

$Q$ -value for MDP  
 $V$ -value

$\Delta$  max  
 $\nabla$  min  
○ chance

# Learning algorithm

Episode:

$$s_0; a_1, r_1, s_1; a_2, r_2, s_2, a_3, r_3, s_3; \dots, a_n, r_n, s_n$$

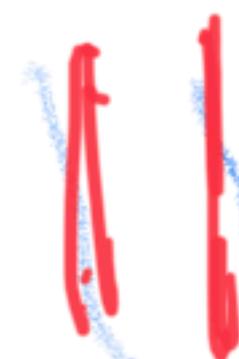
A small piece of experience:

$$(s, a, r, s')$$

Prediction:

$$V(s; \mathbf{w})$$

Target:



$$r + \gamma V(s'; \mathbf{w})$$

$$V(s) = r_{i+1} + \gamma \cdot V(s_{i+1})$$

# Learning algorithm

Episode:

$$s_0; a_1, r_1, s_1; a_2, r_2, s_2, a_3, r_3, s_3; \dots, a_n, r_n, s_n$$

A small piece of experience:

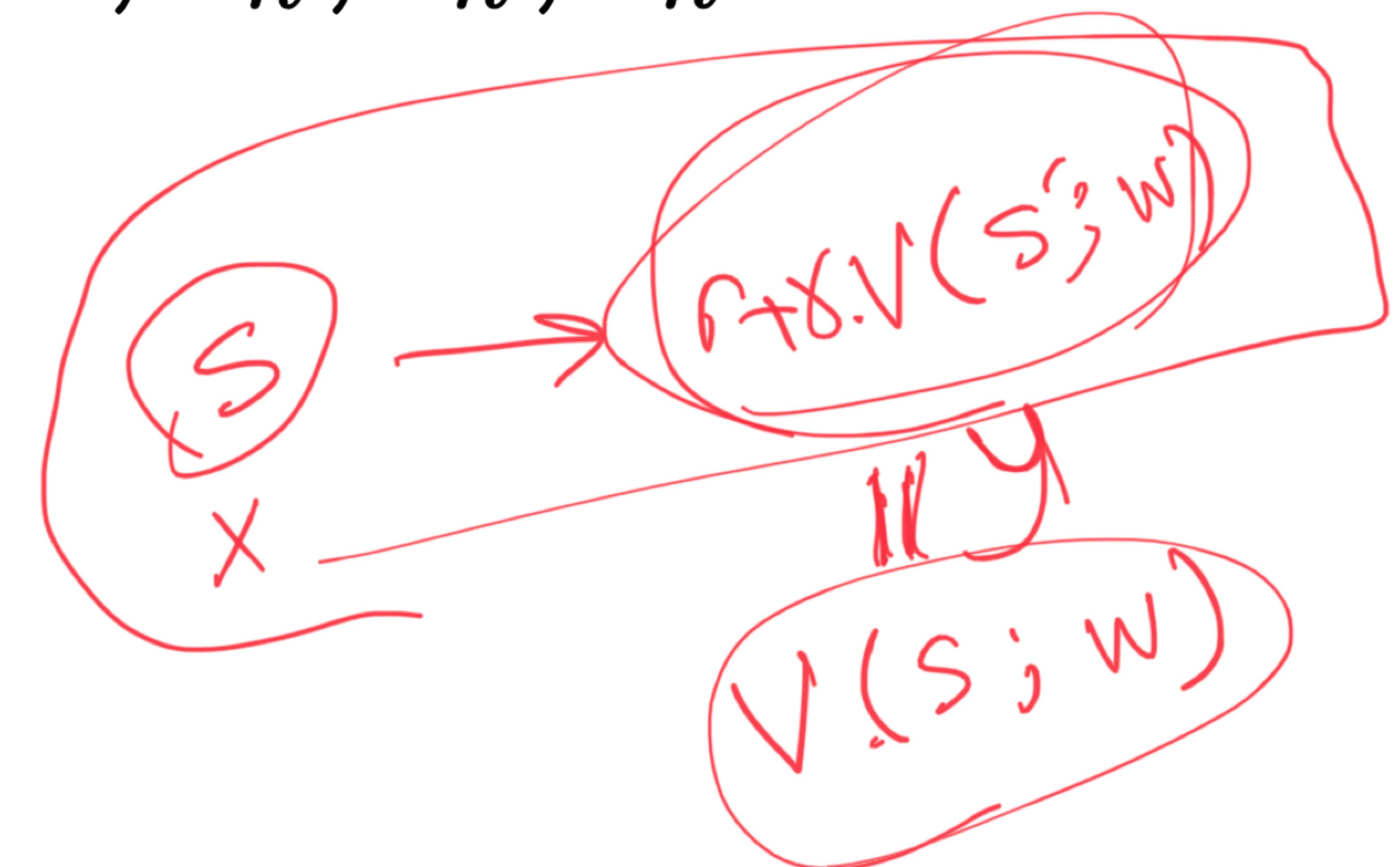
$$(s, a, r, s')$$

Prediction:

$$V(s; \mathbf{w})$$

Target:

$$r + \gamma V(s'; \mathbf{w})$$



# Learning algorithm

Episode:

$$s_0; a_1, r_1, s_1; a_2, r_2, s_2, a_3, r_3, s_3; \dots, a_n, r_n, s_n$$

A small piece of experience:

$$(s, a, r, s')$$

Prediction:

$$V(s; \mathbf{w})$$

Target:

$$r + \gamma V(s'; \mathbf{w})$$

# Example of TD learning

Step size  $\eta = 0.5$ , discount  $\gamma = 1$ , reward is end utility

Example: TD learning						
<b>S1</b>	r:0	<b>S4</b>	r:0	<b>S8</b>	r:1	<b>S9</b>
$\phi: \begin{pmatrix} 0 \\ 1 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 2 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:0 p-t:0	$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:0 p-t:0	$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:1 p-t:-1	$w: \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$
<b>S1</b>	r:0	<b>S2</b>	r:0	<b>S6</b>	r:0	<b>S10</b>
$\phi: \begin{pmatrix} 0 \\ 1 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
$w: \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$	p:1 t:0.5 p-t:0.5	$w: \begin{pmatrix} 0.5 \\ 0.75 \end{pmatrix}$	p:0.5 t:0 p-t:0.5	$w: \begin{pmatrix} 0.25 \\ 0.75 \end{pmatrix}$	p:0 t:0.25 p-t:-0.25	

$$\begin{aligned}
 S &\rightarrow S' \\
 P &= w \cdot \Phi(S) = V(S; w) \\
 t &= r + V(S'; w) \\
 &= r + w \cdot \Phi(S')
 \end{aligned}$$

