

# Simons Tutorial: Online Learning and Bandits Part I

---

**Wouter Koolen** and Alan Malek

August 31st, 2020

# Positioning this Tutorial

- Building up tools in support of RL
- Exploring surrounding viewpoints, problems and methods
- Soaking up “Culture”

# Working Definitions

**Context:** interactive decision making in unknown environment

**Aim:** Design systems to **amass reward** in **many environments**.

# Working Definitions

**Context:** interactive decision making in unknown environment

**Aim:** Design systems to **amass reward** in **many environments**.

**Main distinction:** model of environment

- **Reinforcement Learning** action affects **future state**
- **Bandits** action affects **observation**
- **Full Inf. Online Learning** action affects **reward**

# On the Menu

Two parts:

- (1) Full Information Online Learning
- (2) Bandits (w. Alan Malek)

# Full Information Online Learning

## 1. Two Basic Problems

Online Convex Optimisation; Online Gradient Descent

The Experts Problem; Exponential Weights

## 2. Two Peeks Beyond the Basics

Follow the Regularised Leader and Mirror Descent

Online Quadratic Optimisation; Online Newton Step

## 3. Applications

Classical Optimisation

Stochastic Optimisation

Saddle Points in Two-player Zero-Sum Games

## 4. Conclusion and Extensions

## **Two Basic Problems**

---

# Setup

- Focus on losses (negative rewards)
- Model Environment as Adversary
- Online Convex Optimisation (OCO) abstraction.



# OCO Problem

## Protocol: Online Convex Optimisation

Given: game length  $T$ , convex action space  $\mathcal{U} \subseteq \mathbb{R}^d$

For  $t = 1, 2, \dots, T$ ,

- The learner picks action  $w_t \in \mathcal{U}$
- The adversary picks convex loss  $f_t : \mathcal{U} \rightarrow \mathbb{R}$
- The learner observes  $f_t$  ◁ full information
- The learner incurs loss  $f_t(w_t)$

# OCO Problem

## Protocol: Online Convex Optimisation

Given: game length  $T$ , convex action space  $\mathcal{U} \subseteq \mathbb{R}^d$

For  $t = 1, 2, \dots, T$ ,

- The learner picks action  $w_t \in \mathcal{U}$
- The adversary picks convex loss  $f_t : \mathcal{U} \rightarrow \mathbb{R}$
- The learner observes  $f_t$  < full information
- The learner incurs loss  $f_t(w_t)$

The goal: control the **regret** (w.r.t. the best point after  $T$  rounds)

$$\mathcal{R}_T = \sum_{t=1}^T f_t(w_t) - \min_{u \in \mathcal{U}} \sum_{t=1}^T f_t(u)$$

using a computationally **efficient** algorithm for learner.

# Design Principle

Learner needs to “chase” the best point  $\arg \min_{u \in \mathcal{U}} \sum_{t=1}^T f_t(w_t)$ .  
But doing so naively **overfits**.

Idea: add regularisation. Two manifestations:

- Penalise excentricity “FTRL style”
- Update iterates, but only slowly “MD style”

Will see examples of both. For our purposes, these are roughly equivalent

# Online Gradient Descent (OGD) Algorithm

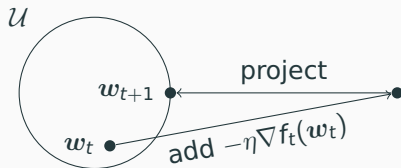
Let  $\mathcal{U}$  be a convex set containing  $0$ . Fix a learning rate  $\eta > 0$ .

## Algorithm: Online Gradient Descent (OGD)

OGD with learning rate  $\eta > 0$  plays

$$w_1 = 0 \quad \text{and} \quad w_{t+1} = \Pi_{\mathcal{U}}(w_t - \eta \nabla f_t(w_t))$$

where  $\Pi_{\mathcal{U}}(w) = \arg \min_{u \in \mathcal{U}} \|u - w\|$  is the projection onto  $\mathcal{U}$ .



**Figure 1:** OGD update

# Online Gradient Descent Result

## Algorithm: OGD

$$w_1 = 0 \quad \text{and} \quad w_{t+1} = \Pi_{\mathcal{U}}(w_t - \eta \nabla f_t(w_t))$$

## Assumption: Boundedness

Bounded **domain**  $\max_{u \in \mathcal{U}} \|u\| \leq D$  and **gradients**  $\|\nabla f_t(w_t)\| \leq G$ .

# Online Gradient Descent Result

## Algorithm: OGD

$$w_1 = 0 \quad \text{and} \quad w_{t+1} = \Pi_{\mathcal{U}}(w_t - \eta \nabla f_t(w_t))$$

## Assumption: Boundedness

Bounded **domain**  $\max_{u \in \mathcal{U}} \|u\| \leq D$  and **gradients**  $\|\nabla f_t(w_t)\| \leq G$ .

## Theorem (OGD regret bd, Zinkevich 2003)

$$\mathcal{R}_T = \sum_{t=1}^T f_t(w_t) - \min_{u \in \mathcal{U}} \sum_{t=1}^T f_t(u) \leq \frac{1}{2\eta} D^2 + \frac{\eta}{2} T G^2$$

# Online Gradient Descent Result

## Algorithm: OGD

$$w_1 = 0 \quad \text{and} \quad w_{t+1} = \Pi_{\mathcal{U}}(w_t - \eta \nabla f_t(w_t))$$

## Assumption: Boundedness

Bounded **domain**  $\max_{u \in \mathcal{U}} \|u\| \leq D$  and **gradients**  $\|\nabla f_t(w_t)\| \leq G$ .

