum_{t=0}^{T-1} r(t)^T y(t) + g^T A y(T)$

$\rightarrow x^+ = Ax + Bu$

$y(t)$: mass dist over actions at time t

$r(t)$: rewards at time t

g_s : final reward at states

s.t. $A y(0) = \varphi(0)$

initial state distribution (given)

MASS
CONS.
OVER
time

$A y(t+1) = P y(t)$
 $y(t) \geq 0$

$t=0, \dots, T-1$
"
 $\varphi(t+1) = P y(t)$

Dual variables

$A y(0) = \varphi(0) \Rightarrow v(0) \in \mathbb{R}^{|S|}$

$A y(t+1) = P y(t) \Rightarrow v(t+1) \in \mathbb{R}^{|S|}$

$y(t) \geq 0 \Rightarrow \mu(t) \in \mathbb{R}_+^{|A|}$

Lagrangian

$\mathcal{L}(y, v, \mu) = \sum_{t=0}^{T-1} r(t)^T y(t) + g^T A y(T)$

$- v(0)^T (A y(0) - \varphi(0)) - \sum_{t=0}^{T-1} v(t+1)^T (A y(t+1) - P y(t)) + \sum_{t=0}^T \mu(t)^T y(t)$

$$\max_y \min_{v, \mu \geq 0} \mathcal{L} \leq \min_{v, \mu \geq 0} \max_y \mathcal{L}$$

$$\frac{\partial \mathcal{L}}{\partial y} = 0$$

$$\frac{\partial \mathcal{L}}{\partial y(t)} = r(t)^T + v(t+1)^T P - v(t)^T A + \mu(t)^T = 0$$

$t = 0, \dots, T-1$

$$\frac{\partial \mathcal{L}}{\partial y(T)} = g^T A - v(T)^T A + \mu(T)^T = 0$$

Dual.

$$\min_{v, \mu} \underline{v(0)^T \varphi(0)} \rightarrow \text{"expected total reward if our initial mass dist is } \varphi(0)\text{"}$$

$$\text{s.t. } v(T)^T A = g^T A + \mu(T)^T, \mu(T) \geq 0$$

$$\rightarrow \underline{v(t)^T A} = \underline{r(t)^T + v(t+1)^T P + \mu(t)^T} \quad t=0, \dots, T-1$$

$\mu(t) \geq 0$

↳ Bellman Egn.

element wise

$$V_s(t) = r_a(t) + v(t+1)^T P[:, a] + M_a(t)$$

value at next time
" "

$$\Rightarrow V_s(t) = r_a(t) + \underbrace{v(t+1)^T P[:, a]}_{\substack{\text{reward for} \\ \text{taking action} \\ a}} + \underbrace{\mu_a(t)}_{\substack{\text{inefficiency} \\ \text{of action} \\ a}}$$

\downarrow
 value function
 at state s
 time t

\downarrow
 reward for
 taking action
 a

\downarrow
 $q_a(t+1) =$
 "expected
 reward to go
 for taking
 action a "

\downarrow
 inefficiency
 of action
 a

$$v(t+1)^T P := q(t+1)^T$$

$q(t+1) \in \mathbb{R}^{|A|}$

\rightarrow "Q-value from
 Q learning"

$$\underline{V_s(t)} \geq \underline{r_a(t)} + v(t+1)^T P[:, a]$$

$$V_s(T) = g_s + \mu_a(t)$$

$$\rightarrow V_s(T) \geq g_s \leftarrow$$

$$v^T A = q^T$$

RL / Q learning

$$\underline{v^T A} = r^T + \underline{q^T} + \mu^T$$



Infinite Horizon - Average Reward MDP LP

assume steady state -- \uparrow

$y \in \mathbb{R}^{|A|}$: steady state distribution

$r \in \mathbb{R}^{|A|}$: reward vector

$$\begin{aligned} \max_{y \in \mathbb{R}^{|A|}} & r^T y \rightarrow \text{average expected reward} \\ \text{s.t.} & Ay = Py, \mathbb{1}^T y = 1, y \geq 0 \end{aligned}$$

$$- E_0 w_y = E_i w_y$$

$$- (E_i - E_0) w_y = 0$$

Dual variables:

$$Ay = Py \quad v \in \mathbb{R}^{|S|}$$

$$\mathbb{1}^T y = 1 \quad \lambda \in \mathbb{R}$$

$$y \geq 0 \quad \mu \in \mathbb{R}_+^{|A|}$$

min λ

v, λ, μ

$$\text{s.t.} \quad \lambda \mathbb{1}^T + v^T A = r^T + v^T P + \mu^T$$

$$\mu \geq 0$$

steady state Bellman eqn

λ : overall average cost

v : "value function" tells you how much

ea. r_a differs from λ

μ_a : inefficiency of ea. action

$$\mu_a y_a = 0$$

Note: special conditions
on P required
for this average reward
simulation.