# Example of TD learning

$$\text{LOSS} = \frac{1}{2} \frac{(\text{Pred}(w) - \text{target})^2}{\text{target}}$$

Step size  $\eta = 0.5$ , discount  $\gamma = 1$ , reward is end utility

Example: TD learning						
<b>S1</b>	r:0	<b>S4</b>	r:0	<b>S8</b>	r:1	<b>S9</b>
$\phi: \begin{pmatrix} 0 \\ 1 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 2 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:0 p-t:0	$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:0 p-t:0	$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:1 p-t:-1	$w: \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$
<b>S1</b>	r:0	<b>S2</b>	r:0	<b>S6</b>	r:0	<b>S10</b>
$\phi: \begin{pmatrix} 0 \\ 1 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
$w: \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$	p:1 t:0.5 p-t:0.5	$w: \begin{pmatrix} 0.5 \\ 0.75 \end{pmatrix}$	p:0.5 t:0 p-t:0.5	$w: \begin{pmatrix} 0.25 \\ 0.75 \end{pmatrix}$	p:0 t:0.25 p-t:-0.25	

# Example of TD learning

$$V(s; w) \approx V(s'; w) + r$$

Step size  $\eta = 0.5$ , discount  $\gamma = 1$ , reward is end utility

Example: TD learning						
S1	r:0	S4	r:0	S8	r:1	S9
$\phi: \begin{pmatrix} 0 \\ 1 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 2 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:0 p-t:0	$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:0 p-t:0	$w: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$	p:0 t:1 p-t:-1	$w: \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$
S1	r:0	S2	r:0	S6	r:0	S10
$\phi: \begin{pmatrix} 0 \\ 1 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 0 \\ 0 \end{pmatrix}$		$\phi: \begin{pmatrix} 1 \\ 0 \end{pmatrix}$
$w: \begin{pmatrix} 0.5 \\ 1 \end{pmatrix}$	p:1 t:0.5 p-t:0.5	$w: \begin{pmatrix} 0.5 \\ 0.75 \end{pmatrix}$	p:0.5 t:0 p-t:0.5	$w: \begin{pmatrix} 0.25 \\ 0.75 \end{pmatrix}$	p:0 t:0.25 p-t:-0.25	

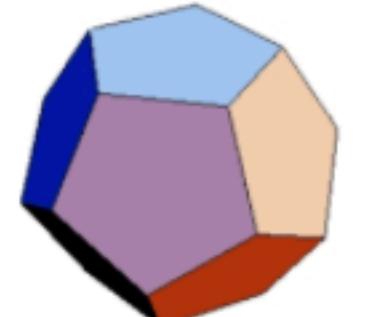
# Game evaluation



## Definition: game evaluation

The **value** of the game if player A follows  $\pi_A$  and player B follows  $\pi_B$  is

$$V(\pi_A, \pi_B) = \sum_{a,b} \pi_A(a)\pi_B(b)V(a, b)$$



## Example: two-finger Morra

Player A always chooses 1:  $\pi_A = [1, 0]$

Player B picks randomly:  $\pi_B = [\frac{1}{2}, \frac{1}{2}]$

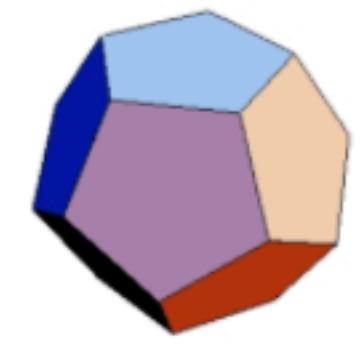
Value:  $-\frac{1}{2}$

$$\begin{matrix} & \begin{matrix} 2 & -3 \\ -3 & 4 \end{matrix} \end{matrix}$$

$$\begin{aligned} & 2 \cdot \pi_A(1) \cdot \pi_B(1) + (-3) \cdot \pi_A(1) \cdot \pi_B(0) \\ & + (-3) \cdot \pi_A(0) \cdot \pi_B(1) + 4 \cdot \pi_A(0) \cdot \pi_B(0) \end{aligned}$$

[whiteboard: matrix]