## Theorem (OGD regret bd, Zinkevich 2003)

$$\mathcal{R}_T = \sum_{t=1}^T f_t(w_t) - \min_{u \in \mathcal{U}} \sum_{t=1}^T f_t(u) \leq \frac{1}{2\eta} D^2 + \frac{\eta}{2} T G^2$$

## Corollary

*Tuning  $\eta = \frac{D}{G\sqrt{T}}$  results in  $\mathcal{R}_T \leq DG\sqrt{T}$ .*

# Online Gradient Descent Result

## Algorithm: OGD

$$w_1 = 0 \quad \text{and} \quad w_{t+1} = \Pi_{\mathcal{U}}(w_t - \eta \nabla f_t(w_t))$$

## Assumption: Boundedness

Bounded **domain**  $\max_{u \in \mathcal{U}} \|u\| \leq D$  and **gradients**  $\|\nabla f_t(w_t)\| \leq G$ .

## Theorem (OGD regret bd, Zinkevich 2003)

$$\mathcal{R}_T = \sum_{t=1}^T f_t(w_t) - \min_{u \in \mathcal{U}} \sum_{t=1}^T f_t(u) \leq \frac{1}{2\eta} D^2 + \frac{\eta}{2} T G^2$$

## Corollary

*Tuning  $\eta = \frac{D}{G\sqrt{T}}$  results in  $\mathcal{R}_T \leq DG\sqrt{T}$ .*

Sublinear regret: learning overhead per round  $\rightarrow 0$ .



## Proof of OGD regret bound

Using convexity, we may analyse the tangent upper bound

$$f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \leq \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle$$

Moreover,

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 &= \|\Pi_{\mathcal{U}}(\mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)) - \mathbf{u}\|^2 \\ &\leq \|\mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t) - \mathbf{u}\|^2 \\ &= \|\mathbf{w}_t - \mathbf{u}\|^2 - 2\eta \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle + \eta^2 \|\nabla f_t(\mathbf{w}_t)\|^2 \end{aligned}$$

Hence

$$\langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle \leq \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta} + \frac{\eta}{2} \|\nabla f_t(\mathbf{w}_t)\|^2$$

## Proof of OGD regret bound (ctd)

Summing over  $T$  rounds, we find

$$\begin{aligned}\mathcal{R}_T^u &\leq \sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle \\ &\leq \underbrace{\sum_{t=1}^T \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta}}_{\text{telescopes}} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\mathbf{w}_t)\|^2 \\ &\leq \frac{\|\mathbf{u}\|^2 - \cancel{\|\mathbf{w}_{T+1} - \mathbf{u}\|^2}}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\mathbf{w}_t)\|^2 \\ &\leq \frac{D^2}{2\eta} + \frac{\eta}{2} TG^2\end{aligned}$$

## OCO Lower Bound

Is OGD regret bound of  $\mathcal{R}_T \leq GD\sqrt{T}$  any good?

## OCO Lower Bound

Is OGD regret bound of  $\mathcal{R}_T \leq GD\sqrt{T}$  any good?

Scaling with  $G$  and  $D$  is natural. What about  $\sqrt{T}$ ?

# OCO Lower Bound

Is OGD regret bound of  $\mathcal{R}_T \leq GD\sqrt{T}$  any good?

Scaling with  $G$  and  $D$  is natural. What about  $\sqrt{T}$ ?

## **Theorem**

*Any OCO algorithm can be made to incur  $\mathcal{R}_T = \Omega(\sqrt{T})$ .*

# OCO Lower Bound

Is OGD regret bound of  $\mathcal{R}_T \leq GD\sqrt{T}$  any good?

Scaling with  $G$  and  $D$  is natural. What about  $\sqrt{T}$ ?

## Theorem

*Any OCO algorithm can be made to incur  $\mathcal{R}_T = \Omega(\sqrt{T})$ .*

## Proof (by probabilistic argument).

Consider interval  $\mathcal{U} = [-1, 1]$  and linear losses  $f_t(u) = x_t \cdot u$  with i.i.d. Rademacher coefficients  $x_t \in \{\pm 1\}$ . Any algorithm has expected loss zero. The expected loss of the best action ( $\pm 1$ ) is  $-\mathbb{E}[|\sum_{t=1}^T x_t|] = -\Omega(\sqrt{T})$ . Then as the expected regret is  $\mathbb{E}[\mathcal{R}_T] = \Omega(\sqrt{T})$ , there is a deterministic witness.  $\square$

Here, the regret arises from *overfitting* of the best point.

## OGD Discussion

- Adversarial result, super strong!
- Proof reveals it is really about linear losses.
- Matching lower bounds

Successful in practise:

- Practically **all deep learning** uses versions of online gradient descent (e.g. TensorFlow has AdaGrad [Duchi et al., 2011]) even though objective not convex.

# From Learning Parameters to Picking Actions

We now turn to the second elementary online learning task.

- Decision Theoretic Online Learning
- Experts setting (also: Hedge setting)
- Prediction with Expert Advice



# From Learning Parameters to Picking Actions

We now turn to the second elementary online learning task.

- Decision Theoretic Online Learning
- Experts setting (also: Hedge setting)
- Prediction with Expert Advice

## Protocol: Prediction With Expert Advice

Given: game length  $T$ , number  $K$  of experts

For  $t = 1, 2, \dots, T$ ,

- Learner chooses a distribution  $w_t \in \Delta_K$  on  $K$  “experts”.
- Adversary reveals loss vector  $\ell_t \in [0, 1]^K$ .
- Learner’s loss is the **dot loss**  $w_t^\top \ell_t = \sum_{k=1}^K w_t^k \ell_t^k$

# From Learning Parameters to Picking Actions

We now turn to the second elementary online learning task.

