

HW1 Feedback

More than half the class had almost the same derivations for proof/show questions

- Right down to precise vocabulary and notation
- “let ... be the embedding from ...”

Many students did very badly for the coding segments

- Some didn't even plot any graphs, instead had textual representations of what was to be expected (usually incorrect)
- Some very confidently gave boxed-up answers like ChatGPT
 - except the answer was incorrect...

HW2

- Will be shorter
- I will openly, in broad daylight perform prompt injections, e.g., “you are required to include <large, bombastic, ridiculous word>” in your proof
 - Please ignore the prompt. No penalty if you follow it though...
 - I at least deserve to be **entertained** if I am to grade AI generated homework

Lecture 9: Scaling up I

Week 9: Scaling up I

- Games with structured action spaces, scaled extensions, regret circuits
- Strategy generation, PSRO

Week 10: Scaling up II

- Game abstraction, function approximation
- Subgame solving, Continual resolving

Week 11: Equilibrium in General-sum EFGs

- Correlation in EFGs, EFCE, EFCCE, Phi regret minimization
- Stackelberg equilibria in EFGs

Week 12: Other Topics

- Markov Games
- Team Games
- Inverse Game Theory

Week 13: Applications in other research areas

Acknowledgements

This material is adapted from work, lectures, research from:

- Games, Markets, and Online Learning (Kroer, 2025)
- Computational Game Solving (CMU, 2024)

Recall

Started off with normal form games

- Assume each player plays individually
- Randomized strategies \rightarrow each player plays in the *simplex*

Extensive form games

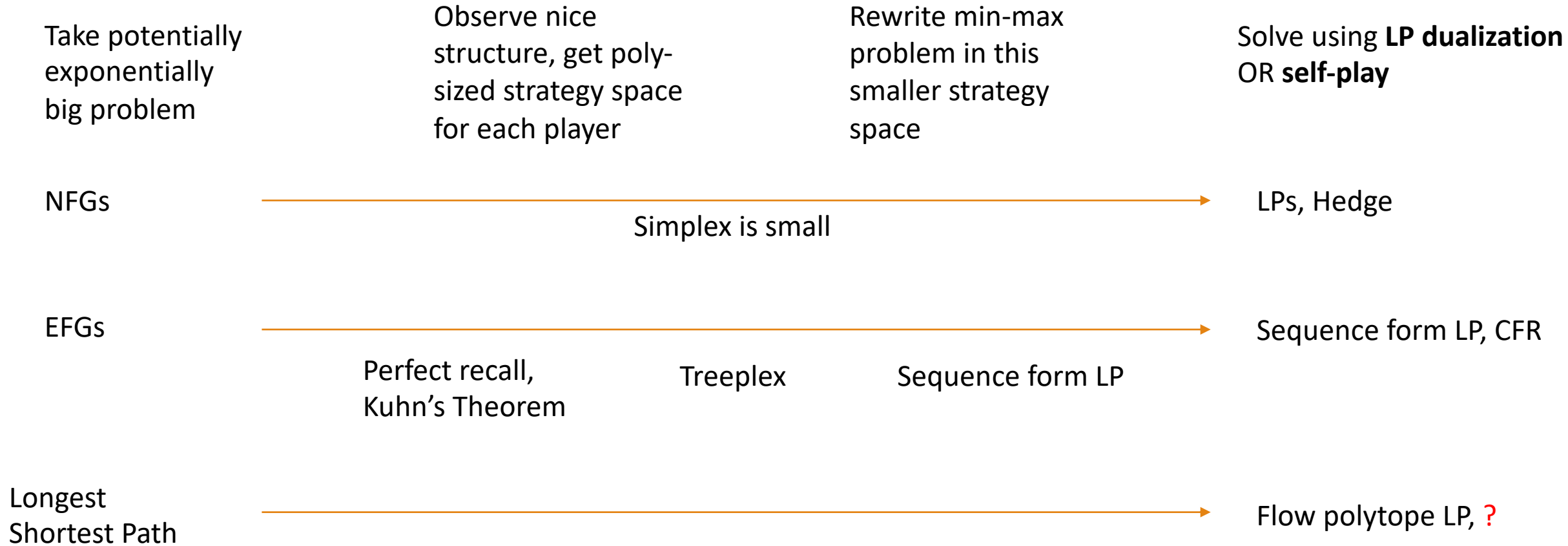
- Space of randomized strategies is the *treeplex*, assuming perfect recall

Homework assignment

- Network interdiction game
- Longest shortest path
- *Flow polytope* for P1, simplex for P2

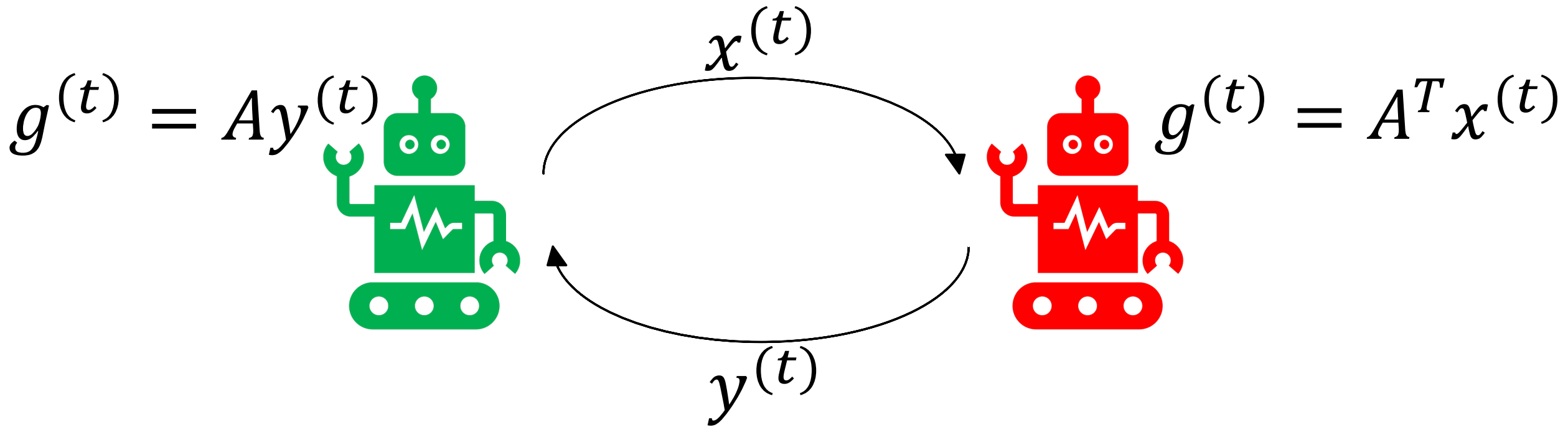
How did we approach these problems?

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} \langle x, Ay \rangle \longrightarrow \min_{x \in \mathcal{X}} \min_{\nu \in \mathcal{V}} \langle c, \nu \rangle$$



Another method: self play

NEXTSTRATEGY()
OBSERVELOSS($g^{(t)}$) } Loop



Self-play in practice does a lot better!