# Mixed strategies



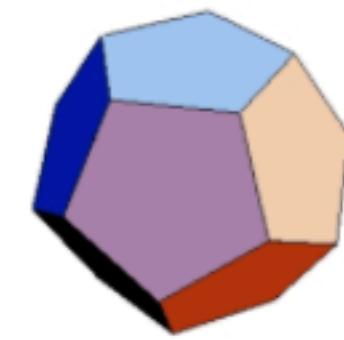
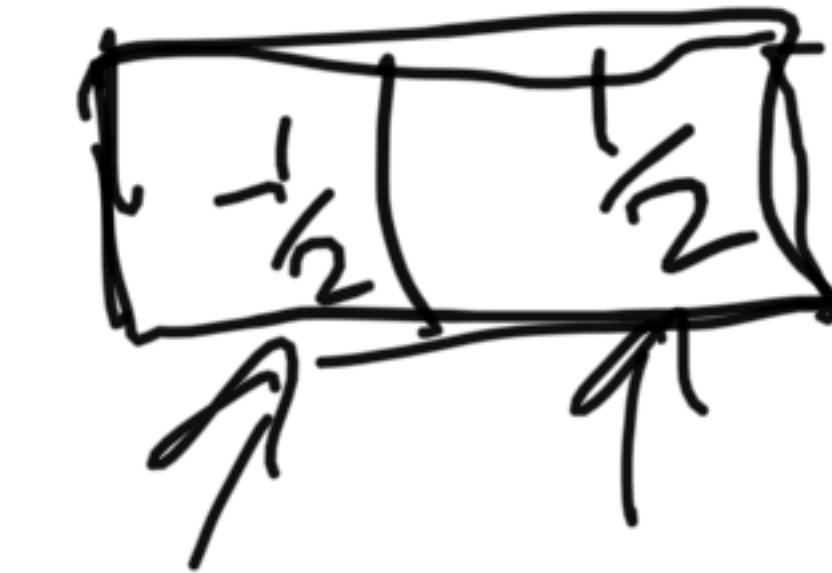
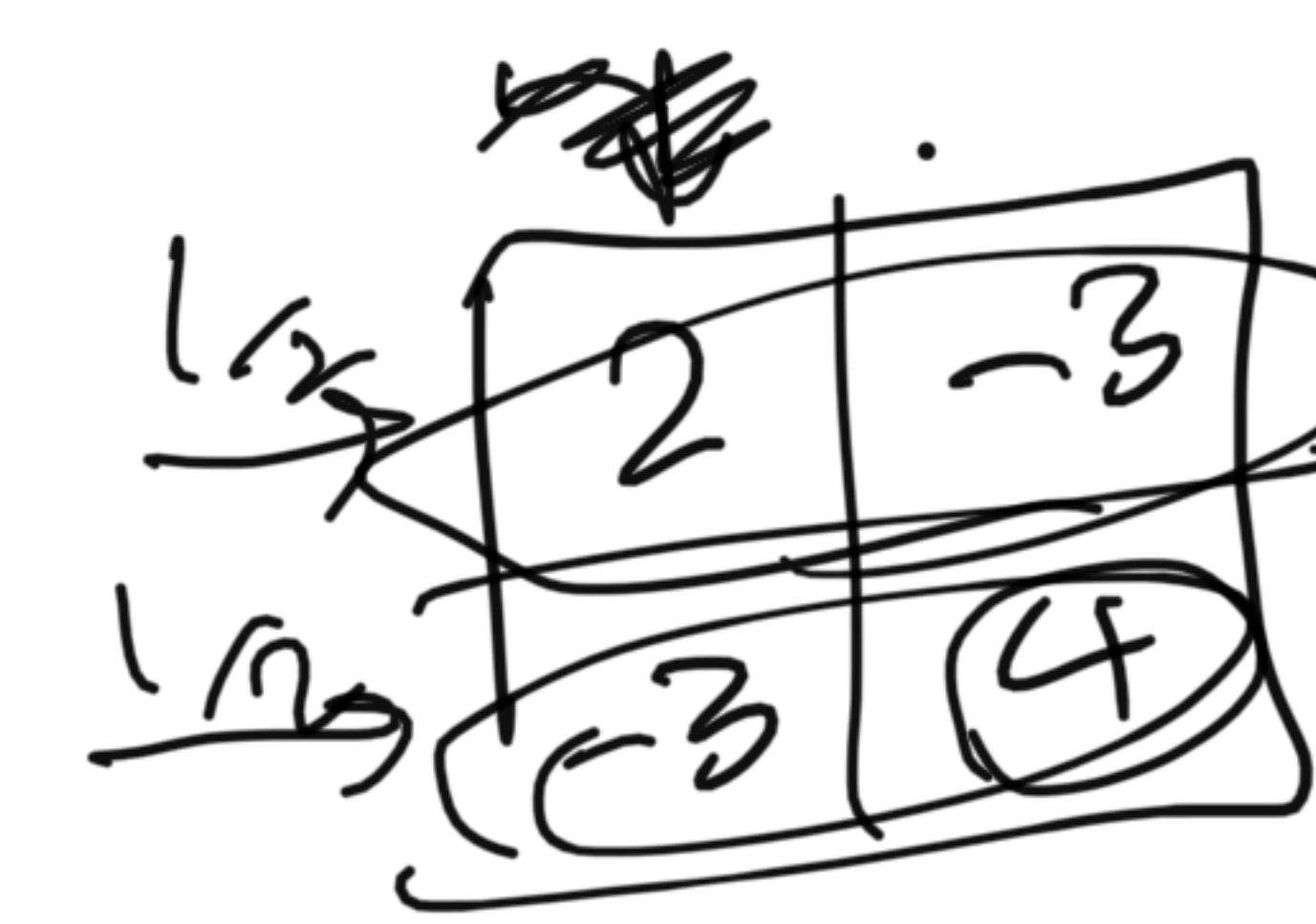
## Example: two-finger Morra

Player A reveals:  $\pi_A = [\frac{1}{2}, \frac{1}{2}]$

Value  $V(\pi_A, \pi_B) = \pi_B(1)(-\frac{1}{2}) + \pi_B(2)(+\frac{1}{2})$

Optimal strategy for player B is  $\pi_B = [1, 0]$  (pure!)

# Mixed strategies

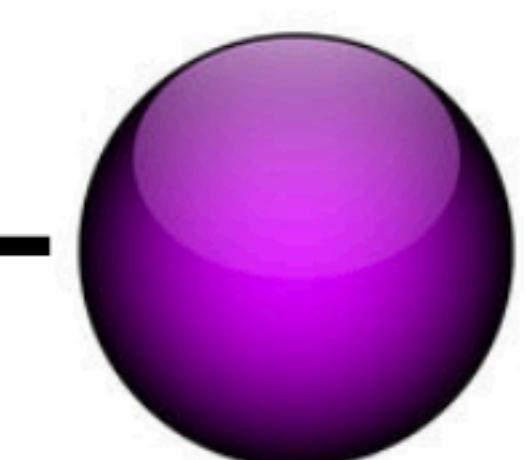


## Example: two-finger Morra

Player A reveals:  $\pi_A = \left[ \frac{1}{2}, \frac{1}{2} \right]$

$$\text{Value } V(\pi_A, \pi_B) = \pi_B(1)\left(-\frac{1}{2}\right) + \pi_B(2)\left(+\frac{1}{2}\right)$$

Optimal strategy for player B is  $\pi_B = [1, 0]$  (pure!)