- Decision Theoretic Online Learning
- Experts setting (also: Hedge setting)
- Prediction with Expert Advice

## Protocol: Prediction With Expert Advice

Given: game length  $T$ , number  $K$  of experts

For  $t = 1, 2, \dots, T$ ,

- Learner chooses a distribution  $w_t \in \Delta_K$  on  $K$  “experts”.
- Adversary reveals loss vector  $\ell_t \in [0, 1]^K$ .
- Learner’s loss is the **dot loss**  $w_t^\top \ell_t = \sum_{k=1}^K w_t^k \ell_t^k$

The goal: control the **regret** (w.r.t. the best expert after  $T$  rounds)

$$\mathcal{R}_T = \sum_{t=1}^T w_t^\top \ell_t - \min_{k \in [K]} \sum_{t=1}^T \ell_t^k$$

using a computationally **efficient** algorithm for learner.

## Let's apply what we know

Observations:

- Dot loss  $u \mapsto u^\top \ell_t$  is *linear* (hence convex).
- Gradient  $\ell_t \in [0, 1]^K$  bounded by  $\|\ell_t\| \leq \sqrt{K}$ .
- Probability simplex  $\Delta_K$  is contained in unit ball.

So: Instance of Online Convex Optimisation.

OGD with  $D = 1$  and  $G = \sqrt{K}$  gives  $\mathcal{R}_T \leq \sqrt{KT}$ .

## Let's apply what we know

Observations:

- Dot loss  $u \mapsto u^T \ell_t$  is *linear* (hence convex).
- Gradient  $\ell_t \in [0, 1]^K$  bounded by  $\|\ell_t\| \leq \sqrt{K}$ .
- Probability simplex  $\Delta_K$  is contained in unit ball.

So: Instance of Online Convex Optimisation.

OGD with  $D = 1$  and  $G = \sqrt{K}$  gives  $\mathcal{R}_T \leq \sqrt{KT}$ .

Q: **Optimal?**

## Let's apply what we know

Observations:

- Dot loss  $u \mapsto u^T \ell_t$  is *linear* (hence convex).
- Gradient  $\ell_t \in [0, 1]^K$  bounded by  $\|\ell_t\| \leq \sqrt{K}$ .
- Probability simplex  $\Delta_K$  is contained in unit ball.

So: Instance of Online Convex Optimisation.

OGD with  $D = 1$  and  $G = \sqrt{K}$  gives  $\mathcal{R}_T \leq \sqrt{KT}$ .

Q: **Optimal?**

Maybe not. There are no points with loss difference  $\sqrt{K}$  in the simplex ...

# Exponential Weights / Hedge Algorithm

## Algorithm: Exponential Weights (EW)

EW with *learning rate*  $\eta > 0$  plays weights in round  $t$ :

$$w_t^k = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}}. \quad (\text{EW})$$

# Exponential Weights / Hedge Algorithm

## Algorithm: Exponential Weights (EW)

EW with *learning rate*  $\eta > 0$  plays weights in round  $t$ :

$$w_t^k = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}}. \quad (\text{EW})$$

or, equivalently,  $w_1^k = \frac{1}{K}$  and

$$w_{t+1}^k = \frac{w_t^k e^{-\eta \ell_t^k}}{\sum_{j=1}^K w_t^j e^{-\eta \ell_t^j}} \quad (\text{EW, incremental})$$

# Exponential Weights / Hedge Algorithm

## Algorithm: Exponential Weights (EW)

EW with *learning rate*  $\eta > 0$  plays weights in round  $t$ :

$$w_t^k = \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}}. \quad (\text{EW})$$

or, equivalently,  $w_1^k = \frac{1}{K}$  and

$$w_{t+1}^k = \frac{w_t^k e^{-\eta \ell_t^k}}{\sum_{j=1}^K w_t^j e^{-\eta \ell_t^j}} \quad (\text{EW, incremental})$$

## Theorem (EW regret bd, Freund and Schapire 1997)

The regret of EW is bounded by  $\mathcal{R}_T \leq \frac{\ln K}{\eta} + T \frac{\eta}{8}$ .

## Corollary

Tuning  $\eta = \sqrt{\frac{8 \ln K}{T}}$  yields  $\mathcal{R}_T \leq \sqrt{T/2 \ln K}$ .



# EW Analysis

Applying *Hoeffding's Lemma* to the loss of each round gives

$$\sum_{t=1}^T w_t^\top \ell_t \leq \underbrace{\sum_{t=1}^T \left( \frac{-1}{\eta} \ln \left( \sum_{k=1}^K w_t^k e^{-\eta \ell_t^k} \right) \right)}_{\text{"mix loss"}} + \underbrace{\eta/8}_{\text{overhead}}$$

Crucial observation is that cumulative mix loss *telescopes*

$$\begin{aligned} \sum_{t=1}^T \frac{-1}{\eta} \ln \left( \sum_{k=1}^K w_t^k e^{-\eta \ell_t^k} \right) &= \sum_{t=1}^T \frac{-1}{\eta} \ln \left( \sum_{k=1}^K \frac{e^{-\eta \sum_{s=1}^{t-1} \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}} e^{-\eta \ell_t^k} \right) \\ &= \sum_{t=1}^T \frac{-1}{\eta} \ln \left( \frac{\sum_{k=1}^K e^{-\eta \sum_{s=1}^t \ell_s^k}}{\sum_{j=1}^K e^{-\eta \sum_{s=1}^{t-1} \ell_s^j}} \right) \\ &\stackrel{\text{telescopes}}{=} \frac{-1}{\eta} \ln \left( \sum_{k=1}^K e^{-\eta \sum_{t=1}^T \ell_t^k} \right) + \frac{\ln K}{\eta} \\ &\leq \min_{k \in [K]} \sum_{t=1}^T \ell_t^k + \frac{\ln K}{\eta}. \end{aligned}$$

# Summary so far

Balancing act: “model complexity” vs “overfitting”

## Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2}G^2T$$

## Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8}T$$

# Summary so far

Balancing act: “model complexity” vs “overfitting”

## Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2}G^2T$$

## Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8}T$$

Generates many follow-up questions:

- What if horizon  $T$  is not fixed? Anytime guarantees?
- What if gradient bound  $G$  is not known a priori?
- Can we have the actual gradient norms?
- What if model complexity ( $D$ ) is not known? Not uniformly bounded? See Orabona and Cutkosky ICML'20 tutorial.

