

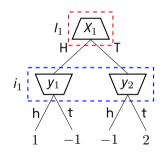
# **Solving Imperfect Information EFGs**

Ondřej Kubíček

Artificial Intelligence Center Faculty of Electrical Engineering Czech Technical University in Prague

# Problems with solving imperfect information games

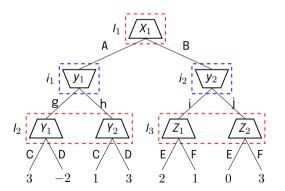




- First player in  $I_1$  chooses between H and T based on the strategy that the player 2 plays.
- Second player in i<sub>1</sub> chooses between h and t based on the probability if it is in y<sub>1</sub> or y<sub>2</sub>, that directly corresponds to the strategy of player 1.
- Backward induction will not work because of this interconnected dependency.
- Generally even in perfect recall imperfect information games, the policy depends on both policy in previous and subsequent parts of the game tree



- Induced normal-form game can include the same leaf multiple times.
- Playing some actions may invalidate different actions in future.

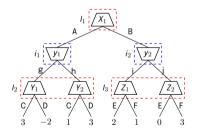


	gi	gj	hi	hj
ACE	3	3	1	1
ACF	3	3	1	1
ADE	-2	-2	3	3
ADF	-2	-2	3	3
BCE	2	0	2	0
BCF	1	3	1	3
BDE	2	0	2	0
BDF	1	3	1	3

#### **Sequences**



- Ordered list of all the actions that may be played in a single playthrough of the game is called Sequence.
- We denote all possible sequences of player *i* as  $\Sigma_i$ .



$\Sigma_1$	$\Sigma_2$
Ø	Ø
Α	g
В	h
AC	i
AD	j
BE	
BF	

#### **Realization plans**

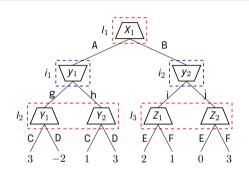


- Realization plan  $r_i(\sigma_i)$  is a probability that  $\sigma_i$  will be played, assuming that the other player plays only actions that allow  $\sigma_i$  to be executed
- Let us assume that  $\sigma_i$  leads to the infoset  $I_i$ . Behavioral strategy  $\pi(I_i, a)$  of playing action a in this infoset is computed as

$$\pi(l_i, a) = \begin{cases} rac{r(\sigma_i a)}{r(\sigma_i)} & \text{if } r(\sigma_i) > 0 \\ 0 & \text{otherwise} \end{cases}$$

•  $\sigma_i a$  represents a extensions of sequence  $\sigma_i$  with action a.





- $r_1(\emptyset) = 1$
- $r_1(A) + r_1(B) = r_1(\emptyset)$
- $r_1(AC) + r_1(AD) = r_1(A)$
- $r_1(BE) + r_1(BF) = r_1(B)$

(d)	1	

- $r_2(\emptyset) = 1$
- $r_2(g) + r_2(h) = r_2(\emptyset)$

 $\Sigma_1$ 

В

AC AD BE BF

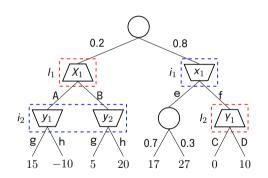
 $\overline{\Sigma_2}$ 0

h

•  $r_2(i) + r_2(j) = r_2(\emptyset)$ 

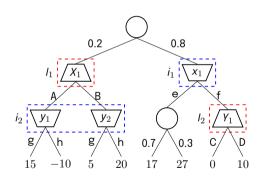
#### **Propagating chance node**

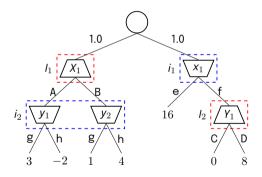




- Chance nodes have fixed probabilities in the game
- These probabilities can be propagated from the chance node to the terminal utilities
- Chance nodes that do not reveal any information until the end of the game can be prunned away completely







### **Extended utility function**



• Extended utility function for sequences is  $g: \Sigma_1 imes \Sigma_2 o \mathbb{R}$ 

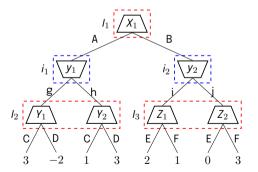
$$g(\sigma_1,\sigma_2) = \sum_{z \in \mathcal{Z}'} \mathcal{C}(z) u(z)$$