## Proposition: second player can play pure strategy

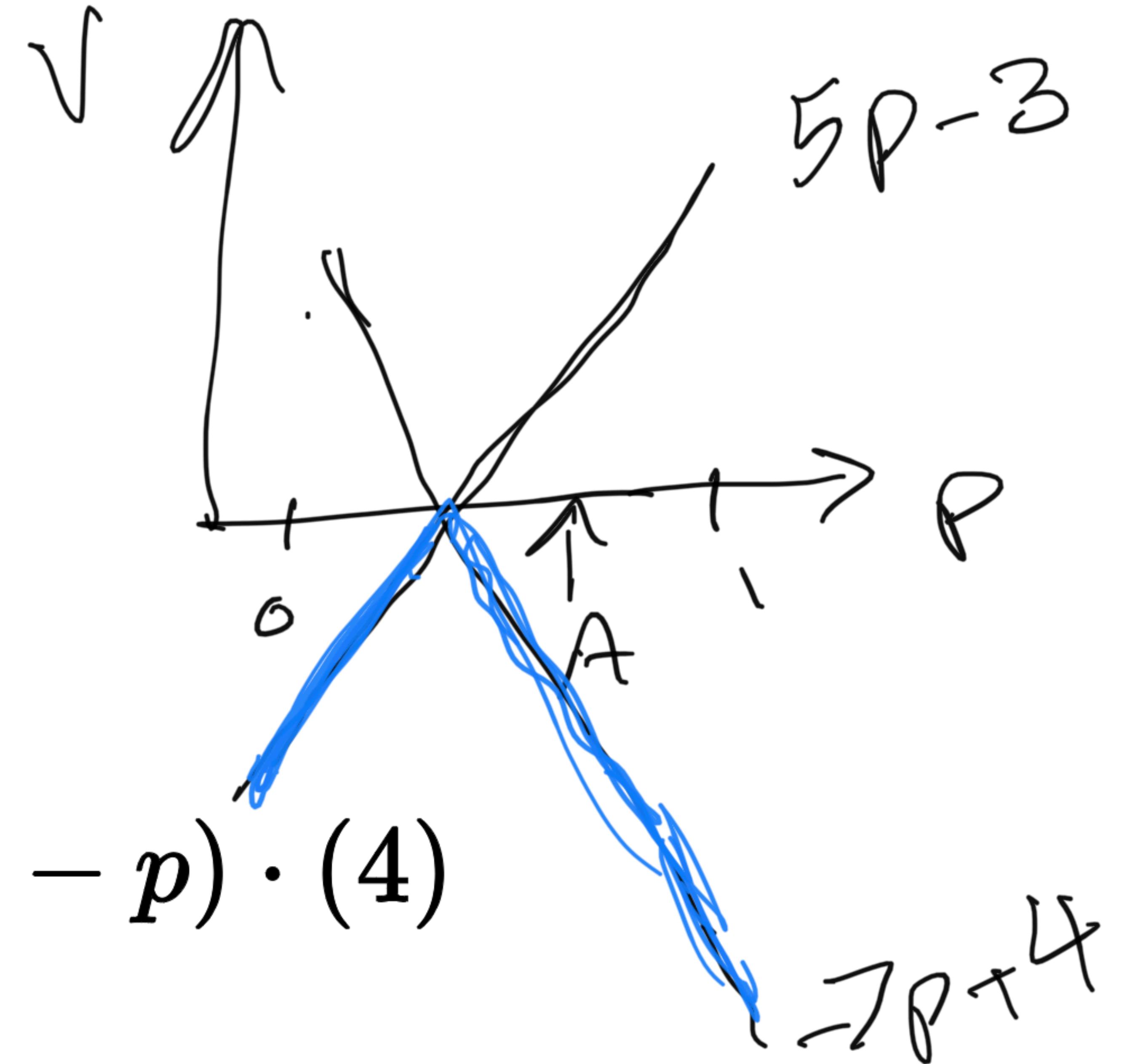