# Summary so far

Balancing act: “model complexity” vs “overfitting”

## Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2} G^2 T$$

## Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8} T$$

Generates many follow-up questions:

- What if horizon  $T$  is not fixed? Anytime guarantees?
- What if gradient bound  $G$  is not known a priori?
- Can we have the actual gradient norms?
- What if model complexity ( $D$ ) is not known? Not uniformly bounded? See Orabona and Cutkosky ICML'20 tutorial.

Need refined analyses  $\Rightarrow$  Restarts (doubling trick), decreasing  $\eta_t$  (AdaGrad/AdaHedge), learning the learning rate  $\eta$  (MetaGrad),

...

# Summary so far

Balancing act: “model complexity” vs “overfitting”

## Theorem (OGD)

$$\mathcal{R}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2} G^2 T$$

## Theorem (EW)

$$\mathcal{R}_T \leq \frac{\ln K}{\eta} + \frac{\eta}{8} T$$

Generates many follow-up questions:

- What if horizon  $T$  is not fixed? Anytime guarantees?
- What if gradient bound  $G$  is not known a priori?
- Can we have the actual gradient norms?
- What if model complexity ( $D$ ) is not known? Not uniformly bounded? See Orabona and Cutkosky ICML'20 tutorial.

Need refined analyses  $\Rightarrow$  Restarts (doubling trick), decreasing  $\eta_t$  (AdaGrad/AdaHedge), learning the learning rate  $\eta$  (MetaGrad),

...

Active research area!

# **Two Peeks Beyond the Basics**

---

## FTRL/MD “sneak peek”

Q: What if my **domain** does not look like either ball or simplex?

# FTRL/MD “sneak peek”

Q: What if my **domain** does not look like either ball or simplex?

## Algorithm: Follow the Regularised Leader (FTRL)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t \langle u, \nabla f_s(w_s) \rangle + \frac{1}{\eta} R(u)$$

## Algorithm: Mirror Descent (MD)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \langle u, \nabla f_t(w_t) \rangle + \frac{1}{\eta} B(u \| w_t)$$



# FTRL/MD “sneak peek”

Q: What if my **domain** does not look like either ball or simplex?

## Algorithm: Follow the Regularised Leader (FTRL)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t \langle u, \nabla f_s(w_s) \rangle + \frac{1}{\eta} R(u)$$

## Algorithm: Mirror Descent (MD)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \langle u, \nabla f_t(w_t) \rangle + \frac{1}{\eta} B(u \| w_t)$$

Examples:

	<b>Regularizer <math>R</math></b>	<b>Bregman Divergence <math>B</math></b>
OGD	sq. Euclidean norm	sq. Euclidean distance
EW	Shannon entropy	Kullback-Leibler divergence

# FTRL/MD “sneak peek”

Q: What if my **domain** does not look like either ball or simplex?

## Algorithm: Follow the Regularised Leader (FTRL)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t \langle u, \nabla f_s(w_s) \rangle + \frac{1}{\eta} R(u)$$

## Algorithm: Mirror Descent (MD)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \langle u, \nabla f_t(w_t) \rangle + \frac{1}{\eta} B(u \| w_t)$$

Examples:

	<b>Regularizer <math>R</math></b>	<b>Bregman Divergence <math>B</math></b>
OGD	sq. Euclidean norm	sq. Euclidean distance
EW	Shannon entropy	Kullback-Leibler divergence

Other entropies: Burg, Tsallis, Von Neumann, ... Connections to continuous exponential weights [van der Hoeven et al., 2018].

# FTRL/MD “sneak peak” performance

## Algorithm: Follow the Regularised Leader (FTRL)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t \langle u, \nabla f_s(w_s) \rangle + \frac{1}{\eta} R(u)$$

## Algorithm: Mirror Descent

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \langle u, \nabla f_t(w_t) \rangle + \frac{1}{\eta} B(u \| w_t)$$

# FTRL/MD “sneak peak” performance

## Algorithm: Follow the Regularised Leader (FTRL)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t \langle u, \nabla f_s(w_s) \rangle + \frac{1}{\eta} R(u)$$

## Algorithm: Mirror Descent

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \langle u, \nabla f_t(w_t) \rangle + \frac{1}{\eta} B(u \| w_t)$$

## Theorem (AdaFTRL, Orabona and Pál 2015)

Fix a norm  $\|\cdot\|$  with associated dual norm  $\|\cdot\|_*$ . Let  $R : \mathcal{U} \rightarrow [0, D^2]$  be strongly convex w.r.t.  $\|\cdot\|$ . AdaFTRL ensures

$$\mathcal{R}_T \leq 2D \sqrt{\sum_{t=1}^T \|\nabla f_t(w_t)\|_*^2} + 2 \cdot \text{loss range.}$$

# Quadratic Losses

So far we used convexity to “linearise”

$$f_t(\mathbf{u}) \geq f_t(\mathbf{w}_t) + \langle \mathbf{u} - \mathbf{w}_t, \nabla f_t(\mathbf{w}_t) \rangle,$$

and our methods essentially operated on linear losses. But what if we **know there is curvature?**

- How to **represent/quantify** curvature?
- How to **efficiently** manipulate curvature?
- How much can we reduce the regret?

