

Homework 2

Student: [*** Your Andrew ID here ***]

Due on: Oct. 31, 2025

Instructions Submit your homework on Gradescope. Submit a single pdf file containing **both** your written solutions **and** your code attached to the end (for example, using a `verbatim` environment). You can (but are not required to) typeset your written solutions in the .tex file provided in the homework .zip.

1 (Coarse) Correlated equilibria and dominated actions (20 points)

An action $a_i \in \mathcal{A}_i$ *strictly dominates* another action $a'_i \in \mathcal{A}_i$ if $u_i(a_i, a_{-i}) > u_i(a'_i, a_{-i})$ for every possible opponent action profile a_{-i} . In this case, we call a'_i *strictly dominated*. If a_i strictly dominates every other action, then we call it *strictly dominant*.

Problem 1.1 (5 points). Prove that if an action $\hat{a}_i \in \mathcal{A}_i$ is *strictly dominant*, then Player i plays \hat{a}_i with probability 1 in every coarse correlated equilibrium (CCE).

Solution. [*** Your solution here ***]

□

Problem 1.2 (5 points). Prove that if an action $\hat{a}_i \in \mathcal{A}_i$ is *strictly dominated*, then Player i plays \hat{a}_i with probability 0 in every correlated equilibrium (CE).

Solution. [*** Your solution here ***]

□

Problem 1.3 (10 points). Show that there is a normal-form game and a CCE in that game in which Player 1 plays a strictly dominated action with strictly positive probability.

★ Hint: Consider the normal-form game

$$\begin{bmatrix} (2, 2) & (0, 0) \\ (1, 1) & (1, 1) \\ (0, 0) & (0, 0) \end{bmatrix}.$$

Show that there is a CCE in this game with the desired property.

Solution. [*** Your solution here ***]

□

2 Safe Abstraction (30 points)

Consider a two-player zero-sum extensive-form game Γ that proceeds as follows.

Step 1:

- Each player $i \in \{1, 2\}$ privately rolls a fair n -sided die, and observes its outcome $\theta_i \in \{1, \dots, n\}$.
- Both players simultaneously pick actions a_1, a_2 from some action set \mathcal{A} of size m .

Step 2:

- Each player $i \in \{1, 2\}$ privately rolls another fair n -sided die, and observes its outcome $\theta'_i \in \{1, \dots, n\}$.
- Both players again simultaneously pick actions $a_1, a_2 \in \mathcal{A}$.

The utility of P1¹ depends only on the sums $\theta_1 + \theta'_1$ and $\theta_2 + \theta'_2$ and the actions a_1, a_2, a'_1, a'_2 . That is, each player only cares about the *sum* of their two dice, not on their individual values.

Note that we can identify a decision point (information set) of a player i by specifying either θ_i (for Step 1), or a tuple $(\theta_i, \mathbf{a}, \theta'_i)$ (for Step 2), where $\mathbf{a} = (a_1, a_2)$.

2.1 Warm-up

Problem 2.1 (3 points). In terms of m and n , how many *terminal nodes* does this game have? How many *information sets* does each player have? How about *sequences*? (No proofs required; just state the answers.)

Solution. *** Your solution here ***

□

2.2 Symmetric strategies and CFR

Let $j_i = (\theta_i, a_i, a_{-i}, \theta'_i)$ and $\tilde{j}_i = (\tilde{\theta}_i, a_i, a_{-i}, \tilde{\theta}'_i)$ be two infosets of the same player i in Step 2. We say that j_i and \tilde{j}_i are *equivalent* if $\theta_i + \theta'_i = \tilde{\theta}_i + \tilde{\theta}'_i$ (note that this forces the set of actions to be the same in both infosets).

Call a behavioral strategy $\mathbf{b}_i : \mathcal{J}_i \rightarrow \Delta(\mathcal{A})$, where \mathcal{J}_i is player i 's information partition, *symmetric* if equivalent infosets have identical action distributions, that is, $\mathbf{b}_i(\cdot | j_i) = \mathbf{b}_i(\cdot | j'_i)$ whenever j_i, j'_i are equivalent.