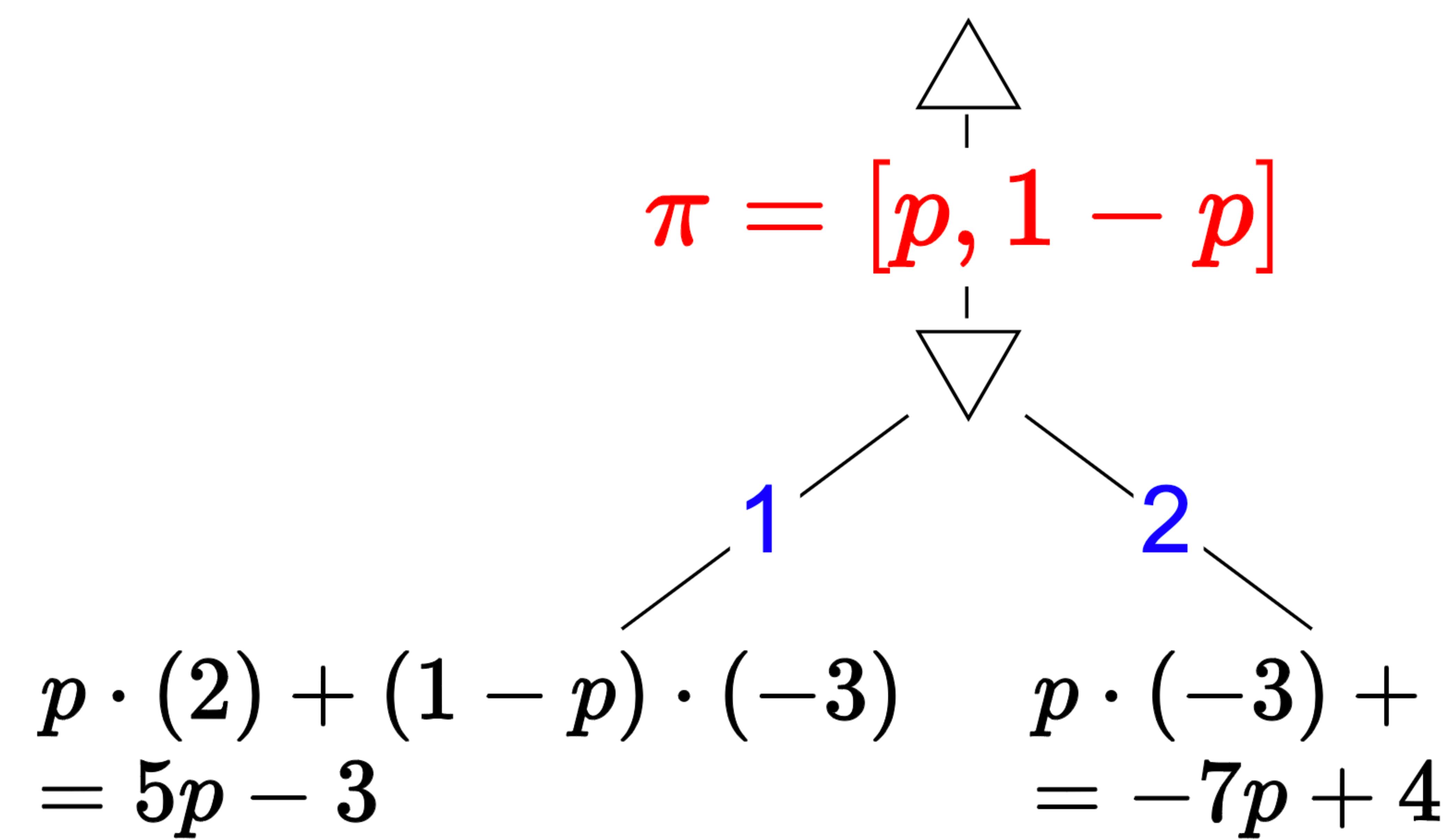
For any fixed mixed strategy  $\pi_A$ :