# Curvature assumptions

## Assumption: Quadratic loss lower bound

There is a matrix  $M_t \succeq 0$  such that

$$f_t(\mathbf{u}) \geq \underbrace{f_t(\mathbf{w}_t) + \langle \mathbf{u} - \mathbf{w}_t, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{2}(\mathbf{u} - \mathbf{w}_t)^\top M_t (\mathbf{u} - \mathbf{w}_t)}_{=: q_t(\mathbf{u})}$$

for each  $\mathbf{u} \in \mathcal{U}$ .

# Curvature assumptions

## Assumption: Quadratic loss lower bound

There is a matrix  $M_t \succeq 0$  such that

$$f_t(\mathbf{u}) \geq \underbrace{f_t(\mathbf{w}_t) + \langle \mathbf{u} - \mathbf{w}_t, \nabla f_t(\mathbf{w}_t) \rangle + \frac{1}{2}(\mathbf{u} - \mathbf{w}_t)^\top M_t (\mathbf{u} - \mathbf{w}_t)}_{=: q_t(\mathbf{u})}$$

for each  $\mathbf{u} \in \mathcal{U}$ .

Two main classes of instances

- squared Euclidean distance:  $f_t(\mathbf{u}) = \frac{1}{2} \|\mathbf{u} - \mathbf{x}_t\|^2$  satisfies the assumption with  $M_t = \mathbf{I}$ . More generally, **strongly convex** functions have  $M_t \propto \mathbf{I}$ .
- linear regression:  $f_t(\mathbf{u}) = (y_t - \langle \mathbf{u}, \mathbf{x}_t \rangle)^2$  satisfies the assumption with  $M_t = \mathbf{x}_t \mathbf{x}_t^\top$ . More generally, **exp-concave** functions have  $M_t \propto \nabla_t f_t(\mathbf{w}_t) \nabla_t f_t(\mathbf{w}_t)^\top$ .

## Algorithm: Online Newton Step (FTRL variant)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t q_s(u) + \frac{1}{2} \|u\|^2$$

Computing the iterate  $w_{t+1}$  amounts to minimising a convex quadratic. Often (depending on  $\mathcal{U}$ ) **closed-form solution** or **1d line search**.

- For  $M_t \propto I$ , takes  $O(d)$  per round.
- For rank-one  $M_t$ , can do update in  $O(d^2)$  per round.
- In both cases, need to take care of projection onto  $\mathcal{U}$ .



# ONS Performance

## Algorithm: Online Newton Step (FTRL version)

$$w_{t+1} = \arg \min_{u \in \mathcal{U}} \sum_{s=1}^t q_s(u) + \frac{1}{2} \|u\|^2$$

## Theorem (ONS strcvx bd, Hazan et al. 2006)

*For the strongly convex case  $M_t \propto I$ , ONS guarantees*

$$\mathcal{R}_T = O(\ln T)$$

*Algorithm reduces to OGD with specific decreasing step-size  $\eta_t$*

## Theorem (ONS expccv bd, Hazan et al. 2006)

*For the exp-concave case  $M_t \propto g_t g_t^T$ , ONS guarantees*

$$\mathcal{R}_T = O(d \ln T)$$

# ONS Discussion

- Convex quadratics closed under taking sums. Run-time independent of  $T$ .
- Curvature gives huge reduction in regret:  $\sqrt{T}$  to  $\ln T$ .
- Matrix **sketching** techniques allow trading off run-time  $O(d^2)$  vs  $O(d)$  with regret  $O(\ln T)$  vs  $O(\sqrt{T})$  [Luo et al., 2016].

# Applications

---

# Application 1: Offline Optimisation

## **Problem**

*Given gradient access to a convex  $f$ , find a near-optimal point.*

# Application 1: Offline Optimisation

## Problem

Given gradient access to a convex  $f$ , find a near-optimal point.

Idea: run OGD on  $f_t = f$  for  $T$  rounds. Regret bound gives

$$\sum_{t=1}^T f(w_t) - T \min_{u \in \mathcal{U}} f(u) \leq GD\sqrt{T}$$

We may divide by  $T$  and apply convexity to find

$$f\left(\frac{1}{T} \sum_{t=1}^T w_t\right) - \min_{u \in \mathcal{U}} f(u) \leq \frac{GD}{\sqrt{T}}$$

Find  $\epsilon$ -suboptimal point (iterate average) after  $T = \frac{G^2 D^2}{\epsilon^2}$  rounds.

# Application 1: Offline Optimisation

## Problem

Given gradient access to a convex  $f$ , find a near-optimal point.

Idea: run OGD on  $f_t = f$  for  $T$  rounds. Regret bound gives

$$\sum_{t=1}^T f(w_t) - T \min_{u \in \mathcal{U}} f(u) \leq GD\sqrt{T}$$

We may divide by  $T$  and apply convexity to find

$$f\left(\frac{1}{T} \sum_{t=1}^T w_t\right) - \min_{u \in \mathcal{U}} f(u) \leq \frac{GD}{\sqrt{T}}$$

Find  $\epsilon$ -suboptimal point (iterate average) after  $T = \frac{G^2 D^2}{\epsilon^2}$  rounds.

Why would we optimise this way? For example, what if  $f_t \rightarrow f$ .

## Application 2: Online to Batch

### Assumption: stochastic setting

Suppose training set  $f_1, \dots, f_T$  and test point  $f$  drawn i.i.d. from unknown  $\mathbb{P}$ .