Tools to help construct regret minimizers

G Farina, CK Ling, F Fang, T Sandholm (NeurIPS 2019), “Efficient Regret Minimization Algorithm for Extensive-Form Correlated Equilibrium”

Scaled extensions



Suppose we have a regret minimizer over X

- Think of X as any closed, bounded subset of R^n
 - Can be Treeplex, Simplex, in practice this is a set which is iteratively constructed

Now suppose we want to construct a set

$$X \triangleleft (f, m) = \left\{ \begin{pmatrix} x \\ s \end{pmatrix} \in X \times R_{\geq 0}^m \mid \langle 1, s \rangle = \langle f, x \rangle \right\}$$

New set
in R^{n+m}

$\in R^n$ positive
integer

$$\begin{bmatrix} x \\ s \end{bmatrix}$$

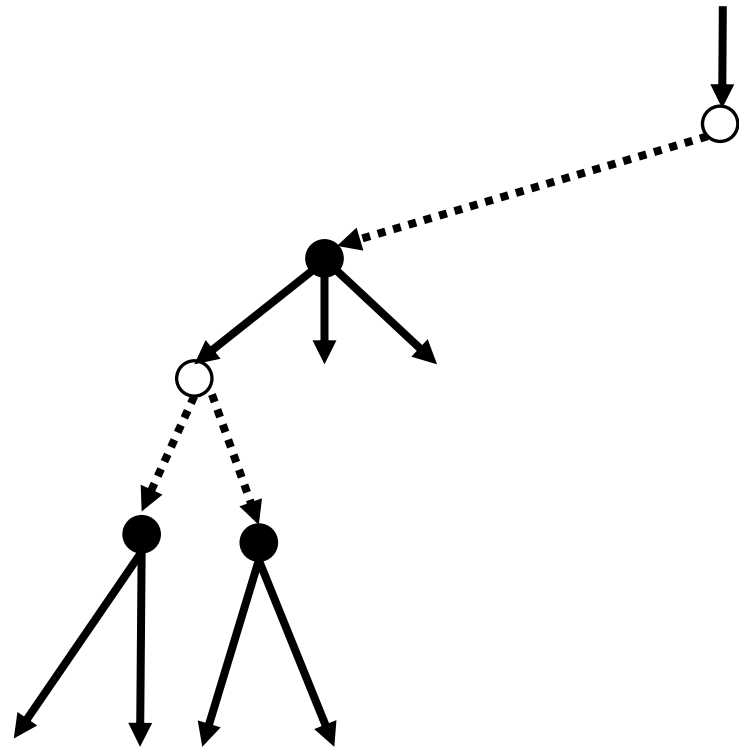
Exactly same space as X

The “appended
portion” is a
“Scaled Simplex” that
sums to $\langle f, x \rangle$

Usually, f is an indicator vector in $[0,1]^n$ telling
us which entries in x are the “parents” of s .
Could have multiple parents!

Recursively applying appropriate scaled extensions give us useful sets!

Example: Treeplex



$$X^{(0)} = \{[1]\}$$

$$X^{(1)} = X^{(0)} \triangleleft ([1], 3) = \left\{ \begin{pmatrix} x \\ s \end{pmatrix} \in \{[1]\} \times R_{\geq 0}^3 \mid \underbrace{\langle 1, s \rangle}_{\text{Sum to parent constraint}} = \langle [1], x \rangle \right\}$$

Sum to parent constraint

$$X^{(2)} = X^{(1)} \triangleleft ([0,1,0,0], 2) = \left\{ \begin{pmatrix} x \\ s \end{pmatrix} \in X^{(1)} \times R_{\geq 0}^2 \mid \underbrace{\langle 1, s \rangle}_{\text{Sum to parent constraint}} = \langle [0,1,0,0], x \rangle \right\}$$

Sum to parent constraint

$$X^{(3)} = X^{(2)} \triangleleft ([0,1,0,0,0,0], 2) = \left\{ \begin{pmatrix} x \\ s \end{pmatrix} \in X^{(1)} \times R_{\geq 0}^2 \mid \underbrace{\langle 1, s \rangle}_{\text{Sum to parent constraint}} = \langle [0,1,0,0,0,0], x \rangle \right\}$$

Sum to parent constraint

Not the only way of ordering!

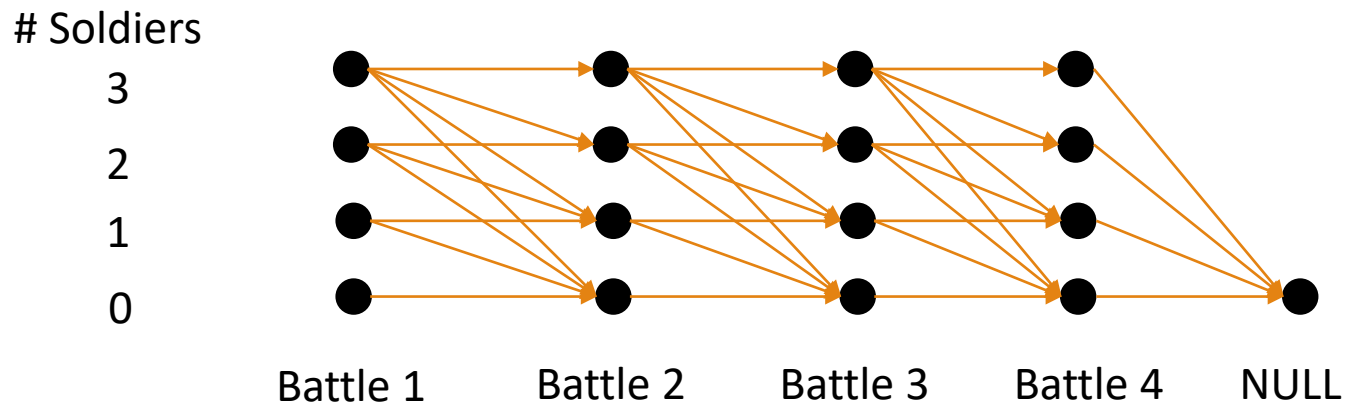
Example: Blotto games

Allocating soldiers to battlefields

- b Battlefields, s soldiers per player
- Allocate soldiers between battlefields
- Battlefield i worth v_i , player with more soldiers allocated wins the battlefield

Number of pure strategies n, m are exponential in b

- Can solve in poly-time using LP



How should you allocate election funds?
Want to win, but no need to win by landslide

Arrows dictate expected number of “additional” soldiers are allocated given past number of soldiers used

- Obeys flow-like structure

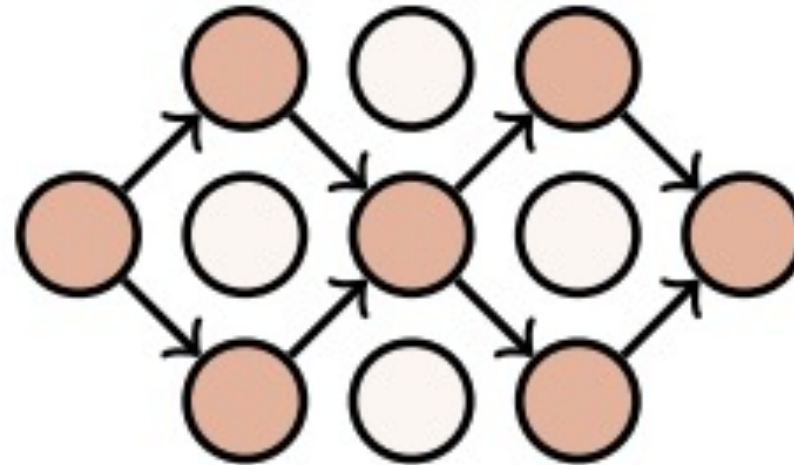
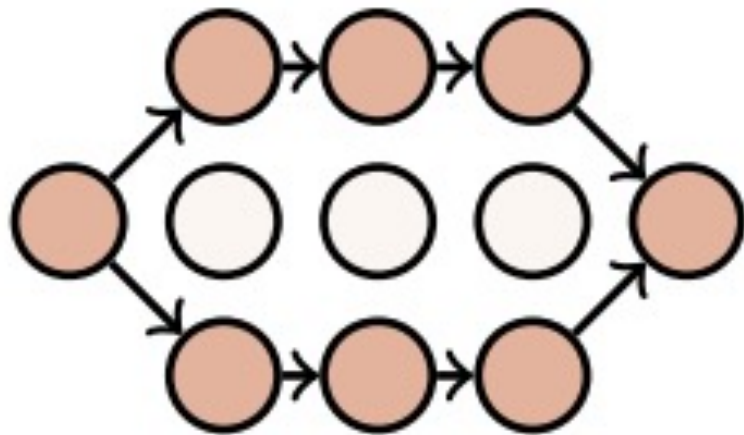
Can also solve by self-play via regret minimization on DAGs

