

Approachability in unknown games

Gilles Stoltz (CNRS — HEC Paris)

Joint work with

Shie Mannor (Technion)

Vianney Perchet (Univ. Paris-Diderot — INRIA)



Statement of the problem

Regret can be minimized whether the game is known or not.

Can also approachability theory be extended to unknown games?

Classical approachability theory (Blackwell, 1956)

Finite sets of actions \mathcal{X} and \mathcal{Y} , actions taken $x_t \in \mathcal{X}$ and $y_t \in \mathcal{Y}$

Payoff function $r : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$

Aim: closed convex set \mathcal{C} , with first player to ensure

$$\frac{1}{T} \sum_{t=1}^T r(x_t, y_t) \longrightarrow \mathcal{C}$$

and second player to prevent this convergence

Characterization: $\forall \mathbf{y} \in \Delta(\mathcal{Y}), \exists \mathbf{x} \in \Delta(\mathcal{X}) : r(\mathbf{x}, \mathbf{y}) \in \mathcal{C}$

If this condition fails, the smallest approachable blow-up of \mathcal{C} is

$$\mathcal{C}_{\alpha_{\text{unif}}} = \left\{ c' \in \mathbb{R}^d : d_2(c', \mathcal{C}) \leq \alpha_{\text{unif}} \right\}$$

where $\alpha_{\text{unif}} = \max_{\mathbf{y} \in \Delta(\mathcal{Y})} \min_{\mathbf{x} \in \Delta(\mathcal{X})} d_2(r(\mathbf{x}, \mathbf{y}), \mathcal{C})$

Classical approachability theory (Blackwell, 1956)

Closed convex set \mathcal{C} and $\forall \mathbf{y} \in \Delta(\mathcal{Y}), \exists \mathbf{x} \in \Delta(\mathcal{X}) : r(\mathbf{x}, \mathbf{y}) \in \mathcal{C}$

Associated strategy

- Compute $\bar{\mathbf{c}}_t = \Pi_{\mathcal{C}}(\bar{\mathbf{r}}_t)$, the projection onto \mathcal{C} of

$$\bar{\mathbf{r}}_t = \frac{1}{t} \sum_{s=1}^t r(\mathbf{x}_s, \mathbf{y}_s)$$

- Draw $\mathbf{x}_{t+1} \sim \mathbf{x}_{t+1}$ such that

$$\forall \mathbf{y} \in \Delta(\mathcal{Y}), \quad \langle \bar{\mathbf{r}}_t - \bar{\mathbf{c}}_t, r(\mathbf{x}_{t+1}, \mathbf{y}) - \bar{\mathbf{c}}_t \rangle \leq 0$$

What does the player need to know/observe?

- **Bandit monitoring** enough: observe $r(\mathbf{x}_t, \mathbf{y}_t)$, not necessarily \mathbf{y}_t
- **Game** r needs to be **known** in general

Can the **game** be **unknown**?

Hope arises from the special case of **regret minimization**

Action sets \mathcal{X} and $[0, 1]^{\mathcal{X}}$

At each round and simultaneously,

- player 1 chooses $x_t \in \mathcal{X}$,
- player 2 picks $(g_{x',t})_{x' \in \mathcal{X}}$

Player 1 gets $g_{x_t,t}$ and ensures

$$\frac{1}{T} \sum_{t=1}^T g_{x_t,t} - \max_{x' \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T g_{x',t} \longrightarrow \mathbb{R}_-$$

No underlying game $g : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ needs to exist!

Can the **game** be **unknown**?

Hope arises from the special case of **regret minimization**

Player 1 ensures

$$\frac{1}{T} \sum_{t=1}^T g_{x_t, t} - \max_{x' \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T g_{x', t} \longrightarrow \mathbb{R}_-$$

Proof (full monitoring, known game)

Consider the vector payoff function

$$r(x_t, y_t) = (g(x_t, y_t) - g(x', y_t))_{x' \in \mathcal{X}}$$

Convex set to approach $\mathcal{C} = (\mathbb{R}_-)^{\mathcal{X}}$

Solution: choose \mathbf{x}_{t+1} proportional to $(\bar{r}_t)_+$, then

$$\forall \mathbf{y} \in \Delta(\mathcal{Y}), \quad \langle \bar{r}_t - \bar{c}_t, r(\mathbf{x}, \mathbf{y}) - \bar{c}_t \rangle = \left\langle (\bar{r}_t)_+, r(\mathbf{x}_{t+1}, \mathbf{y}) + (\bar{r}_t)_- \right\rangle = 0$$

because of the definition of r as a vector of differences

Remark: Structure of the game not used, the game could be **unknown**

Can the **game** be **unknown**?