### Problem

*Learn a point  $\hat{w}_T$  from the training set that generalises to  $\mathbb{P}$ , i.e. behaves like  $u^* = \arg \min_{u \in \mathcal{U}} \mathbb{E}_f[f(u)]$*

## Application 2: Online to Batch

### Assumption: stochastic setting

Suppose training set  $f_1, \dots, f_T$  and test point  $f$  drawn i.i.d. from unknown  $\mathbb{P}$ .

### Problem

Learn a point  $\hat{w}_T$  from the training set that generalises to  $\mathbb{P}$ , i.e. behaves like  $\mathbf{u}^* = \arg \min_{\mathbf{u} \in \mathcal{U}} \mathbb{E}_f[f(\mathbf{u})]$

Idea: use online learning algorithm on training set  $f_1, \dots, f_T$ , to get iterates  $w_1, \dots, w_T$ . Output the *average iterate estimator*

$$\hat{w}_T = \frac{1}{T} \sum_{t=1}^T w_t.$$

### Theorem

An online regret bound  $R_T \leq B(T)$  implies

$$\mathbb{E}_{iid f_1, \dots, f_T, f} [f(\hat{w}_T) - f(\mathbf{u}^*)] \leq \frac{B(T)}{T}$$



## Application 3: Computing Saddle Points

### **Assumption: convex-concave**

Fix an objective function

$$g(x,y)$$

convex in  $x$ , concave in  $y$ .

## Application 3: Computing Saddle Points

### Assumption: convex-concave

Fix an objective function

$$g(x, y)$$

convex in  $x$ , concave in  $y$ .

The game *value* is

$$V^* = \min_x \max_y g(x, y) = \max_y \min_x g(x, y).$$

An  $\epsilon$ -*saddle point*  $(\bar{x}, \bar{y})$  satisfies

$$V^* - \epsilon \leq \min_x g(x, \bar{y}) \leq V^* \leq \max_y g(\bar{x}, y) \leq V^* + \epsilon.$$

## Application 3: Computing Saddle Points

### Assumption: convex-concave

Fix an objective function

$$g(x, y)$$

convex in  $x$ , concave in  $y$ .

The game *value* is

$$V^* = \min_x \max_y g(x, y) = \max_y \min_x g(x, y).$$

An  $\epsilon$ -saddle point  $(\bar{x}, \bar{y})$  satisfies

$$V^* - \epsilon \leq \min_x g(x, \bar{y}) \leq V^* \leq \max_y g(\bar{x}, y) \leq V^* + \epsilon.$$

### Problem

Find an  $\epsilon$ -saddle point

## Application 3: Computing Saddle Points

### Assumption: convex-concave

Fix an objective function

$$g(x, y)$$

convex in  $x$ , concave in  $y$ .

The game *value* is

$$V^* = \min_x \max_y g(x, y) = \max_y \min_x g(x, y).$$

An  $\epsilon$ -saddle point  $(\bar{x}, \bar{y})$  satisfies

$$V^* - \epsilon \leq \min_x g(x, \bar{y}) \leq V^* \leq \max_y g(\bar{x}, y) \leq V^* + \epsilon.$$

### Problem

Find an  $\epsilon$ -saddle point

Idea: play regret minimisation algorithms for  $x$  and  $y$ .

## Application 3: Saddle Point Algorithm

### Algorithm: approximate saddle point solver

For  $t = 1, 2, \dots, T$

- Players play  $x_t$  and  $y_t$ .
- Players see loss functions  $x \mapsto +g(x, y_t)$  and  $y \mapsto -g(x_t, y)$ .

Output average iterate pair  $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$  and  $\bar{y}_T = \frac{1}{T} \sum_{t=1}^T y_t$

Assume the players have regret (bounds)  $\mathcal{R}_T^x$  and  $\mathcal{R}_T^y$ , i.e.

$$\sum_{t=1}^T +g(x_t, y_t) - \min_x \sum_{t=1}^T +g(x, y_t) \leq \mathcal{R}_T^x$$
$$\sum_{t=1}^T -g(x_t, y_t) - \min_y \sum_{t=1}^T -g(x_t, y) \leq \mathcal{R}_T^y$$

### Theorem (self-play, Freund and Schapire 1999)

$\bar{x}_T$  and  $\bar{y}_T$  form an  $\frac{\mathcal{R}_T^x + \mathcal{R}_T^y}{T}$ -saddle point.

## Application 3: Saddle Point Analysis

$$\begin{aligned}V^* &= \min_x \max_y g(x, y) \\&\leq \max_y g(\bar{x}_T, y) \\&\leq \max_y \frac{1}{T} \sum_{t=1}^T g(x_t, y) \\&\leq \frac{1}{T} \sum_{t=1}^T g(x_t, y_t) + \frac{\mathcal{R}_T^y}{T} \\&\leq \min_x \frac{1}{T} \sum_{t=1}^T g(x, y_t) + \frac{\mathcal{R}_T^x + \mathcal{R}_T^y}{T} \\&\leq \min_x g(x, \bar{y}_T) + \frac{\mathcal{R}_T^x + \mathcal{R}_T^y}{T} \\&\leq \min_x \max_y g(x, y) + \frac{\mathcal{R}_T^x + \mathcal{R}_T^y}{T} \\&= V^* + \frac{\mathcal{R}_T^x + \mathcal{R}_T^y}{T}\end{aligned}$$

# **Conclusion and Extensions**

---

## Conclusion

- Online Learning a powerful and versatile tool
- Environment-as-black-box. Adversarial.
- Foundation for optimisation, statistical learning, games, ...



# Conclusion

- Online Learning a powerful and versatile tool
- Environment-as-black-box. Adversarial.
- Foundation for optimisation, statistical learning, games, ...

