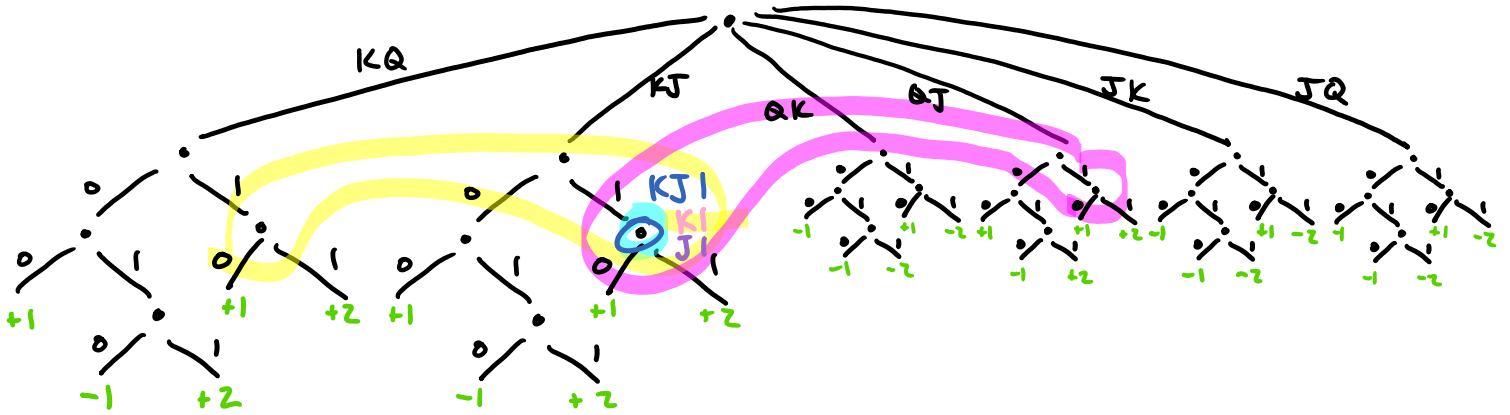


Extensive Form



history: sequence of actions leading to a state
 information set: set of states consistent w/ private history
 strategy: function from info sets to prob dist over actions possible at info sets

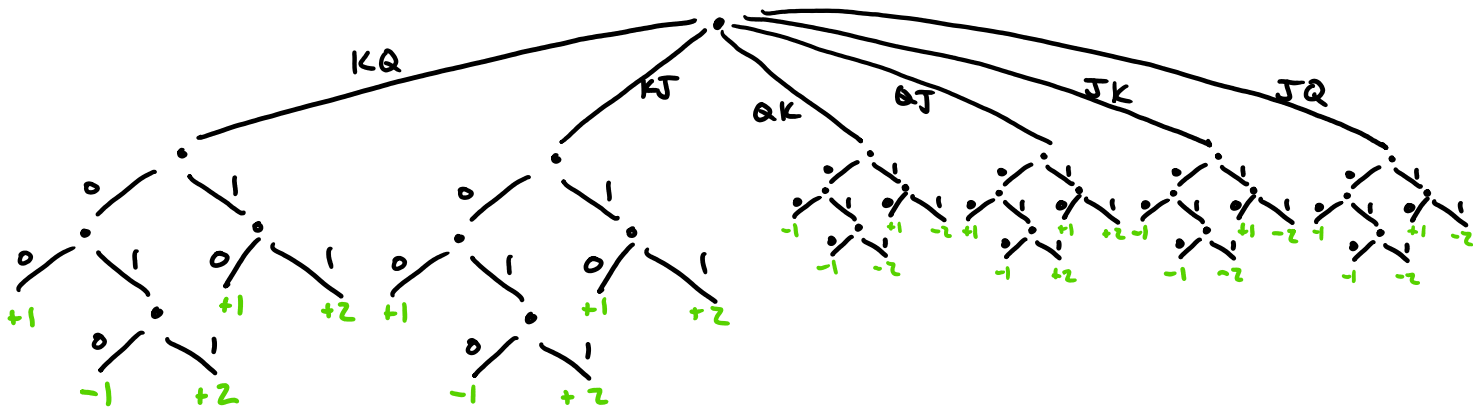
private history: for each p , what p has observed