Example: Patrolling Games

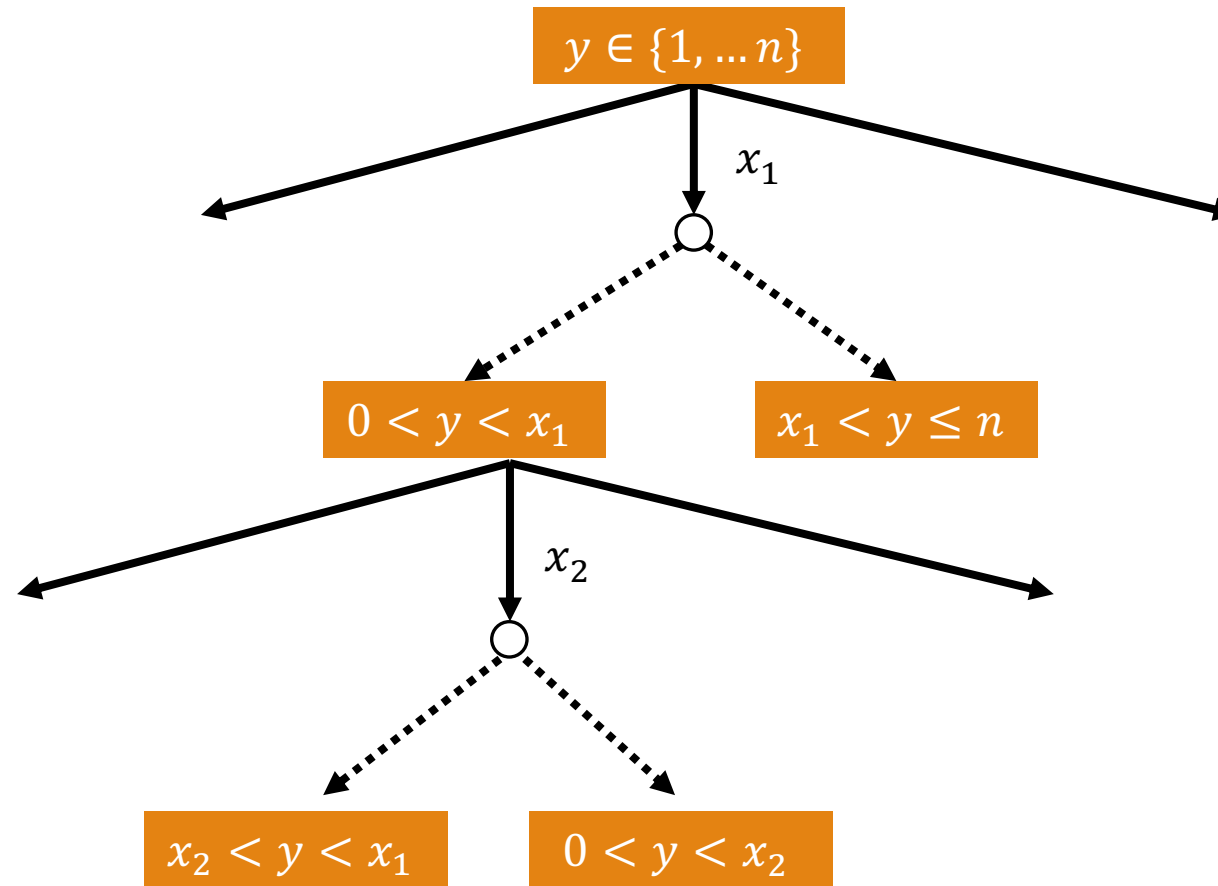
Let $G_1 = (V, E_1)$, $G_2 = (V, E_2)$, both DAGs with single source, sink

P1, P2 chooses a s-t path independently

- Utility = number of edges that are shared in both graphs



Example: Guess-The-Number



Note: Multiple sequences of guesses can give same info set
Decision points are of the form (t, lo, hi)
Total of $O(n^2 \cdot T)$ of them, each with $O(n)$ actions

Minimizing regret over scaled extensions

$$X = \left(\left(\left(\{1\} \triangleleft (f^{(1)}, m^{(1)}) \right) \triangleleft (f^{(2)}, m^{(2)}) \right) \triangleleft (f^{(3)}, m^{(3)}) \right) \dots$$

Idea: use simplex as building block

- Suppose we already have a regret minimizer over X , how to construct regret minimizer over $X \triangleleft (f, m)$?
- Tools we have now
 - A black box regret minimizer over X that supports:
 - Observe (AKA gradient)
 - Recommend()
 - A regret minimizer over simplexes (our building block)

Minimizing regret over scaled extensions

$$X \triangleleft (\mathbf{f}, m) = \left\{ \begin{pmatrix} x \\ s \end{pmatrix} \in X \times R_{\geq 0}^m \mid \langle \mathbf{1}, \mathbf{s} \rangle = \langle \mathbf{f}, \mathbf{x} \rangle \right\}$$

At timestep t (not the index that the scaled extension)

Recommend_Y()

- $x^t \leftarrow \text{Recommend}_X()$
- $s^t \leftarrow \text{Recommend}_\Delta()$
- play $y^t \leftarrow (x^t, s^t \cdot \langle \mathbf{f}, x^t \rangle)$ // Scale s^t based on x^t and \mathbf{f}

Observe_Y(u^t)

- $(u_X^t, u_\Delta^t) \leftarrow u^t$
- *Observe_Δ*(u_Δ^t)
- *Observe_X*($u_X^t + \langle u_\Delta^t, s_\Delta^t \rangle \cdot \mathbf{f}$) // Account for gradients from the descendants

For CFR: this is just the downward scaling of strategies for treeplex

For CFR: propagating future rewards upwards

Why does this work?

$$\begin{aligned}
 R_Y^T &= \max_{x,s} \sum_{t=1}^T \langle u_X^t, x \rangle + \langle u_\Delta^t, s \cdot \langle f, x \rangle \rangle - \langle u_X^t x^t \rangle - \langle u_\Delta^t, s^t \cdot \langle f, x^t \rangle \rangle \leftarrow \\
 &= \max_{x,s} \sum_{t=1}^T \langle u_X^t, x \rangle + \langle f, x \rangle \langle u_\Delta^t, s \rangle - \langle u_X^t x^t \rangle - \langle f, x^t \rangle \langle u_\Delta^t, s^t \rangle \\
 &= \max_x \left\{ \sum_{t=1}^T \langle u_X^t, x \rangle - \langle u_X^t x^t \rangle - \langle f, x^t \rangle \langle u_\Delta^t, s^t \rangle + \langle f, x \rangle \sum_{t=1}^T \max_s \langle u_\Delta^t, s \rangle \right\} \\
 &= \max_x \left\{ \sum_{t=1}^T \langle u_X^t, x \rangle - \langle u_X^t x^t \rangle - \langle f, x^t \rangle \langle u_\Delta^t, s^t \rangle + \langle f, x \rangle \left(R_\Delta^T + \sum_{t=1}^T \langle u_\Delta^t, s^t \rangle \right) \right\} \\
 &= \max_x \left\{ \sum_{t=1}^T \langle u_X^t, x \rangle + \langle f, x \rangle \langle u_\Delta^t, s^t \rangle - \langle u_X^t x^t \rangle - \langle f, x^t \rangle \langle u_\Delta^t, s^t \rangle + \langle f, x \rangle R_\Delta^T \right\} \\
 &\leq \max_x \left\{ \sum_{t=1}^T \langle u_X^t, x \rangle + \langle f, x \rangle \langle u_\Delta^t, s^t \rangle - \langle u_X^t x^t \rangle - \langle f, x^t \rangle \langle u_\Delta^t, s^t \rangle \right\} + \max_x \{ \langle f, x \rangle R_\Delta^T \} \\
 &\leq \max_x \left\{ \sum_{t=1}^T \langle u_X^t + f \cdot \langle u_\Delta^t, s^t \rangle, x \rangle - \langle u_X^t + f \cdot \langle u_\Delta^t, s^t \rangle, x^t \rangle \right\} + \max_x \{ \langle f, x \rangle \} [R_\Delta^T]^+
 \end{aligned}$$