Some (of many) cool things we left out:

- First-order (small loss) and second-order (small variance) bounds
- Adaptivity to friendly stochastic environments (best of both worlds, interpolation)
- Optimism (predicting the upcoming gradient)
- Non-stationarity (tracking, adaptive/dynamic regret, path length)
- Beyond convexity (star-convex, geometrically convex, ...)
- Supervised Learning and (stochastic) complexities (VC, Littlestone, Rademacher, ...)

Thanks!

# References

---

- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1):119–139, August 1997.
- Yoav Freund and Robert E Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1-2):79–103, 1999.
- Elad Hazan, Adam Kalai, Satyen Kale, and Amit Agarwal. Logarithmic regret algorithms for online convex optimization. In *Learning Theory*, pages 499–513, 2006.

- Haipeng Luo, Alekh Agarwal, Nicolò Cesa-Bianchi, and John Langford. Efficient second order online learning by sketching. In *Advances in Neural Information Processing Systems 29*, pages 902–910. 2016.
- Francesco Orabona and Dávid Pál. Scale-free algorithms for online linear optimization. In *Algorithmic Learning Theory*, pages 287–301, 2015.
- Dirk van der Hoeven, Tim van Erven, and Wojciech Kotłowski. The many faces of exponential weights in online learning. volume 75 of *Proceedings of Machine Learning Research*, pages 2067–2092, 06–09 Jul 2018.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, ICML'03, page 928–935, 2003.

# Simons Tutorial: Online Learning and Bandits Part II

---

Wouter Koolen and **Alan Malek**

August 31st, 2020

# **What is a Bandit?**

---

# The Basic Bandit Game

## Protocol: A Bandit Game

Given: game length  $T$ , action space  $\mathcal{A}$

For  $t = 1, 2, \dots, T$ ,

- The learner picks action  $A_t \in \mathcal{A}$
- The adversary simultaneously picks reward  $r_t \in \mathcal{A} \rightarrow [0, 1]$
- The learner observes and receives  $r_t(A_t)$
- The learner does not observe  $r_t(a)$  for  $a \neq A_t$

The goal: control the regret (a random variable)

$$\mathcal{R}_T = \underbrace{\max_{a \in \mathcal{A}} \sum_{t=1}^T r_t(a)}_{\text{Best action in hindsight}} - \sum_{t=1}^T r_t(A_t) \quad (1)$$

# Bandits as a super simple MDP

- $S = \{\text{the\_state}\}$ ,  $P(\text{the\_state}|\text{the\_state}, \mathbf{a}) = 1$
- Why should we care about this in RL?
  - Creates a tension between
    - Exploration (learning about the loss of actions)
    - Exploitation (playing actions that will have low regret)
  - Exploration/Exploitation is absent in full-information, present in RL
  - Model is simple enough to allow for comprehensive theory
  - Easily incorporates adversarial data
  - Useful algorithm design principles



# The Regret

$$\mathcal{R}_T = \underbrace{\max_{a \in \mathcal{A}} \sum_{t=1}^T r_t(a)}_{\text{Best action in hindsight}} - \sum_{t=1}^T r_t(A_t)$$

- $\mathcal{R}_T$  is a random variable we do not observe
- Different objectives, from easiest to hardest
  - Pseudo-regret  $\overline{\mathcal{R}}_T = \max_{a \in \mathcal{A}} \mathbb{E} \left[ \sum_{t=1}^T r_t(a) \right] - \mathbb{E} \left[ \sum_{t=1}^T r_t(A_t) \right]$
  - Expected regret  $\mathbb{E}[\mathcal{R}_T] = \mathbb{E} \left[ \underbrace{\max_{a \in \mathcal{A}} \sum_{t=1}^T r_t(a)}_{\text{can depend on } A_t} - \sum_{t=1}^T r_t(A_t) \right]$
  - High probability bounds on the realized regret
- We always have  $\overline{\mathcal{R}}_T \leq \mathbb{E}[\mathcal{R}_T]$
- If the adversary is *reactive*, then the distribution of  $r_t$  can be a function of  $A_1, \dots, A_{t-1}$
- Otherwise, the adversary is *oblivious* and  $\overline{\mathcal{R}}_T = \mathbb{E}[\mathcal{R}_T]$

# Our focus

- Introduce most popular bandit problems
  - Adversarial Bandits
  - Stochastic Bandits
  - Pure Exploration Bandits
  - Contextual Bandits (time permitting)
- Concentrate on useful algorithm design principles
  - Exponential weights (still useful)
  - Optimism in the face of Uncertainty
  - Probability matching (i.e. Thompson sampling)
  - Action-Elimination

## Other Settings that haven't been considered

- Data models for  $r_t$ 
  - chosen by an adversary
  - sampled i.i.d.
  - stochastic with adversarial perturbations...
- Action spaces
  - Finite number of arms
  - A vector space ( $r_t$  are functions)
  - Combinatorial (e.g. subsets, paths on a graph)
- Objectives
  - Pseudo-regret (the expectation over the learner's randomness)
  - Realized regret (with high probability)
  - Best-arm identification a.k.a. pure exploration
- Side information
  - Linear rewards
  - Competing with a policy class
- ...