Normal Form

- 1) Players reveal strategies
- 2) Play game according to strategies

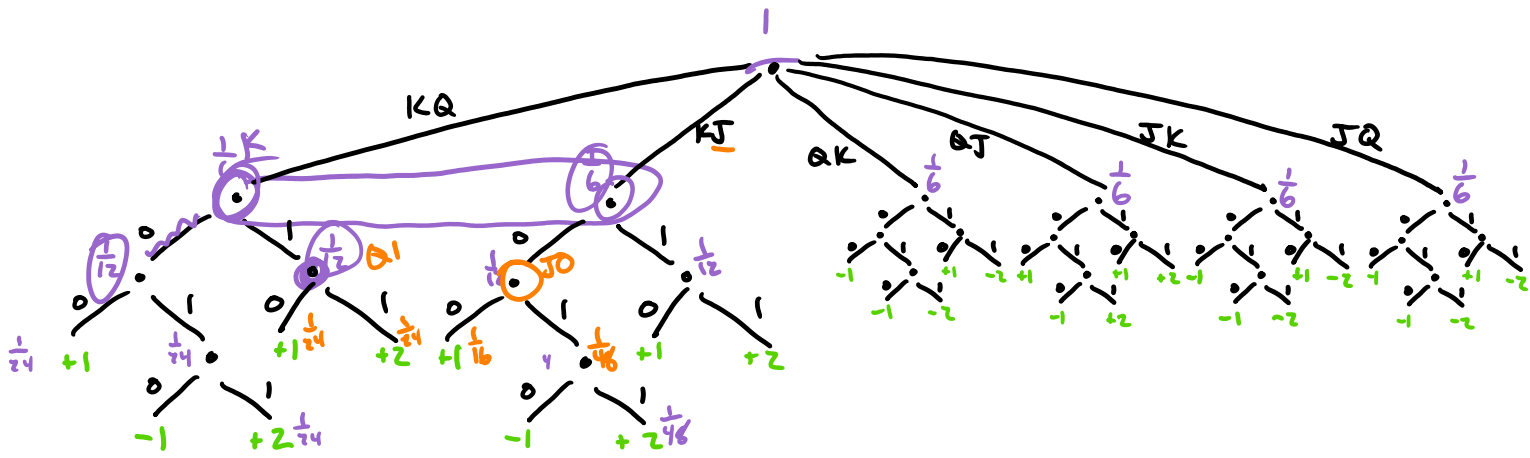
PO : $\overline{K} \quad \overline{Q} \quad \overline{J} \quad \overline{KOI} \quad \overline{QOI} \quad \overline{JOI}$

PI : $\overline{KO} \quad \overline{QO} \quad \overline{JO} \quad \overline{KI} \quad \overline{QI} \quad \overline{JI}$