$$\min_{\pi_B} V(\pi_A, \pi_B)$$

can be attained by a pure strategy.

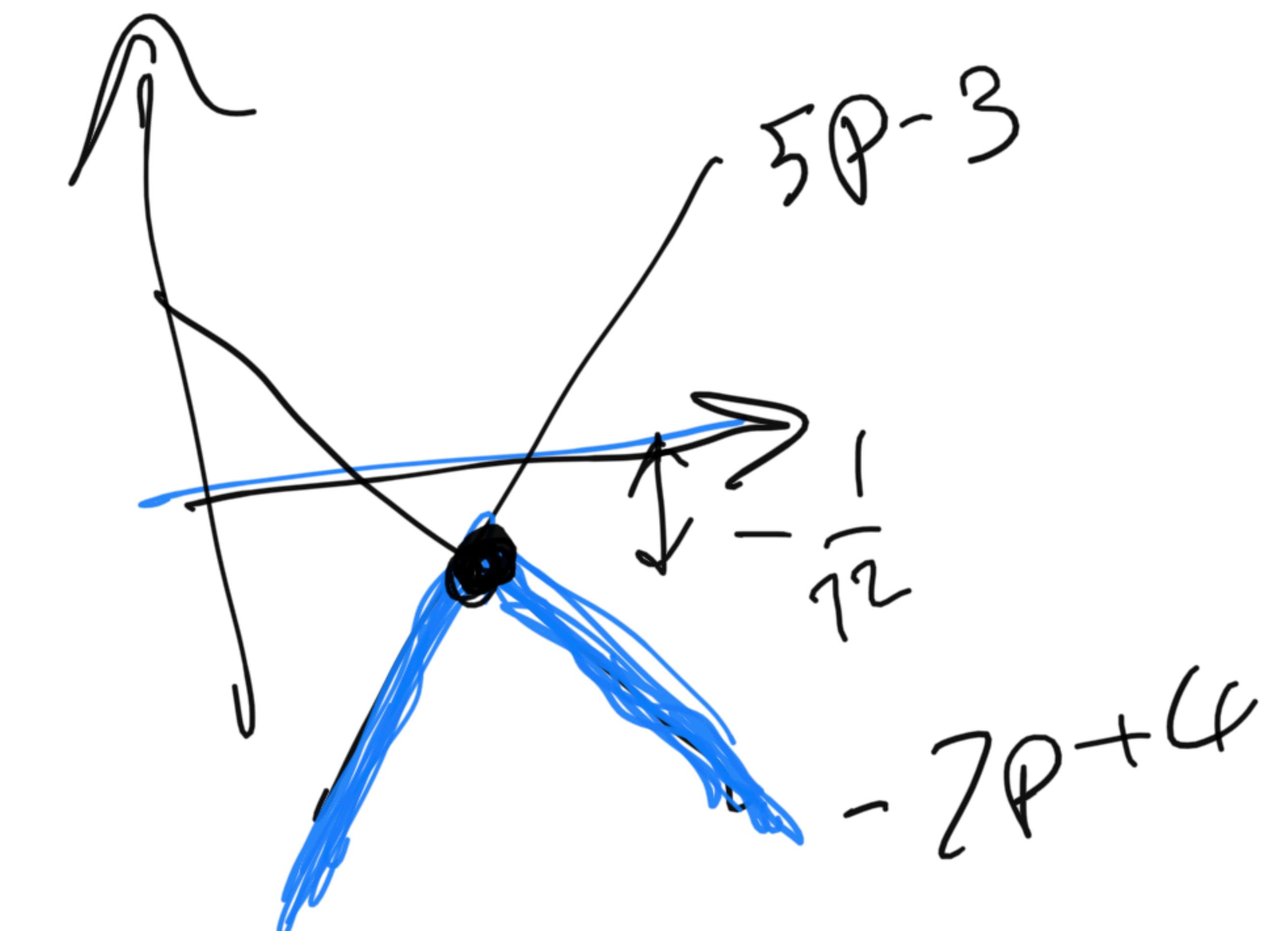
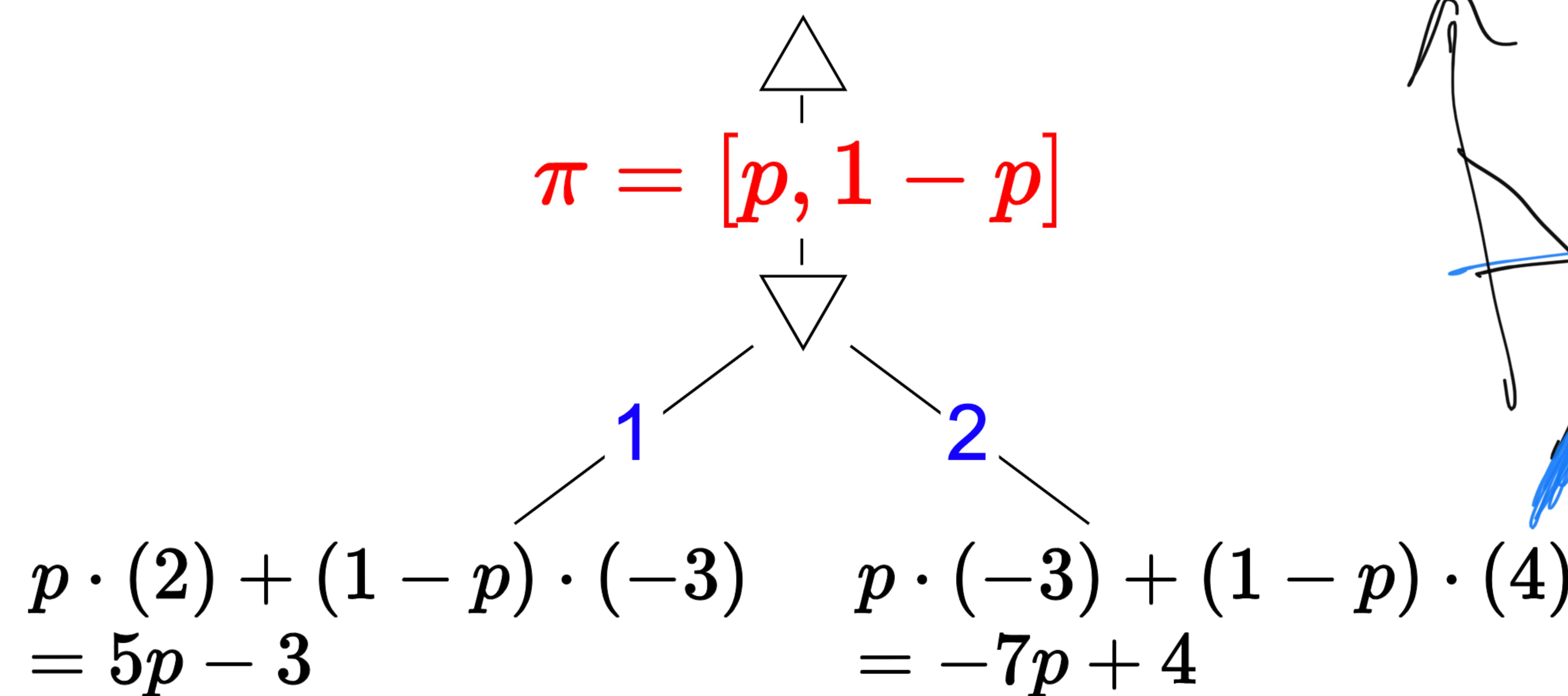
# Mixed strategies

