

# ① Game Theory

Q1

	A $q$	B $1-q$
A $p$	<u>2, 2</u>	-1, -1
B $1-p$	-1, -1	<u>1, 1</u>

PNE:

(A, A), (B, B)

MNE:

$$EU_2(p, A) = EU_2(p, B)$$

$$2p - (1-p) = -p + (1-p)$$

$$3p - 1 = -2p + 1$$

$$5p = 2 \Rightarrow p^* = \frac{2}{5}$$

$$q^* = \frac{2}{5}$$

(symmetric game)

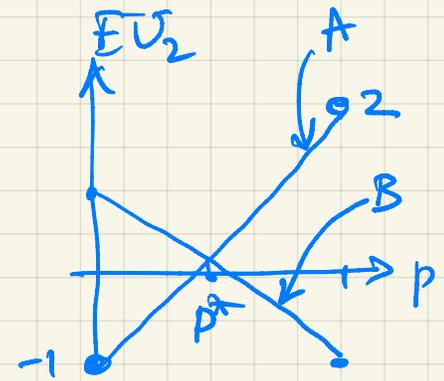
Corresponding utility:

$$EU_2(p^*, A) = EU_2(p^*, B) = 2 \cdot \frac{2}{5} - \frac{3}{5} = \frac{1}{5}$$

$$EU_1(A, q^*) = EU_1(B, q^*) = \frac{1}{5}$$

Q2: MNE is Pareto dominated by PNE

$$MNE < (B, B) < (A, A)$$



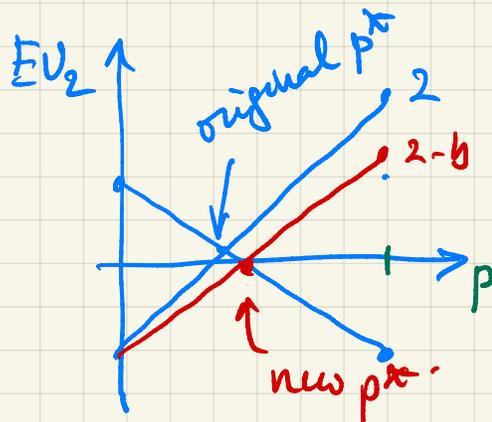
Q3,

$$u_2(A, A) = 2 - b$$

$$(0 < b < 1)$$

	A $q$	B $1-q$
A $p$	2, 2-b	-1, -1
B $1-p$	-1, -1	1, 1

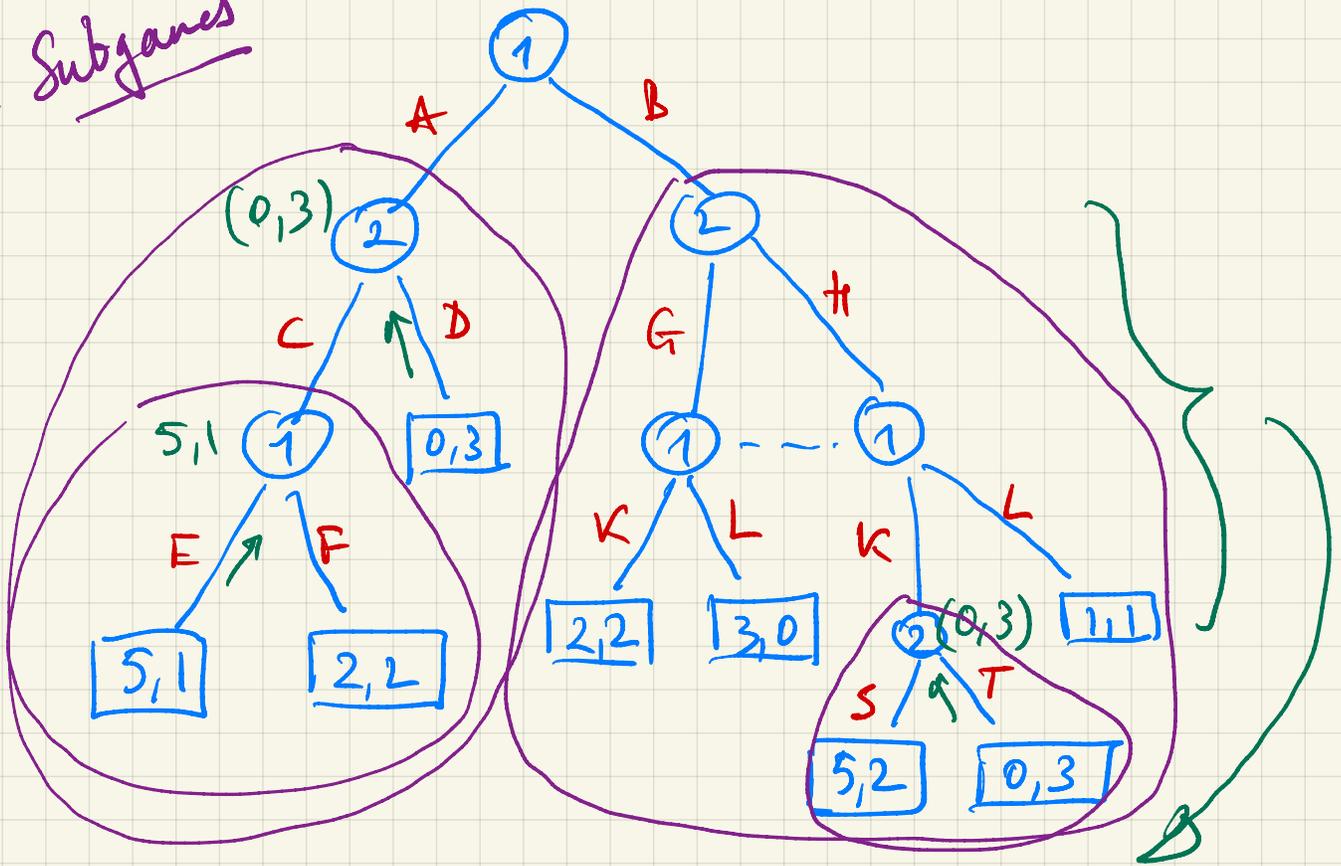
PNE unchanged.