Recommend_Y()

- $x^t \leftarrow \text{Recommend}_X()$
- $s^t \leftarrow \text{Recommend}_\Delta()$
- play $y^t \leftarrow (x^t, s^t \cdot \langle f, x^t \rangle)$

Observe_Y(u^t)

- $(u_X^t, u_\Delta^t) \leftarrow u^t$
- *Observe_Δ*(u_Δ^t)
- *Observe_X*($u_X^t + \langle u_\Delta^t, s^t \rangle \cdot f$)

Constant, equal
to 1 for Treeplex

Combining everything

$$X = \left(\left(\left(\{1\} \triangleleft (f^{(1)}, m^{(1)}) \right) \triangleleft (f^{(2)}, m^{(2)}) \right) \triangleleft (f^{(3)}, m^{(3)}) \right) \dots$$

$$R^T \leq J \cdot \sum_k R_k^T$$

Depends on f 's

Over all "simplexes" added
into the scaled extension.

Use any regret minimizer you
want! For RM, $\leq O(\sqrt{AT})$

For CFR, get $R^T \leq |\text{num sequences}| \sqrt{T}$

For DAGs, $R^T \leq |\text{num edges}| \sqrt{T}$

Many CFR-related speedups also work!

The rest follows from definition of saddlepoint residual, ϵ -Nash etc

Summary

Regret circuits are an easy recipe to construct new regret minimizers from small ones

- Can use **any** regret minimizer over the simplex! RM, RM+, Hedge or others not covered

Related work

- Regret circuits (Farina et. al., 2021)
 - Allows for operations like set intersection, albeit costly
- Kernelized Hedge (Farina et. al., 2024)
 - Earlier work Takimoto and Warmuth, 2003

An optimization angle

Games, Markets, and Online Learning (Kroer, 2025)

Online Convex Optimization

Let X be a compact convex set (e.g., simplex, treeplex)

Every iteration

- Choose $x_t \in X$
- Suffer loss $f_t(x_t)$
- f_t is restricted to be convex (we'll see more later)

Goal is to minimize regret

- $R_T = \sum_{t=1}^T f_t(x_t) - \min_x \sum_{t=1}^T f_t(x)$
- Similar to before: what loss we got minus what we could get in hindsight
- Note: you don't know what f_t is the future

Fact: Straightforward “follow-the-leader” algorithm does not work. Same argument as a long time ago

Idea: try to “regularize” things, smoothen things out (e.g., Hedge)

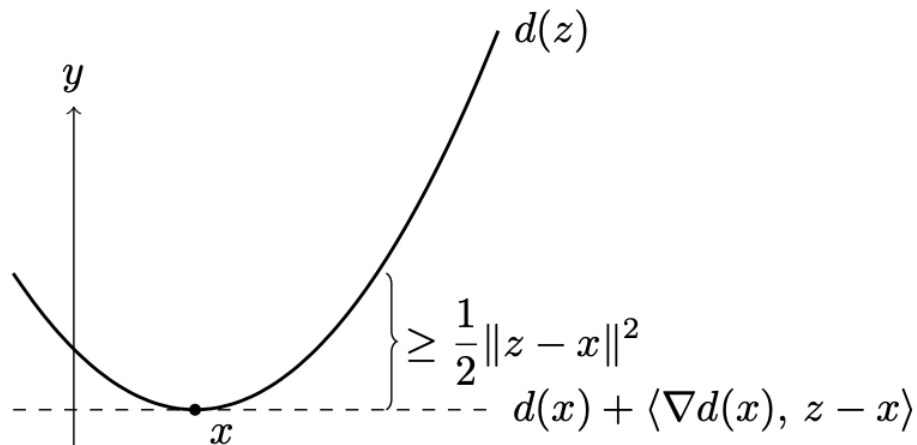
Background: Strong convexity

We will need to define

- A norm $\| \cdot \|$ in the primal space and a dual norm $\| \cdot \|_*$ in the *dual space*
- A **distance generating function** (DGF) d that is 1 –strongly convex w.r.t. $\| \cdot \|$

Strong convexity

- $d(z) \geq d(x) + \langle \nabla d(x), z - x \rangle + \frac{1}{2} \|z - x\|^2$ for all $x, z \in X$



Intuition: d is sufficiently “bent” upwards

Background: Bregman Divergences

Defined using 1-strongly convex DGF d

- $D(z||x) = d(z) - d(x) - \langle \nabla d(x), z - x \rangle$
- Always positive (why?), Not a metric (asymmetric)

Example:

- Euclidean DGF: $d(x) = 0.5\|x\|^2$ gives ℓ_2 distance $D(z||x) = \|x - z\|^2$
 - Strongly convex w.r.t. to ℓ_2 -norm
- Entropic DGF: $d(x) = \sum_i^n x_i \log x_i$ gives **KL divergence** $D(z||x) = \sum_i^n z_i (\log z_i / x_i)$
 - Strongly convex w.r.t ℓ_1 -norm (not 2 norm!)

Width w.r.t. d

- $\Omega_X = \max_{x,z} d(x) - d(z)$, basically the “diameter” of X

Dual norm $\|\cdot\|_*$

- $\|g\|_* = \max_{\|x\| \leq 1} \langle x, g \rangle$, Note: this is indeed a norm, dual of dual is original norm
- Dual norm of ℓ_2 is ℓ_2 norm (self dual). Dual norm of ℓ_1 is ℓ_∞ norm (max norm).

Online mirror descent

A general recipe for no-regret learning

$$x_{t+1} = \arg \min_{x \in X} \langle \eta \nabla f_t^{g_t}(x_t), x \rangle + D(x || x_t)$$

Intuition: we are selecting next iterate greedily based on previous gradient observed, but smoothing out using $D(\cdot || x_t)$

Example: Entropic DGF (leads to KL divergence)

- $x_{t+1} = \arg \min_{x \in X} \langle \eta g_t, x \rangle + D(x || x_t) = \frac{x \exp(-\eta g_t)}{\sum_i x_i \exp(-\eta g_{t,i})}$, (prove it!)
- $w_{t+1,i} = w_{t,i} \exp(-\eta g_{t,i}) \Rightarrow p_t = \frac{w_t \exp(-\eta g_t)}{\sum_i w_{t,i} \exp(-\eta g_{t,i})}$

Hedge is exactly OMD with entropic DGF!