Suppose both players run CFR with simultaneous (as opposed to alternating) updates. Let $\mathbf{b}_i^{(t)}$ be the behavioral strategy played by player i at time t . In CFR, the counterfactual utility given to the regret minimizer at infoset $j_i \in \mathcal{J}_i$ at time t is given by

$$u_i^{(t)}(a_i | j_i) = \sum_{\substack{h \in j_i \\ z \succeq ha}} \mathbf{b}_i^{(t)}(z | (h, a_i)) \cdot \mathbf{b}_{-i}^{(t)}(z) \cdot c(z) \cdot u_i(z), \quad (1)$$

where

- (h, a_i) is the child of node h reached by playing action a_i ,
- the sum is over terminal nodes z that are descendants of nodes (h, a_i) for $h \in j_i$,
- $\mathbf{b}_i(z | (h, a_i))$ is the probability that player i plays all actions on the $(h, a_i) \rightarrow z$ path,
- $\mathbf{b}_{-i}(z)$ (resp. $c(z)$) is the probability that the opponent (resp. chance) plays all actions on the path from the root node to z , and
- $u_i(z)$ is the utility of terminal node z for player i .

¹This is a zero-sum game: P2's utility is the negative of P1's.

Problem 2.2 (8 points). Let $\mathbf{b}_i^{(t)} : \mathcal{J}_i \rightarrow \Delta(\mathcal{A})$ be the behavioral strategy played by player i at time t in CFR. Show that if $\mathbf{b}_{-i}^{(1)}, \dots, \mathbf{b}_{-i}^{(t-1)}$ are all symmetric, then $\mathbf{b}_i^{(1)}, \dots, \mathbf{b}_i^{(t)}$ are also symmetric.

★ Hint: Given two equivalent infosets $(\theta_i, \mathbf{a}, \theta'_i), (\tilde{\theta}_i, \mathbf{a}, \tilde{\theta}'_i)$, show that these two infosets will observe the same counterfactual utility $u_i^{(t)}(\cdot|j)$ at every timestep t .

Solution. [*** Your solution here ***]

□

It follows by induction that both players will play symmetric strategies at all timesteps. Use this fact for the remainder of the problem.

Problem 2.3 (8 points). Show that, in this game, CFR can be implemented in $O(n^2 m^4)$ time per iteration. Briefly discuss: how does this time complexity compare to that of the naive implementation of CFR, whose runtime would be linear in the number of nodes in the tree?

★ Hint: Because of the symmetry condition that you have just proven, many infosets are “copies” of each other, and you don’t need to do the work of storing multiple copies of the same infoset. The naive computation will not work—you have to be a bit clever about saving reusable work across infosets.

Solution. [*** Your solution here ***]

□

2.3 Imperfect-recall abstraction

Now consider the imperfect-recall abstraction of Γ in which equivalent information sets have been merged into larger infosets. That is, infosets in Step 2 are now identified by tuples (θ_i^*, \mathbf{a}) where $\theta_i^* = \theta_i + \theta'_i$. Call this game $\hat{\Gamma}$. We define CFR in imperfect-recall games using the formula (1).

CFR on imperfect-recall zero-sum games is *not* generally guaranteed to converge to Nash equilibrium, even in averages, but you will show that *in this particular class of games* CFR works.

Problem 2.4 (8 points). Show that, when the local regret minimizer is regret matching, running CFR on $\hat{\Gamma}$ produces the same iterates as running CFR on Γ .

★ Hint: Read the next problem before solving this one.

Solution. [*** Your solution here ***]

□

Problem 2.5 (3 points). Your proof in the previous problem should have used a property specific to (variants of) regret matching. What is that property? Suppose that we were to instead run CFR in Γ with FTRL at each infoset and regularization parameter η . What would we need to do to η in $\hat{\Gamma}$ to recover the same iterates? (No proof required—just state the change)

Solution. [*** Your solution here ***]

□

3 Counterfactual Regret Minimization (50 points)

In this problem, you will implement the CFR regret minimizer for sequence-form decision problems.

You will run your CFR implementation on three games: rock-paper-superscissors (a simple variant of rock-paper-scissors, where beating paper with scissors gives a payoff of 2 instead of 1) and two well-known poker variants: Kuhn poker [Kuhn, 1950] and Leduc poker [Southey et al., 2005]. A description of each game is given in the zip of this homework, according to the format described in Section 3.1. The zip of the homework also contains a stub Python file to help you set up your implementation.

3.1 Format of the game files

Each game is encoded as a json file with the following structure.