- C(z) is a probability that leaf z was reached due to chance nodes along the way.
- $\mathcal{Z}'\subseteq\mathcal{Z}$  are all the terminal histories that could be reached by sequences  $\sigma_1,\sigma_2$



$$\begin{array}{ll} {\it g}(\emptyset,\emptyset) = 0 & {\it g}({\rm BF},{\rm j}) = 3 \\ {\it g}({\rm AC},{\rm i}) = 0 & {\it g}({\rm AD},\emptyset) = 0 \end{array}$$

$$g(\emptyset, g) = 0$$
$$g(A, g) = 0$$



$\Sigma_1$	$\Sigma_2$
Ø	Ø
Α	g
В	h
AC	i
AD	j
BE	
BF	

# **Sequence-form Linear Program (SQF)**



$$\max_{r_1, v} v(I_{\mathsf{root}}) \tag{1}$$

$$s.t. r_1(\emptyset) = 1$$
 (2)

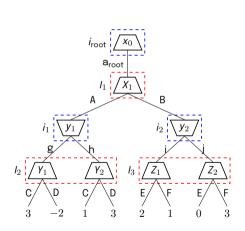
$$r_1(\sigma_1) \ge 0$$
  $\forall \sigma_1 \in \Sigma_1$  (3)

$$\sum_{\boldsymbol{a} \in A_1(I_1)} r(\sigma_1 \boldsymbol{a}) = r_1(\sigma_1) \qquad \forall I_1 \in \mathcal{I}_1, \sigma_1 = \varsigma_1(I_1) \quad (4)$$

$$\sum_{l_2' \in \mathcal{I}_2: \sigma_2 \boldsymbol{\sigma} = \varsigma_2(l_2')} \boldsymbol{v}(l_2') + \sum_{\sigma_1 \in \Sigma_1} \boldsymbol{g}(\sigma_1, \sigma_2 \boldsymbol{\sigma}) r_1(\sigma_1) \geq \boldsymbol{v}(l_2) \quad \forall l_2 \in \mathcal{I}_2, \sigma_2 = \varsigma_2(l_2), \forall \boldsymbol{\sigma} \in \mathcal{A}(l_2) \quad (5)$$

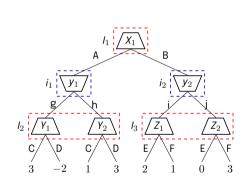
- Variables are realization plans  $r_1$  for all the sequences of player 1 and expected values v for each opponents infoset, if it plays a best response.
- $\varsigma_i$  returns for a given infoset  $l_i$  a sequence  $\sigma_i$  that leads to this infoset.





$$\begin{aligned} \max_{r_1, v} & (v(i_{\mathsf{root}})) \\ & r_1(\emptyset) = 1 \\ & r_1(\mathsf{A}) + r_1(\mathsf{B}) = r_1(\emptyset) \\ & r_1(\mathsf{AC}) + r_1(\mathsf{AD}) = r_1(\mathsf{A}) \\ & r_1(\mathsf{BE}) + r_1(\mathsf{BF}) = r_1(\mathsf{B}) \\ & v(i_1) + v(i_2) \ge v(i_{\mathsf{root}}) \\ & 3r_1(\mathsf{AC}) - 2r_1(\mathsf{AD}) \ge v(i_1) \\ & 1r_1(\mathsf{AC}) + 3r_1(\mathsf{AD}) \ge v(i_1) \\ & 2r_1(\mathsf{BE}) + 1r_1(\mathsf{BF}) \ge v(i_2) \\ & 0r_1(\mathsf{BE}) + 3r_1(\mathsf{BF}) \ge v(i_2) \end{aligned}$$





$$\min_{r_2, v} (v(l_1))$$

$$r_2(\emptyset) = 1$$

$$r_2(g) + r_2(h) = r_2(\emptyset)$$

$$r_2(i) + r_2(j) = r_2(\emptyset)$$

$$v(l_2) \le v(l_1)$$

$$v(l_3) \le v(l_1)$$

$$3r_2(g) + 1r_2(h) \le v(l_2)$$

$$-2r_2(g) + 3r_2(h) \le v(l_2)$$

$$2r_2(i) + 0r_2(j) \le v(l_3)$$

$$1r_2(i) + 3r_2(j) < v(l_3)$$

# **Sequence-form properties**



- The fastest exact algorithm
- Easy to implement
- Poor scaling due to memory requirements of the linear program
- Hard to fine-tune for specific domains to increase performance
- Cannot be used to solve the game only partially

#### **Double Oracle Algorithm**



- Not all strategies in the game are necessary to find Nash equilibrium.
- Double oracle creates smaller game with less strategies and iteratively adds new strategies until it finds Nash equilibrium of the underlying game.
- In the worst case scenario the double oracle algorithm has to add all the strategies from the original game.
- However, in most cases, the restricted game is much smaller than the original game.