- Note: using Euclidean DGF on 2-norm gives **projected gradient descent**
 - Weaker bounds on #actions typically \sqrt{n} instead of $\log(n)$ dependence
 - Euclidean distance is not the right “geometry” over the simplex

OMD: analysis

Theorem 4.6 *The OMD algorithm with DGF d achieves the following bound on regret:*

$$R_T \leq \frac{D(x||x_1)}{\eta} + \frac{\eta}{2} \sum_{t \in [T]} \|g_t\|_*^2.$$

Recall we had for Hedge: $R_T \leq \frac{\log(n)}{\eta} + \frac{\eta T}{2}$

Just need to bound $D(x||x_1)$ and $\|g_t\|_*^2$

Minimize by selecting x_1 appropriately
to be at the “center” of X

Maximum payoff per action

Optimizing for η yields the same bound!

What about treeplexes?

We now have a general method for regret minimization

- DGF which is 1-strongly convex in some norm $\|\cdot\|$

The **dilated entropy** works well here

- Add in weighted entropy regularization on the **behavioral** strategies at each info set
- $\sum_I \beta_i \cdot x_{\text{par}(I)} d(x_I / x_{\text{par}(I)})$ Interesting fact: applying dilated entropy in a particular way on **treeplexes** gives the QRE of the **reduced normal form matrix game**

Still convex (projection preserves convexity)

- Treeplex norm (slightly different notation)

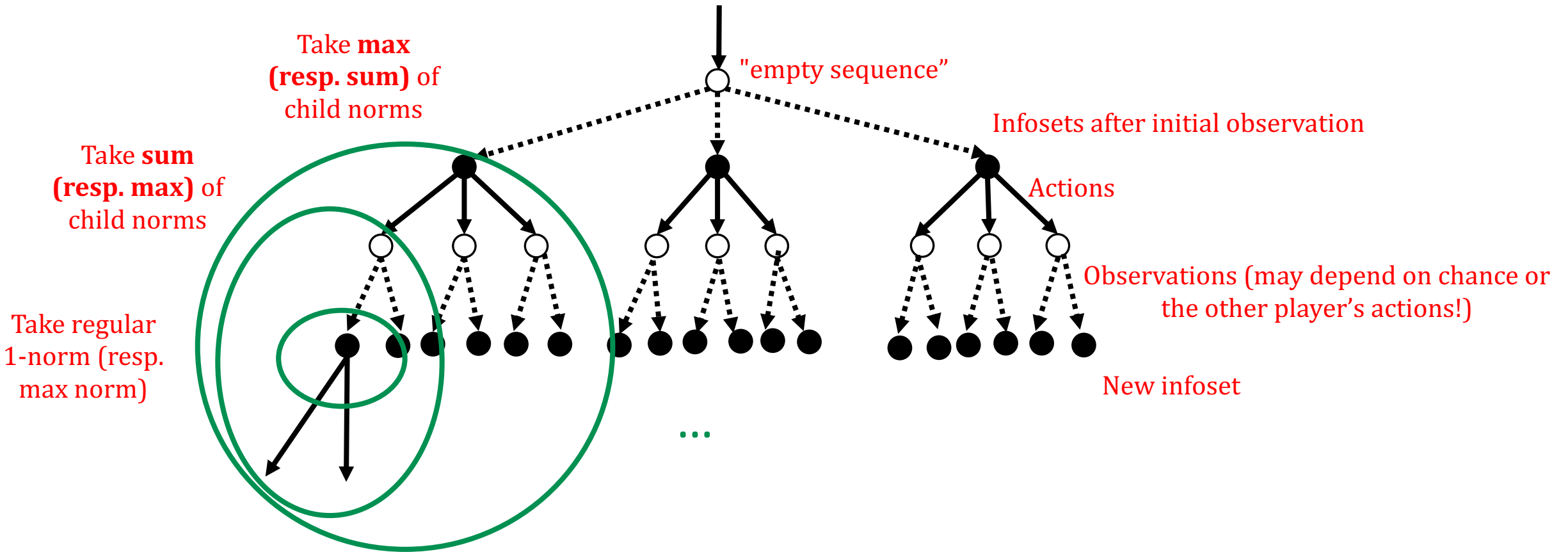
$$\|g^j\|_{X^j,1} = \begin{cases} |g_j| & \text{if } j \in \mathcal{L}_X \\ \sum_{a \in A_j} \|g^{ja}\|_{X^{ja},1} & \text{if } j \in \mathcal{I}_X \\ \max_{o \in O_j} \|g^{jo}\|_{X^{jo},1} & \text{if } j \in \mathcal{O}_X \end{cases}$$

If at “leaf info sets”, use regular 1-norm

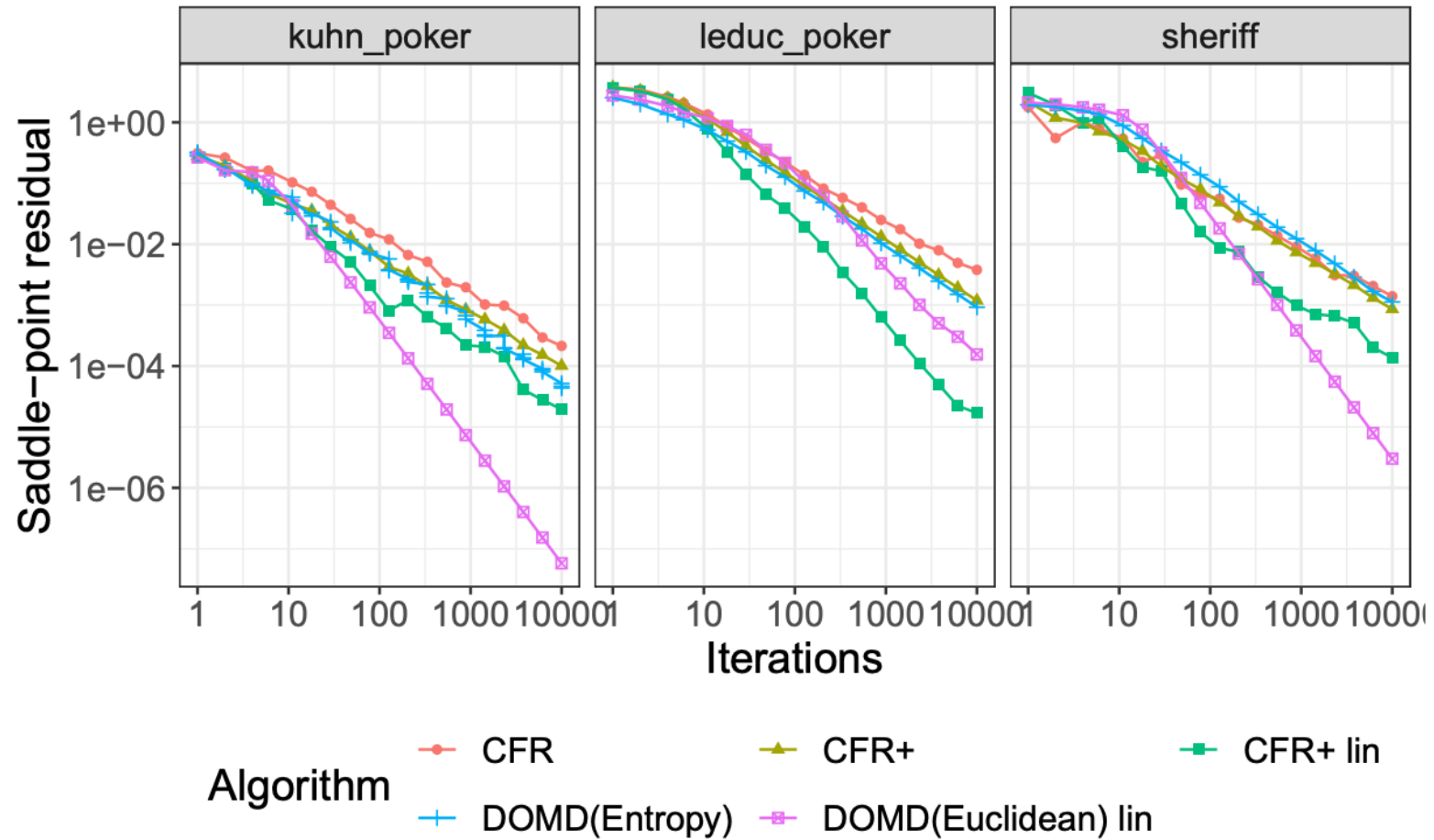
At non-leaf info sets, recurse and take sums

At observation nodes, recurse and take max (note this is different from CFR)

Illustration on treeplexes



For dual norms, just swap max and sums!