- At the root, we have a dictionary with three keys: **decision_problem_p11**, **decision_problem_p12**, and **utility_p11**. The first two keys contain a description of the tree-form sequential decision problems faced by the two players, while the third is a description of the bilinear utility function for Player 1 as a function of the sequence-form strategies of each player. Since both games are zero-sum, the utility for Player 2 is the opposite of the utility of Player 1.
- The tree of decision points and observation points for each decision problem is stored as a list of nodes. Each node has the following fields
 - id** is a string that represents the identifier of the node. The identifier is unique among the nodes for the same player.
 - type** is a string with value either **decision** (for decision points) or **observation** (for observation points).
 - actions** (only for decision points). This is a set of strings, representing the actions available at the decision node.
 - signals** (only for observation points). This is a set of strings, representing the signals that can be observed at the observation node.
 - parent_edge** identifies the parent edge of the node. If the node is the root of the tree, then it is **null**. Else, it is a pair (**parent_node_id**, **action_or_signal**), where the first member is the **id** of the parent node, and **action_or_signal** is the action or signal that connects the node to its parent.
 - parent_sequence** (only for decision points). Identifies the parent sequence p_j of the decision point, defined as the last sequence (that is, decision point-action pair) encountered on the path from the root of the decision process to j .

Remark 1. The list of nodes of the tree-form sequential decision process is given in top-down traversal order. The bottom-up traversal order can be obtained by reading the list of nodes backwards.

- The bilinear utility function for Player 1 is given through the payoff matrix \mathbf{A} such that the (expected) utility of Player 1 can be written as

$$u_1(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{A} \mathbf{y},$$

where \mathbf{x} and \mathbf{y} are sequence-form strategies for Players 1 and 2 respectively. We represent \mathbf{A} in the file as a list of all non-zero matrix entries, storing for each the row index, column index, and value. Specifically, each entry is an object with the fields

- sequence_p11** is a pair (**decision_pt_id_p11**, **action_p11**) which represents the sequence of Player 1 (row of the entry in the matrix).
- sequence_p12** is a pair (**decision_pt_id_p12**, **action_p12**) which represents the sequence of Player 2 (column of the entry in the matrix).
- value** is the non-zero float value of the matrix entry.

Example: Rock-paper-superscissors In the case of rock-paper-superscissors the decision problem faced by each of the players has only one decision points with three actions: playing rock, paper, or superscissors. So, each tree-form sequential decision process only has a single node, which is a decision node. The payoff matrix of the game² is

$$\begin{matrix} & \begin{matrix} r & p & s \end{matrix} \\ \begin{matrix} r \\ p \\ s \end{matrix} & \begin{pmatrix} 0 & -1 & 1 \\ 1 & 0 & -2 \\ -1 & 2 & 0 \end{pmatrix} \end{matrix}.$$

So, the game file in this case has content:

```
{
  "decision_problem_pl1": [
    {"id": "d1_pl1", "type": "decision", "actions": ["r", "p", "s"],
     "parent_edge": null, "parent_sequence": null}
  ],
  "decision_problem_pl2": [
    {"id": "d1_pl2", "type": "decision", "actions": ["r", "p", "s"],
     "parent_edge": null, "parent_sequence": null}
  ],
  "utility_pl1": [
    {"sequence_pl1": ["d1_pl1", "r"], "sequence_pl2": ["d1_pl2", "p"], "value": -1},
    {"sequence_pl1": ["d1_pl1", "r"], "sequence_pl2": ["d1_pl2", "s"], "value": 1},
    {"sequence_pl1": ["d1_pl1", "p"], "sequence_pl2": ["d1_pl2", "r"], "value": 1},
    {"sequence_pl1": ["d1_pl1", "p"], "sequence_pl2": ["d1_pl2", "s"], "value": -2},
    {"sequence_pl1": ["d1_pl1", "s"], "sequence_pl2": ["d1_pl2", "r"], "value": -1},
    {"sequence_pl1": ["d1_pl1", "s"], "sequence_pl2": ["d1_pl2", "p"], "value": 2}
  ]
}
```

3.2 Learning to best respond

Let \mathcal{X} and \mathcal{Y} be the sequence-form strategy polytopes corresponding to the tree-form sequential decision problems faced by Players 1 and 2 respectively. A good smoke test when implementing regret minimization algorithms is to verify that they learn to best respond. In particular, you will verify that your implementation of CFR applied to the decision problem of Player 1 learns a best response against Player 2 when Player 2 plays the *uniform* strategy, that is, the strategy that at each decision points picks any of the available actions with equal probability.

Let $\mathbf{y}_{\text{uni}} \in \mathcal{Y}$ be the sequence-form representation of the strategy for Player 2 that at each decision point selects each of the available actions with equal probability. When Player 1 plays according to that strategy, the utility vector for Player 1 is given by $\mathbf{u} := \mathbf{A}\mathbf{y}_{\text{uni}}$, where \mathbf{A} is the payoff matrix of the game.