### **Double Oracle Algorithm in Extensive-form**



- Using sequences instead of pure strategies is more complicated.
- With limited amount of sequences, some reachable parts of the game tree may not have any sequence to be played.
- To avoid this, in each information set, there is some default action that should be played.
- Default action does not have to be defined explicitly.
- Instead of keeping the full tree, the nodes with default action can be replaced by terminal node.
- The value of this terminal node corresponds to the value if opponent picks best response against the default action.

#### **Double Oracle properties**



- Can solve larger games than SQF.
- Without any additional information the algorithm identifies, which strategies are not important for good solution.
- Computing best response is faster than solving the Linear program and can be improved with some heuristics for specific problems.
- Harder to implement, due to the need to construct the valid restricted game.
- Still requires SQF to solve the restricted game, which is the primary limitation of the method.
- In games where all sequences have to be considered, it is slower than the SQF.

### Policy-space response oracles



- Framework that generalizes double oracle algorithm with empirical game thoeretical analysis.
- It was designed to use double oracle with deep learning with application to very large games.
- Instead of using pure strategies, it uses the full policies.
- The utility of 2 policies against each other can be computed by just letting them play in heads-to-heads.
- This creates matrix game that can be solved through usual means
- Best response can be computed through reinforcement learning, since policy of 1 player is fixed.

# Summary

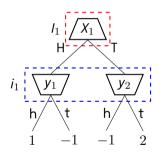


- Mixed strategies are not well suited for solving the imperfect information EFGs.
- Sequences can be used only in perfect recall games, but avoid necessity to define each action for each information set.
- Probability from chance nodes can be propagated up to the leaf utilities.
- If chance node does not reveal any information until the end of the game, it can be removed from the game completely.
- Sequence-form linear program for solving EFGs expands the ideas from linear program for normal-form games to the extensive-form setting.
- It uses the realization plans instead of mixed strategies and expected values for each opponent's infoset.
- Double oracle algorithm can be generalized to a very large games.



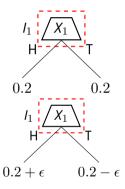
- What if we do not want to solve the whole game, but we only want the strategy in some part of the game.
- In perfect information games, this is done by Monte-Carlo Tree Search.
- What about imperfect information games?
- There are several problems not present in perfect information games.
- One is that the algorithm cannot assume the current state of the game since it is not observable.
- Least amount of considered states has to be closed on the information known by all players.





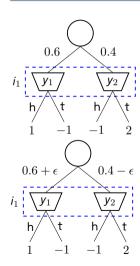
- Optimal strategy in this game is by both players playing heads with 60%.
- What if we are first player and we want to solve this game with depth-limit?
- What if we are the second player and we want to solve this game from our decision onward?
- Let us assume we know opponent's optimal strategy.





- We assume optimal opponent.
- Value of playing heads is  $0.6 \cdot 1 0.4 \cdot 1 = 0.2$
- Value of playing tails is  $-0.6 \cdot 1 + 0.4 \cdot 2 = 0.2$
- Is solution of such a game optimal in the original game?
- What happens with the solution if  $\epsilon > 0$  or  $\epsilon < 0$ ?





- We still assume optimal opponent
- Value of playing heads is  $-0.6 \cdot 1 + 0.4 \cdot 1 = -0.2$
- Value of playing tails is  $0.4 \cdot 1 0.6 \cdot 2 = -0.2$
- What would be a solution of such a game?
- What happens with the solution if  $\epsilon > 0$  or  $\epsilon < 0$ ?



- Both of these problems happen because of the indifference principal.
- In order to have safe depth-limited solving these has to be addressed.
- One possible solution to the first problem is to allow players use several strategies in the future.
- Second is to use iterative approach to converge to a Nash, where after the depth-limit the value depends on the policy in preceding parts
- Second problem is solved by allowing opponent to pick any strategy in the previous part of the game.