*
see below

\Rightarrow every choice of actions
results in a steady state
distribution

Connections w Markov Chains

normally solving an MDP
selecting a "policy"

policy: mixed strategy at ea. state
(feedback) pick \downarrow probability
distribution over actions

$$\downarrow \downarrow \pi_s \in \mathbb{R}^{|A_s|} \quad \mathbb{1}^T \pi_s = 1 \quad \pi_s \geq 0 \quad (A\pi = I)$$

$$(\pi_s)_a = \text{Prob}(a|s) \quad \varphi = Ay$$

$$\underline{(\pi_s)_a} = \frac{y_a}{\varphi_s} = \frac{y_a}{\sum_{a \in A_s} y_a}$$

$$\Pi \in \mathbb{R}^{|A| \times |S|} \quad \Pi = \begin{bmatrix} \pi_1 & & 0 \\ & \ddots & \\ 0 & & \pi_{|S|} \end{bmatrix}$$

state distribution φ :

$$\vec{y} = \Pi \varphi = \begin{bmatrix} \pi_1 & & \\ & \ddots & \\ & & \pi_{|S|} \end{bmatrix} \begin{bmatrix} \varphi_1 \\ \vdots \\ \varphi_{|S|} \end{bmatrix} = \begin{bmatrix} \pi_1 \varphi_1 \\ \vdots \\ \pi_{|S|} \varphi_{|S|} \end{bmatrix}$$

Markov Chain:

$$M \in \mathbb{R}^{|S| \times |S|} \leftarrow$$

$$M_{s's} = \text{Prob}(s' | s)$$

Steady state dist:

$$\vec{f} = M \vec{f} \Rightarrow \vec{f} \text{ steady state state dist.}$$

When you pick Π or Π \rightarrow selecting Markov chain

$$M = \begin{matrix} & \begin{matrix} |A| \\ \hline |S| \end{matrix} \\ \begin{matrix} |A| \\ \hline |S| \end{matrix} & \begin{bmatrix} P & \\ & \Pi \end{bmatrix} \end{matrix} \leftarrow \quad \underline{\underline{I = A \Pi}}$$

$$\underline{\underline{A y = P y}} \Rightarrow \underline{\underline{A \Pi \varphi = P \Pi \varphi}} \Rightarrow \underline{\underline{\varphi = M \varphi}}$$

$\underline{p} = M\underline{p} \Rightarrow p$ is an ^{right} eigenvector of M
with eigenvalue 1

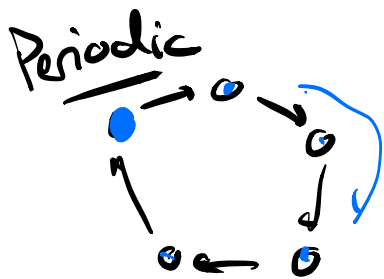
how do we know \underline{p} exists ...

$\underline{1}^T M = \underline{1}^T \rightarrow \underline{1}^T$ is a left eigenvector of M

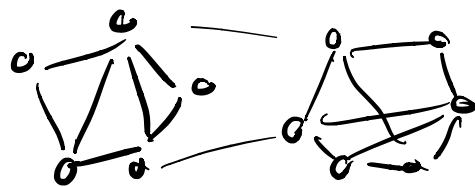
* For any Π :

the resulting Markov matrix $M = P\Pi$
is a periodic & irreducible *

✓ \rightarrow unique solution to $\underline{p} = M\underline{p}$



Not irreducible



Perron Frobenius Thm

$$y \iff \varphi, \Pi$$

$$y \implies \varphi = Ay$$

$$(\Pi_s)_a = y_a / \varphi_s \quad \curvearrowright \quad \checkmark$$

$$\Pi = dg(y) A^T dg(\varphi)^{-1}$$

$$dg(z) = \begin{bmatrix} z_1 & 0 \\ \vdots & \vdots \\ 0 & z_n \end{bmatrix}$$



$$y = \Pi \varphi$$

Discounted Infinite Horizon MDP LP

γ : discount factor $0 \leq \gamma \leq 1$



$$\begin{aligned} \max & \quad \underline{r}^T \underline{y} \\ \text{s.t.} & \quad \underline{A} \underline{y} = \gamma \underline{P} \underline{y} + (1-\gamma) \underline{r}^{(0)} \quad \underline{\mathbb{1}}^T \underline{y} = 1 \quad \underline{y} \geq 0 \end{aligned}$$

$$\begin{aligned} \min & \quad (1-\gamma) \underline{v}^T \underline{r}^{(0)} \quad \leftarrow \begin{array}{l} \text{expected} \\ \text{discounted} \\ \text{reward} \end{array} \\ \text{s.t.} & \quad \underline{v}^T \underline{A} \geq \underline{r}^T + \gamma \underline{v}^T \underline{P} \\ & \quad \text{discounted Bellman equation} \end{aligned}$$

always solvable for any P .