Hope arises from the special case of **regret minimization**

Player 1 ensures

$$\frac{1}{T} \sum_{t=1}^T g_{x_t, t} - \max_{x' \in \mathcal{X}} \frac{1}{T} \sum_{t=1}^T g_{x', t} \longrightarrow \mathbb{R}_-$$

Proof (bandit monitoring, unknown game)

Consider the unbiased estimates

$$\tilde{g}_{x', t} = \frac{g_{x', t}}{\mathbf{x}_t(x')} \mathbb{I}_{\{x_t = x'\}}$$

and the associated vector payoff

$$\tilde{r}_t = (\tilde{g}_{x_t, t} - \tilde{g}_{x', t})_{x' \in \mathcal{X}}$$

Convex set to approach $\mathcal{C} = (\mathbb{R}_-)^{\mathcal{X}}$: again, doable even without knowing the structure

Approachability theory for unknown games

Statement of the problem

At each round and simultaneously,

- player 1 draws $x_t \in \mathcal{X}$ according to $\mathbf{x}_t \in \Delta(\mathcal{X})$
- player 2 picks $\mathbf{m}_t = (m_{x',t})_{x' \in \mathcal{X}} \in K$,

where $K \subset (\mathbb{R}^d)^{\mathcal{X}}$ is compact

Aim: force (player 1) or prevent (player 2) convergence of

$$\bar{r}_T = \frac{1}{T} \sum_{t=1}^T m_{x_t, t}$$

to **some** neighborhood of \mathcal{C}

To do:

- Indicate the targeted neighborhood
- Provide a strategy for player 1

Approachability theory for unknown games

First answers: probably not the end of the story...

Our COLT'14 paper — one solution to the problem

Recall that $\bar{r}_T = \frac{1}{T} \sum_{t=1}^T m_{x_t, t}$

Target sets of the form $\mathcal{C}_\varphi(\bar{\mathbf{m}}_T)$ where $\bar{\mathbf{m}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{m}_t$

Ensure $d_2(\bar{r}_T, \mathcal{C}_\varphi(\bar{\mathbf{m}}_T)) \rightarrow 0$

Why a function of the mean $\bar{\mathbf{m}}_T$ and not of the entire path?

Maybe too restrictive; and we will see that we anyway need to decompose $\bar{\mathbf{m}}_T$

What do we get by **calibration**?

On “average” good predictions $\widehat{\mathbf{m}}_t$ of the \mathbf{m}_t
(up to some grouping rounds according to the values of the $\widehat{\mathbf{m}}_t$)

Randomized strategy given by $\Psi : K \rightarrow \Delta(\mathcal{X})$

Average payoff close to

$$\frac{1}{T} \sum_{t=1}^T \Psi(\widehat{\mathbf{m}}_t) \odot \mathbf{m}_t$$

where $\mathbf{x} \odot \mathbf{m} = \mathbb{E}[m_X]$ when $X \sim \mathbf{x}$

But

- some grouping is needed (because of calibration)
- the guarantee needs to hold along the whole path (for all T)

What do we get by **calibration**?

Final guarantee in terms of convex decompositions:

$$\varphi(\bar{\mathbf{m}}_T) = \sup \left\{ d_2 \left(\sum_i \lambda_i \Psi(\mathbf{m}^{(i)}) \odot \mathbf{m}^{(i)}, \mathcal{C} \right) : \sum_i \lambda_i \mathbf{m}^{(i)} = \bar{\mathbf{m}}_T \right\}$$

Still a big problem to solve:

Which Ψ should be chosen?

There can be “**compensations**” and there are sometimes better choices than

$$\Psi(\mathbf{m}) \in \arg \min_{\mathbf{x} \in \Delta(\mathcal{X})} d_2(\mathbf{x} \odot \mathbf{m}, \mathcal{C})$$

An efficient strategy

We tackled the efficiency issue and offer a strategy that

- minimizes some regret in rounds of lengths 1, 2, 3, ...;
- only calls Ψ once in a round;
- performs no projection;
- ensures $d_2\left(\bar{r}_T, \mathcal{C}_\varphi(\bar{\mathbf{m}}_T)\right) = \mathcal{O}(T^{-1/4})$

This strategy has consequences in classical approachability as well!

1. Classical approachability without projecting onto \mathcal{C}

I.e., with $\varphi \equiv 0$ as a target

Only by exploiting the dual condition

$$\Psi(\mathbf{m}) \odot \mathbf{m} \in \mathcal{C} \quad \text{where} \quad \Psi(\mathbf{m}) \in \arg \min_{\mathbf{x} \in \Delta(\mathcal{X})} d_2(\mathbf{x} \odot \mathbf{m}, \mathcal{C})$$

Note: By a clever trick using that \mathcal{C} is approachable, Bernstein and Shimkin (2014) recover the classical $\mathcal{O}(T^{-1/2})$ rate

2. Convergence to the smallest approachable expansion $\mathcal{C}_{\alpha_{\text{unif}}}$ of a known game

Without even knowing it! Same principle

Does not solve the NP-hard optimization problem of determining

$$\alpha_{\text{unif}} = \max_{\mathbf{y} \in \Delta(\mathcal{Y})} \min_{\mathbf{x} \in \Delta(\mathcal{X})} d_2(r(\mathbf{x}, \mathbf{y}), \mathcal{C})$$