Player A first reveals his/her mixed strategy



# Mixed strategies

Player A first reveals his/her mixed strategy

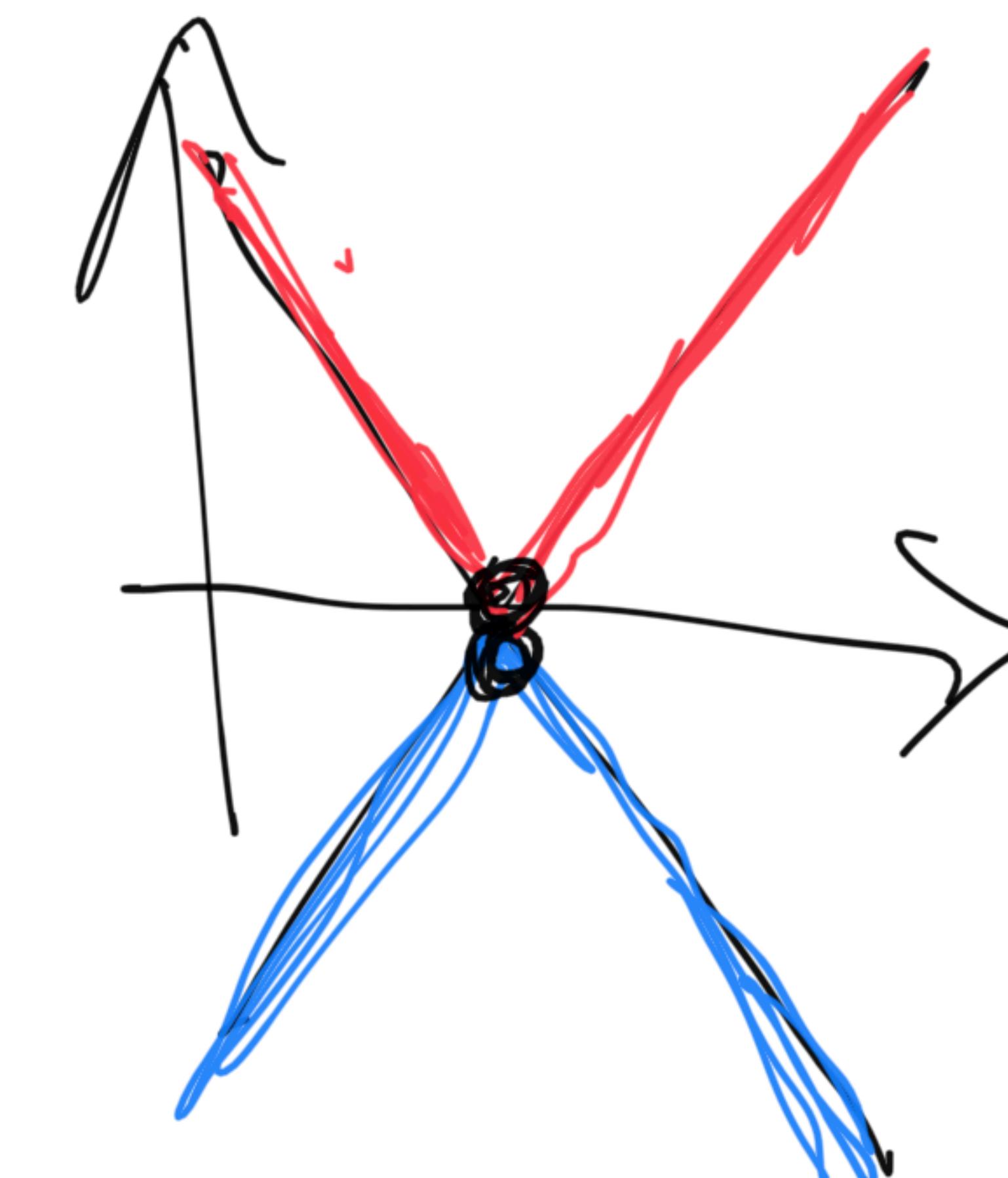
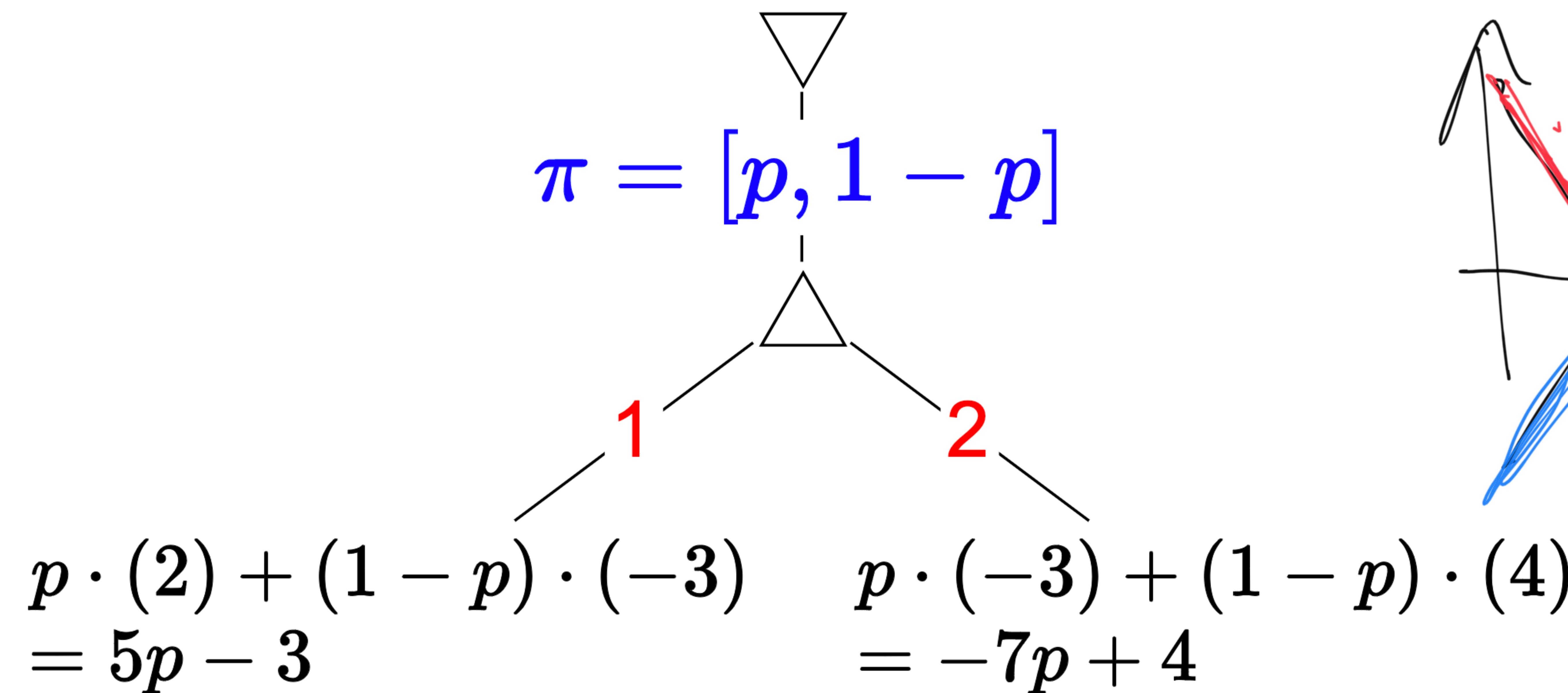


# Minimax value of game:

$$\max_{0 \leq p \leq 1} \min\{5p - 3, -7p + 4\} = \boxed{-\frac{1}{12}} \quad (\text{with } p = \frac{7}{12})$$