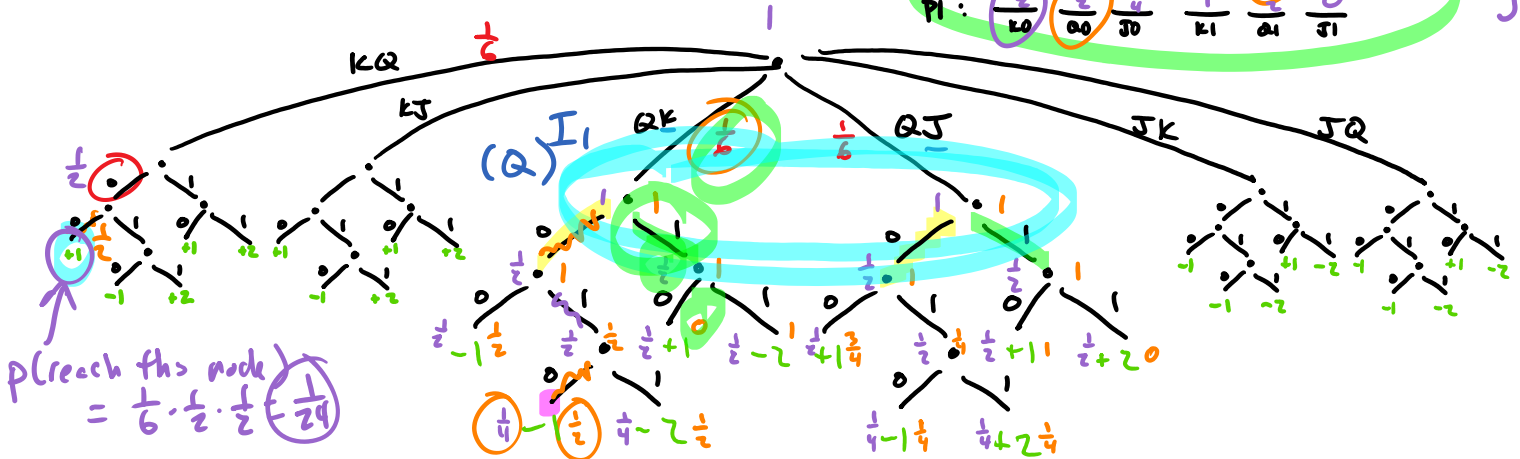
PO : $\frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{2} \quad 1 \quad \frac{1}{2} \quad 0$ ← P(bet)

PI : $\frac{1}{2} \quad \frac{1}{2} \quad \frac{1}{4} \quad 1 \quad \frac{1}{2} \quad 0$



Regret

PO :	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	← probs of betting
PI :	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	



p(reach this node) = $\frac{1}{6} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{24}$

value of σ contributed by I_1 : $\frac{1}{6}(\frac{1}{4} \cdot -1 + \frac{1}{8} \cdot -1 + \frac{1}{8} \cdot -2 + \frac{1}{2} \cdot -2) + \frac{1}{6}(\frac{3}{8} \cdot 1 + \frac{1}{16} \cdot -1 + \frac{1}{16} \cdot 2 + \frac{1}{2} \cdot 1) = \frac{-11}{96}$

value of checking at I_1 :

value of betting at I_1 :

Regret for action = value of action - value of strategy
 positive for good actions
 negative for bad actions

CFR (history, σ , π) \leftarrow players' strategies $[\pi_0, \pi_1, \pi_c]$

if history is terminal
return value for current player
else

CHANCE SAMPLING

OR sample 1 outcome according to prob. dist. of outcomes, let p_a = prob of sampled a
return CFR(.....)

if chance node
for each chance outcome a w/ prob p_a
compute $p_a \cdot \text{CFR}(\text{history} + a, \sigma, [\pi_0, \pi_1, \pi_c \cdot p_a])$
sum those up, return result

else
info \leftarrow info set for current player

$$\begin{matrix} -19 & 10 & 9 \\ \hline a_0 & a_1 & a_2 \end{matrix}$$

keeps track of cumulative regret for each action

$\sigma_{\text{curr}} \leftarrow$ strategy proportional to positive observed regret
compute value of info under each action $\sigma_{\text{curr}} = 0 \quad \frac{10}{19} \quad \frac{9}{19}$

compute value of info under σ_{curr}

for each action a with $\sigma_{\text{curr}}(a) > 0$

compute $-\sigma_{\text{curr}}(a) \cdot \text{CFR}(\text{hist} + a, \sigma, [\pi_1, \pi_0 \cdot \sigma_{\text{curr}}(a), \pi_c])$

update regret and average σ

\leftarrow weighted by $\pi_1 \cdot \pi_c$

\leftarrow converges to Nash eq.
 \leftarrow weighted by π_0

train(n, players)

players \leftarrow array of players w/ empty stats

for $i=1$ to n initial state of game

CFR(root, players, initial probs 1.0 for each player $(1.0, 1.0, \dots)$)

return avg strategy for each player/info set