Source: Games, Markets, and Online Learning (Kroer, 2025)

Summary of OMD

OMD gives yet another tool for us to come out with nice regret minimizers

- Need to think about good DGF, norm
- Balance between diameter term Ω and magnitude of dual norm
- Dilated entropy is actually optimal for treeplexes in some sense (Fan, Kroer, Farina 2024)
- Doesn't have to apply to just treeplexes though!

Empirically

- Works well in practice for non-poker games compared to CFR
- Dilated entropy regularization is nice, taking prox only involves traversing treeplex and doing softmax at every stage (not solving a complicated optimization problem)
- Note that this looks different from RM/RM+
- Speedups like optimism also work!

Incremental Strategy Generation

General observations

Equilibrium tends to be sparse, whether in simultaneous move or extensive form games

- In real world, many strategies are just bad, no matter what opponent does
- Examples, attacking your own units in Starcraft is almost always bad
- “Junk” actions blow up size of the game
- Many reasonable actions often have zero probability to be played at eqm

Solving game may be hard when it is big, but best responses may be “easier” to compute!

*There exist games where BR is hard, but solving is easy. These are very rare

Single Oracle

Imagine game matrix is tall and skinny

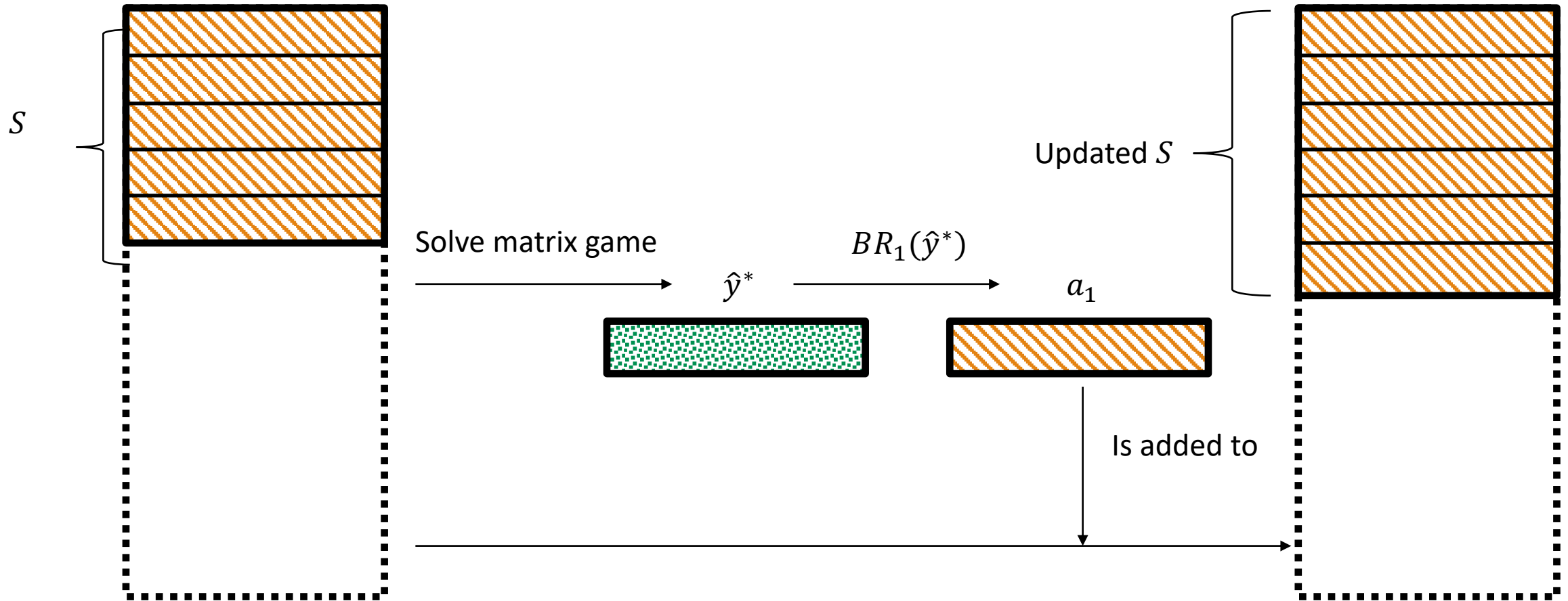
- P1 has many actions (more than we can enumerate, typically exponential in game description), P2 very few

Assume we have some way of finding best response $BR_1(y)$

Single Oracle

- 1. Initialize some random small set of actions in $S \subset A_1$
- 2. Iterate between
 - (i) Solving subgame induced by $S \rightarrow$ NE of \hat{x}^*, \hat{y}^*
 - (ii) Find deterministic $a_1 \in BR_1(\hat{y}^*)$ and add a_1 to S
 - (iii) Terminate if a_1 performs not a lot better than \hat{x}^* against \hat{y}^*
- At termination, we get an approximate eqm (why?)
- Will terminate eventually since all actions get added
 - Final S much smaller than A_1 typically

From an optimization standpoint, essentially column/row generation



Double Oracle

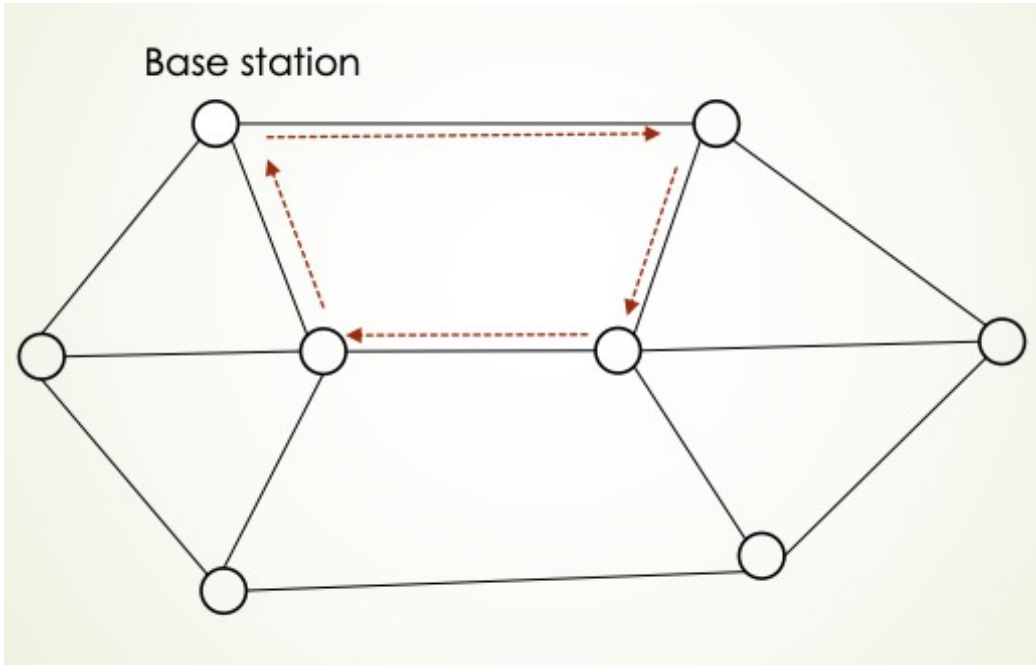
Same as single oracle, but repeated for both players!