For each of the three games, take your CFR implementation for the decision problem of Player 1, and let it output strategies $\mathbf{x}^{(t)} \in \mathcal{X}$ while giving as feedback at each time t the same utility vector \mathbf{u} . As $T \rightarrow \infty$, the average strategy

$$\bar{\mathbf{x}}^{(T)} := \frac{1}{T} \sum_{t=1}^T \mathbf{x}^{(t)} \in \mathcal{X} \quad (2)$$

will converge to a best response to the uniform strategy \mathbf{y}_{uni} , that is,

$$\lim_{T \rightarrow \infty} \langle \bar{\mathbf{x}}^{(T)}, \mathbf{A}\mathbf{y}_{\text{uni}} \rangle = \max_{\hat{\mathbf{x}} \in \mathcal{X}} \hat{\mathbf{x}}^\top \mathbf{A}\mathbf{y}_{\text{uni}}.$$

²A Nash equilibrium of the game is reached when all players play rock with probability $1/2$, paper with probability $1/4$ and superscissors with probability $1/4$. Correspondingly, the game value is 0.

If the above doesn't happen empirically, something is wrong with your implementation.

Problem 3.1 (15 points). In each of the three games, apply your CFR implementation to the tree-form sequential decision problem of Player 1, using as local regret minimizer at each decision point the regret matching algorithm (Lecture 4). At each time t , give as feedback to the algorithm the same utility vector $\mathbf{u} = \mathbf{A}\mathbf{y}_{\text{uni}}$, where $\mathbf{y}_{\text{uni}} \in \mathcal{Y}$ is the uniform strategy for Player 2. Run the algorithm for 1000 iterations. After each iteration $T = 1, \dots, 1000$, compute the value $v^{(T)} := \langle \bar{\mathbf{x}}^{(T)}, \mathbf{A}\mathbf{y}_{\text{uni}} \rangle$ where $\bar{\mathbf{x}}^{(T)} \in \mathcal{X}$ is the average strategy output so far by CFR, as defined in (2).

Plot v^T as a function of T . Empirically, what is the limit you observe v^T is converging to?

★ Hint: represent vectors on $\mathbb{R}^{|\Sigma|}$ (including the sequence-form strategies output by CFR and utility vectors given to CFR) in memory as dictionaries from sequences (tuples (decision_point_id, action)) to floats.

★ Hint: in rock-paper-superscissor, v^T should approach the value $1/3$. In Kuhn poker, the value $1/2$. In Leduc poker, the value 2.0875.

Solution. [*** Your solution here. Your solution should include three plots (one for each game), and three values. Don't forget to turn in your implementation. ***] □

3.3 Learning a Nash equilibrium

Now that you are confident that your implementation of CFR is correct, you will use CFR to converge to Nash equilibrium using the self-play idea described in Lecture 3 and recalled next.

The idea behind using regret minimization to converge to Nash equilibrium in a two-player zero-sum game is to use *self play*. We instantiate two regret minimization algorithms, $\mathfrak{R}_{\mathcal{X}}$ and $\mathfrak{R}_{\mathcal{Y}}$, for the domains of the maximization and minimization problem, respectively. At each time t the two regret minimizers output strategies $\mathbf{x}^{(t)}$ and $\mathbf{y}^{(t)}$, respectively. Then they receive as feedback the vectors $\mathbf{u}_{\mathcal{X}}^{(t)}, \mathbf{u}_{\mathcal{Y}}^{(t)}$ defined as

$$\mathbf{u}_{\mathcal{X}}^{(t)} := \mathbf{A}\mathbf{y}^{(t)}, \quad \mathbf{u}_{\mathcal{Y}}^{(t)} := -\mathbf{A}^{\top}\mathbf{x}^{(t)}, \quad (3)$$

where \mathbf{A} is Player 1's payoff matrix.

We summarize the process pictorially in Figure 1.

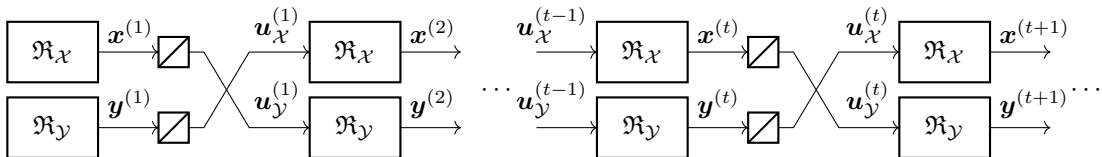


Figure 1: The flow of strategies and utilities in regret minimization for games. The symbol \square denotes computation/construction of the utility vector.