# Mixed strategies

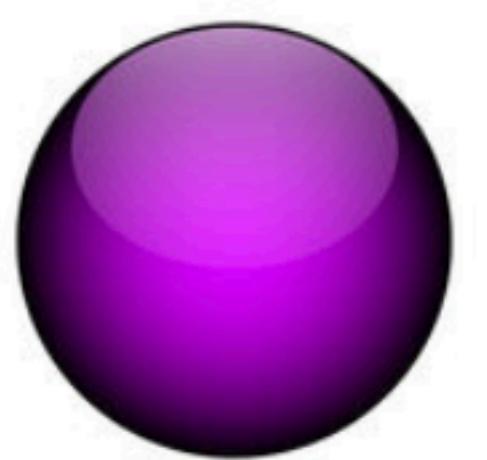
Player B first reveals his/her mixed strategy



Minimax value of game:

$$\min_{p \in [0,1]} \max \{5p - 3, -7p + 4\} = \boxed{-\frac{1}{12}} \text{ (with } p = \frac{7}{12})$$

# General theorem



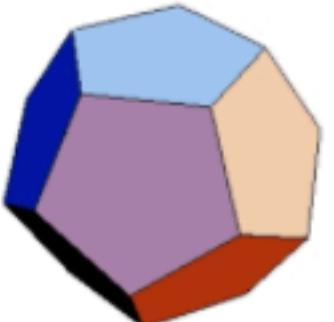
## Theorem: minimax theorem [von Neumann, 1928]

For every simultaneous two-player zero-sum game with a finite number of actions:

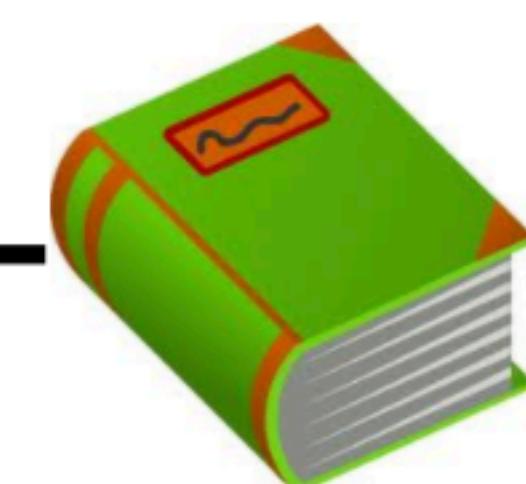
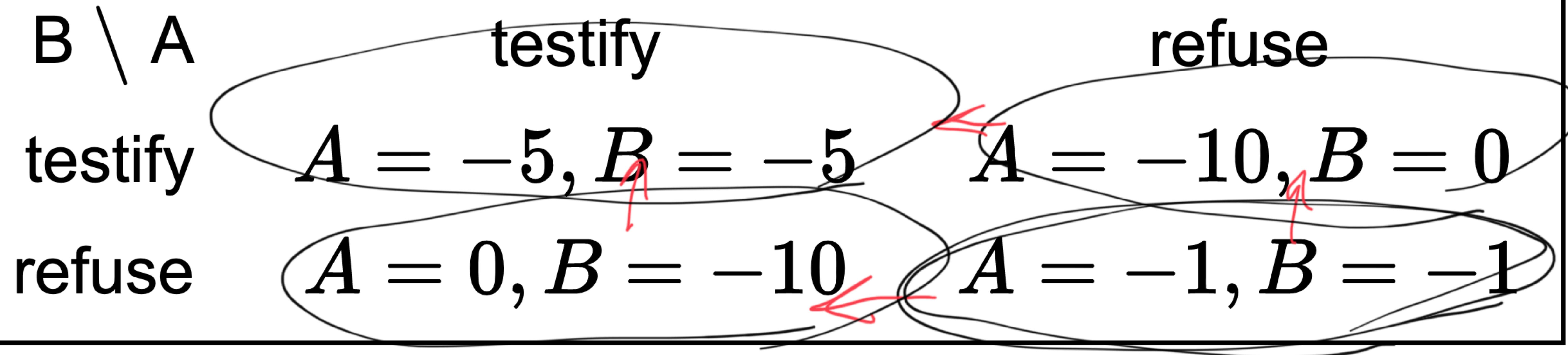
$$\max_{\pi_A} \left( \min_{\pi_B} V(\underline{\pi_A}, \pi_B) \right) = \min_{\pi_B} \left( \max_{\pi_A} V(\pi_A, \underline{\pi_B}) \right),$$

where  $\pi_A, \pi_B$  range over mixed strategies.

# Prisoner's dilemma



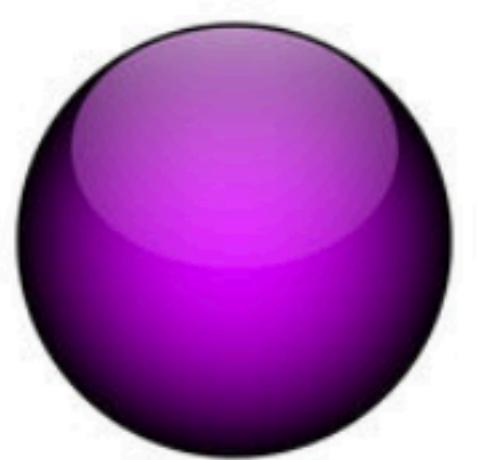
# Example: payoff matrix –



# Definition: payoff matrix-

Let  $V_p(\pi_A, \pi_B)$  be the utility for player  $p$ .

# General theorem



## Theorem: minimax theorem [von Neumann, 1928]

For every simultaneous two-player zero-sum game with a finite number of actions:

$$\max_{\pi_A} \min_{\pi_B} V(\pi_A, \pi_B) \stackrel{?}{=} \min_{\pi_B} \max_{\pi_A} V(\pi_A, \pi_B),$$

where  $\pi_A, \pi_B$  range over **mixed strategies**.

*Mixed*  $\pi_A$  *Pure*  $\pi_B$   
Upshot: revealing your optimal mixed strategy doesn't hurt you!

Proof: linear programming duality