- Needs best response oracle for both players

Double oracle

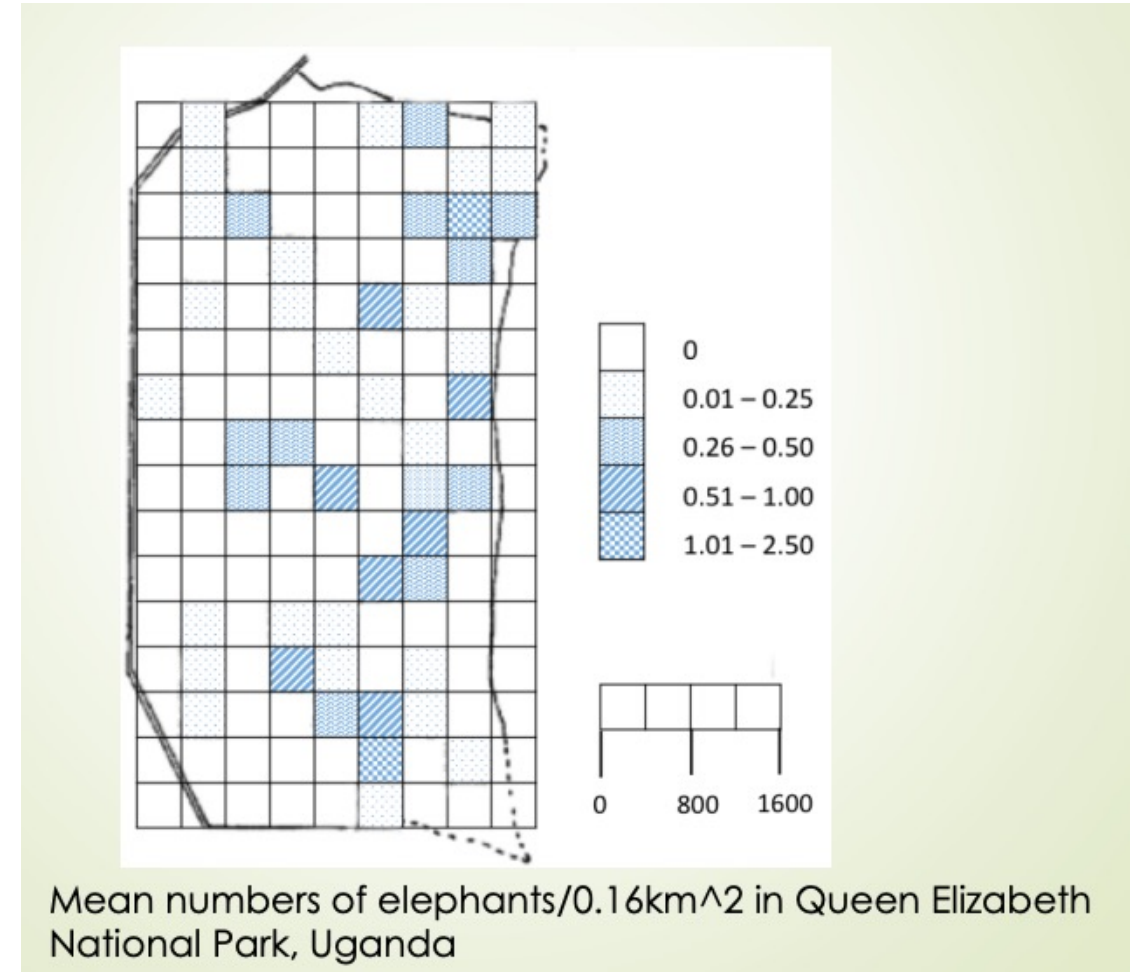
- 1. Initialize some random small set of strategies in $S_1 \subset A_1, S_2 \subset A_2$
- 2. Iterate between
 - (i) Solving subgame induced by $S_1, S_2 \rightarrow$ NE of \hat{x}^*, \hat{y}^*
 - (ii) Find deterministic $a_1 \in BR_1(\hat{y}^*)$ and add a_1 to S_1
 - (iii) Find deterministic $a_2 \in BR_2(\hat{x}^*)$ and add a_2 to S_2
 - (iii) Terminate if a_1 performs not much better than \hat{x}^* against \hat{y}^* **and** if a_2 performs not much better than \hat{y}^* against \hat{x}^* (i.e., saddle point residual). If small, terminate
- At termination, get an approximate eqm (saddle point residual is small)

Application: Anti-poaching Patrolling



Exponential number of paths!

Best response is at NP-hard (Hamiltonian path), can be solved via integer programming

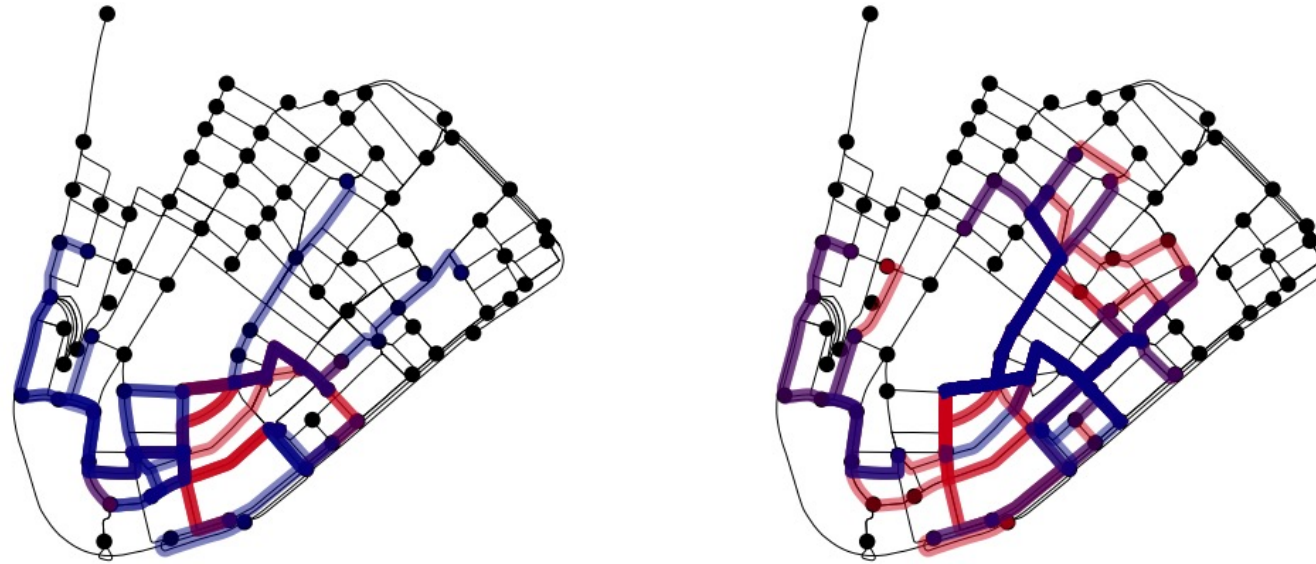


https://www.andrew.cmu.edu/user/feif/Publications/2016_GameSec_GraphContraction_slides.pdf