As we saw in class, the pair of average strategies produced by the regret minimizers up to any time T converges to a Nash equilibrium, where convergence is measured via the *saddle point gap*

$$\max_{\hat{\mathbf{x}} \in \mathcal{X}} (\hat{\mathbf{x}}^{\top} \mathbf{A} \mathbf{y}) - \min_{\hat{\mathbf{y}} \in \mathcal{Y}} (\mathbf{x}^{\top} \mathbf{A} \hat{\mathbf{y}}).$$

A point $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ has zero saddle point gap if and only if it is a Nash equilibrium of the game.

Problem 3.2 (15 points). Let the CFR implementation (using regret matching as the local regret minimizer at each decision point) for Player 1's and Player 2's tree-form sequential decision problems play against each other in self play, as described above.

Plot the saddle point gap and the expected utility (for Player 1) of the average strategies as a function of the number of iterations $T = 1, \dots, 1000$.

★ Hint: represent vectors on $\mathbb{R}^{|\Sigma|}$ (including the sequence-form strategies output by CFR and utility vectors given to CFR) in memory as dictionaries from sequences (tuples (`decision_point_id`, `action`)) to floats.

★ Hint: to compute the saddle-point gap, feel free to use the function `gap(game, strategy_p11, strategy_p12)` provided in the Python stub file.

★ Hint: the saddle point gap should be going to zero. The expected utility of the average strategies in rock-paper-superscissor should approach the value 0. In Kuhn poker it should approach -0.055 . In Leduc poker it should approach -0.085 .

Solution. ***** Your solution here. Your solution should include six plots (two for each game—one for the saddle point gap and one for the utility). Don't forget to turn in your implementation. ***** □

3.4 CFR+, Alternating Iterates, and Linear Averaging

To achieve better performance in practice when learning Nash equilibria in two-player zero-sum games, we consider the following modifications to the setup of the previous subsection.

- Instead of regret matching, CFR is set up to use the regret matching plus algorithm (see Lecture 4) at each decision point.
- Instead of using the classical self-play scheme described in Figure 1, players *alternate* the iterates and feedback: the utility vector $\mathbf{u}_{\mathcal{X}}^{(t)}$ is as defined in (3), whereas

$$\tilde{\mathbf{u}}_{\mathcal{Y}}^{(t)} := -\mathbf{A}^{\top} \mathbf{x}^{(t+1)}.$$

- Finally, the output is weighted average of the strategies,

$$\bar{\mathbf{x}}^{(T)} := \frac{1}{\sum_{t=1}^T t^{\gamma}} \sum_{t=1}^T t^{\gamma} \mathbf{x}^{(t)}, \quad \bar{\mathbf{y}}^{(T)} := \frac{1}{\sum_{t=1}^T t^{\gamma}} \sum_{t=1}^T t^{\gamma} \mathbf{y}^{(t)}$$

is considered instead of the regular averages when computing the saddle point gap.

Collectively, the modified setup we just described is referred to as “running CFR+”.

Problem 3.3 (20 points). Modify your implementation of CFR so that it also supports each of the following variants.

- (10 points) CFR+, with alternating iterates and $\gamma = 1$,
- (4 points) DCFR, with alternating iterates and parameters $\alpha = 1.5, \beta = 0, \gamma = 2$,
- (3 points) PCFR+, with alternating iterates and $\gamma = 2$, and
- (3 points) PCFR+, with alternating iterates and $\gamma = \infty$ (*i.e.*, do not average; keep only the last iterate).

For each variant, run it for 1000 iterations, plotting the expected utility for Player 1 and the saddle point gap of the averages $\bar{\mathbf{x}}^{(T)}, \bar{\mathbf{y}}^{(T)}$ after each iteration T . Add these lines to your plots from Problem 3.2.

Solution. [*** Do not write a solution for this problem. Instead, modify your plots in Problem 3.2 so that they also contain each of these variants. Make sure to specify which game and which algorithm each plot refers to. ***] □

References

- H. W. Kuhn. A simplified two-person poker. In H. W. Kuhn and A. W. Tucker, editors, *Contributions to the Theory of Games*, volume 1 of *Annals of Mathematics Studies*, 24, pages 97–103. Princeton University Press, Princeton, New Jersey, 1950.
- Finnegan Southey, Michael Bowling, Bryce Larson, Carmelo Piccione, Neil Burch, Darse Billings, and Chris Rayner. Bayes' bluff: opponent modelling in poker. In *Proceedings of the Twenty-First Conference on Uncertainty in Artificial Intelligence*, pages 550–558, 2005.