As  $b \uparrow \Rightarrow p^* \uparrow$ .

NB:  $q^*$  is unaffected!

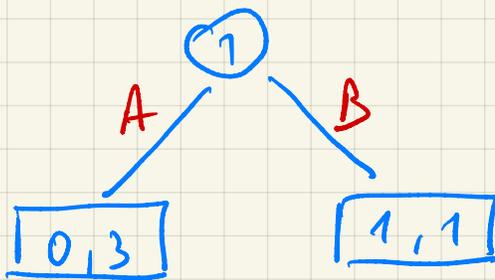
4 Subgames



	G	H
K	2, 2	0, <u>3</u>
L	<u>3</u> , 0	<u>1</u> , <u>1</u>

→ (L, H) is NE

We can reduce the game tree to:

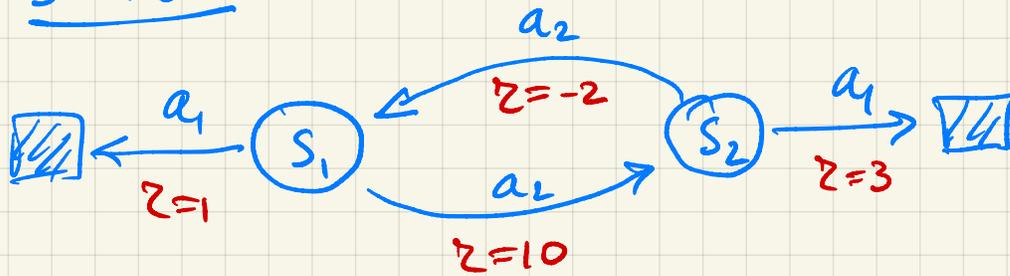


⇒ 1 will choose B  
 2 ——— H  
 1 ——— L

Payoff = (1,1).

Q3: 1 → {A, B} × {EF} × {KL} = {AEK, AE  
 = {AEK, AEL, AFK, AFL, BEK, BEL, BFK, BFL}  
 2 → {C, D} × {GH} × {S, T} = ...

### 3. MDP



Q1: Infinite horizon (could keep on cycling)

Q2:  $\gamma = 0.9 \rightarrow$  optimal policy?

4 possible (deterministic) policies.

$$\pi_1: S_1 \rightarrow a_1, S_2 \rightarrow a_1$$
$$v(S_1) = 1, v(S_2) = 3$$

$$\pi_2: S_1 \rightarrow a_2, S_2 \rightarrow a_1$$
$$v(S_1) = 10 + 3\gamma = 10 + 2.7 = 12.7$$
$$v(S_2) = 3$$

$$\pi_3: S_1 \rightarrow a_1, S_2 \rightarrow a_2$$
$$v(S_1) = 1$$
$$v(S_2) = -2 + \gamma = -2 + 0.9 = -1.1$$

$$\pi_4: S_1 \rightarrow a_2, S_2 \rightarrow a_2$$

$$v(S_1) = 10 + (-2)\gamma + 10\gamma^2 + (-2)\gamma^3 + \dots$$
$$= 10(1 + \gamma^2 + \gamma^4 + \dots) - 2\gamma(1 + \gamma^2 + \dots) \approx \underline{\underline{43.2}}$$

$$\underbrace{43.2}_{(10 - 1.8) \cdot 5.3} \leftarrow \frac{1}{1 - \gamma^2} = \frac{1}{1 - 0.9^2} = \frac{1}{1 - 0.81} = \frac{1}{0.19} \approx 5.3$$

$$\begin{aligned}
v(S_2) &= -2 + \gamma \cdot 10 + (-2)\gamma^2 + 10\gamma^3 + \dots \\
&= -2(1 + \gamma^2 + \dots) + 10\gamma(1 + \gamma^2 + \dots) \\
&= \underbrace{(10\gamma - 2)}_7 \underbrace{(1 + \gamma^2 + \gamma^4 + \dots)}_{\frac{1}{1-\gamma^2} \approx 5.3} \approx \underline{\underline{37.1}}
\end{aligned}$$

Conclusion, for  $\gamma = 0.9$  optimal policy:  $\pi^* = \pi_4$ .

## 4. Reinforcement Learning

① Recall:  $v(s) = \sum_a q(s,a) \pi(a|s)$

$$v(2) = 5 = \underbrace{q(2,R)}_x \underbrace{\pi(R|2)}_{1/4} + \underbrace{q(2,L)}_4 \underbrace{\pi(L|2)}_{3/4}$$

$$\frac{1}{4}x + \frac{3}{4} \cdot 4 = 5 \Rightarrow x = 20 - 12 = 8$$

$$\boxed{q(2,R) = 8}$$

Similarly:  $q(3,L) = 9$

② Q-learning:  $2 \xrightarrow{R} 3$

$$q(2,R) \leftarrow q(2,R) + \alpha \left[ r(2,R) + \gamma \underbrace{\max_{a'} q(3,a')}_{9} - q(2,R) \right]$$

$$\leftarrow 8 + 0.9 \left[ \underbrace{-1 + \frac{2}{3} \cdot 9 - 8}_{-3} \right] = 8 - 2.7 = \underline{\underline{5.3}}$$