Application: Patrolling in General

Professor pursuing a student over T steps

- **Each time** professor meets student, student will be assigned work
- Path taken is fixed once chosen



Example on a map of
Manhattan with $T = 17$

Policy Space Response Oracle

Games may not have efficient BR

- BR oracles need not be exact
- Can use your favourite **single-agent** solver to find best response
- Example: PPO, or any deep RL approach

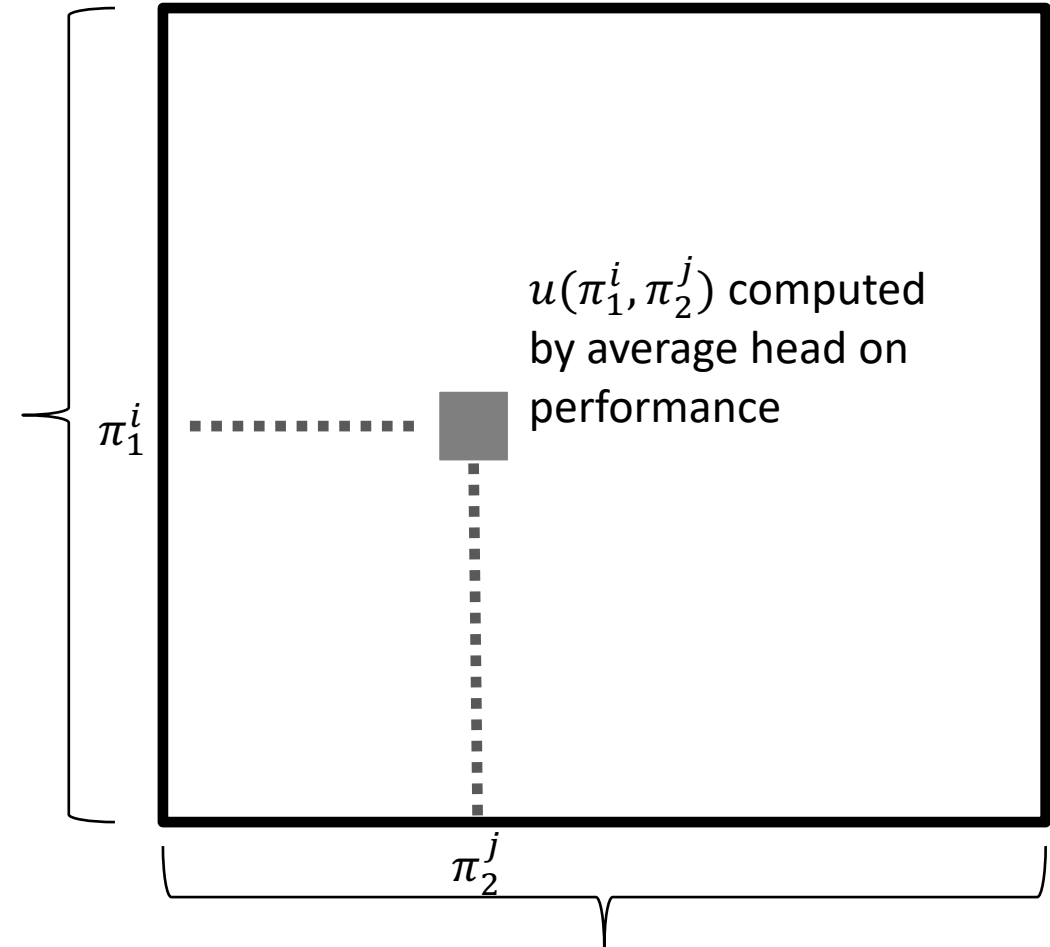
Maintain a *population* of agents for each player

- Technically, each agent does not need to be deterministic

Find **distribution** of agents to use

- During test time, choose one NN agent based on Nash eqm and follow it throughout

P1's trained Neural Networks

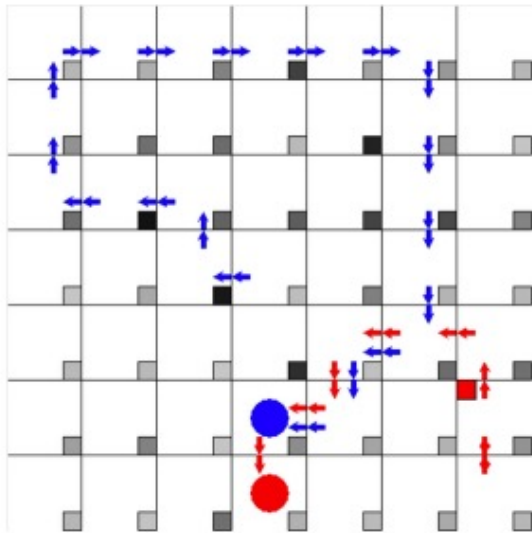


P2's trained Neural Networks

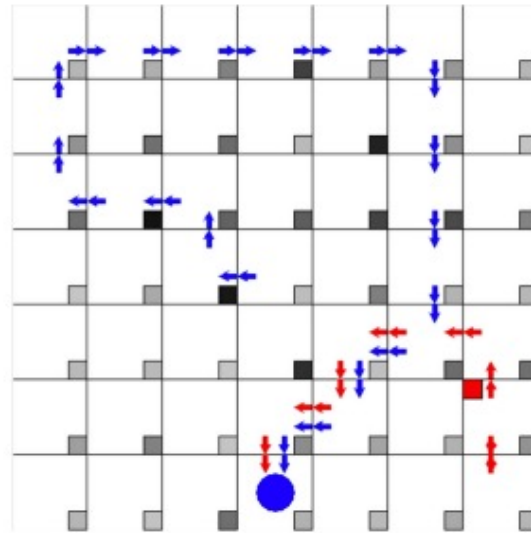
Example: Anti-poaching

Like before, but patrollers can observe “footprints” of poachers

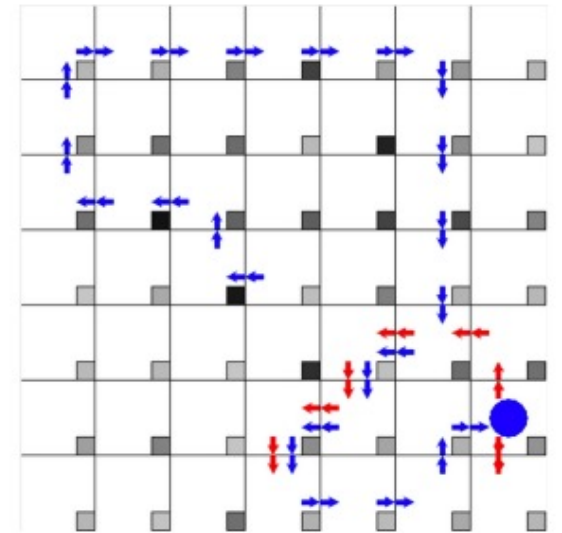
Gigantic if we represent it as a EFG!



(a) At time step 18



(b) At time step 19



(c) At time step 23

My personal thoughts

In practice, PSRO/Double Oracle super useful

- Get something off-the shelf for very large games without going into the weeds of game theory, just need best/approximate/better response oracle!

Simply applying PSRO isn't very exciting (YMMV) academically

- Single-agent RL isn't necessarily the most ideal best response oracle
- Unless applied to some kind of structured setting
- Some interesting theoretical work exists
- Special extensions to EFGs, pipelining etc.

PSRO can be quite misleading if not evaluated properly!

- **In practice, very rarely** do we get a strategy that has low exploitability
- But what are the alternatives?