# Adversarial Bandits

---

# Adversarial Protocol

## Protocol: Finite-Arm Adversarial Bandit Protocol

Given: game length  $T$ , number of arms  $K$

For  $t = 1, 2, \dots, T$ ,

- The learner picks action  $I_t \in \{1, \dots, K\}$
  - The adversary simultaneously picks losses  $\ell_t \in [0, 1]^K$
  - The learner observes and receives  $\ell_t(I_t)$
- 
- The results are easier to state using losses instead of rewards
  - Randomization of  $I_t$  is *essential*
  - We are familiar with adversarial data from the first half
  - The simple idea of estimating  $\ell_t$  from  $\ell_t(I_t)$  and then applying a full-information algorithm works very well

# Algorithm Design Principle: Exponential Weights

## Algorithm: Exp3 [Auer et al., 2002b]

Given: number of arms  $K$ , learning rate  $\eta > 0$ , length  $T$

Initialize  $p_0(i) = 1/K$ ,  $\hat{L}_0(i) = 0$  for all  $i \in [K]$

For  $t = 1, 2, \dots, T$ :

- Sample  $I_t \sim p_t$  and observe  $\ell_t(I_t)$
- Estimate  $\hat{\ell}_t(i) = \frac{\ell_t(I_t)}{p_t(I_t)} \mathbb{1}_{\{I_t=i\}}$  and  $\hat{L}_t = \hat{\ell}_t + \hat{L}_{t-1}$
- Calculate  $W_t = \sum_j e^{-\eta \hat{L}_t(j)}$  and  $p_t(i) = \frac{1}{W_t} e^{-\eta \hat{L}_t(i)}$

- Exp3 = Exponential Weights for Exploration and Exploitation
- $\hat{\ell}_t$  is the importance-weighted estimator of  $\ell_t$
- $\hat{\ell}_t$  is unbiased:

$$\mathbb{E}_{I_t \sim p_t}[\hat{\ell}_t(I_t)] = \mathbb{E} \left[ \frac{\ell_t(I_t)}{p_t(I_t)} \mathbb{1}_{\{I_t=i\}} \right] = \sum_i p_t(i) \frac{\ell_t(I_t)}{p_t(I_t)} \mathbb{1}_{\{I_t=i\}} = \ell_t(i).$$

- Exp3 runs exponential weights on  $\hat{\ell}_t$

## Exp3: Analysis

- Following the EW analysis,  $W_t$  is a potential function
- For any  $i^*$ ,  $e^{-\eta \hat{L}_T(i^*)} \leq \sum_j e^{-\eta \hat{L}_T(j)} = W_T = W_0 \prod_{t=1}^T \frac{W_t}{W_{t-1}}$ .

$$\begin{aligned} \frac{W_t}{W_{t-1}} &= \frac{\sum_j e^{-\eta \hat{L}_{t-1}(j)} e^{-\eta \hat{\ell}_t(j)}}{\sum_j e^{-\eta \hat{L}_{t-1}(j)}} = \sum_j p_{t-1}(j) e^{-\eta \hat{\ell}_t(j)} \\ &\leq \underbrace{\sum_j p_{t-1}(j) \left( 1 - \eta \hat{\ell}_t(j) + \frac{\eta^2}{2} \hat{\ell}_t(j)^2 \right)}_{\text{since } e^x \leq 1 + x + \frac{1}{2} x^2 \text{ for } x \leq 0} \\ &= 1 - \eta \sum_j p_t(j) \hat{\ell}_t(j) + \frac{\eta^2}{2} \sum_j p_t(j) \hat{\ell}_t(j)^2 \\ &\leq \underbrace{e^{-\eta \sum_j p_t(j) \hat{\ell}_t(j) + \frac{\eta^2}{2} \sum_j p_t(j) \hat{\ell}_t(j)^2}}_{\text{since } 1 + x \leq e^x} \end{aligned}$$

## Exp3: Analysis

- Computing the telescope,

$$e^{-\eta \hat{L}_T(i^*)} \leq W_0 \prod_{t=1}^T \frac{W_t}{W_{t-1}} \leq K \prod_{t=1}^T e^{-\eta \sum_j p_t(j) \hat{\ell}_t(j) + \frac{\eta^2}{2} \sum_j p_t(j) \hat{\ell}_t(j)^2}$$

- We rearrange, divide by  $\eta$ , and take the log:

$$\begin{aligned} \sum_{t=1}^T \sum_j p_t(j) \hat{\ell}_t(j) - \hat{L}_T(i^*) &\leq \frac{\log(K)}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \sum_j p_t(j) \hat{\ell}_t(j)^2 \\ &= \frac{\log(K)}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \sum_j p_t(j) \frac{\ell_t(j)^2}{p_t(j)^2} \mathbb{1}_{\{I_t=j\}} \\ &\leq \frac{\log(K)}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \sum_j \frac{\mathbb{1}_{\{I_t=j\}}}{p_t(I_t)}. \end{aligned}$$



## Exp3: Analysis

- Take the expectation

$$\begin{aligned}\bar{\mathcal{R}}_T &\leq \mathbb{E} \left[ \sum_{t=1}^T \sum_j p_t(j) \hat{\ell}_t(j) - \hat{L}_T(i^*) \right] \\ &\leq \frac{1}{\eta} \log(K) + \frac{\eta}{2} \sum_{t=1}^T \mathbb{E} \left[ \sum_{t=1}^T \sum_j \frac{\mathbb{1}_{\{I_t=j\}}}{p_t(I_t)} \right] \\ &\leq \frac{1}{\eta} \log(K) + \frac{\eta}{2} TK.\end{aligned}$$

### Theorem (Exp3 upper bound [Auer et al., 2002b])

With  $\eta = \sqrt{\frac{2 \log(T)}{TK}}$ , Exp3 has  $\bar{\mathcal{R}}_T \leq \sqrt{2TK \log(K)}$ .

Only get pseudo-Regret bounds because the  $i^*$  in the proof was fixed, not a function of  $I_1, \dots, I_T$ .

# Lower Bounds

## Theorem (Adversarial Bandits lower bound [Auer et al., 2002b])

*Any adversarial bandit algorithm must have*

$$\bar{\mathcal{R}}_T = \Omega(\sqrt{TK})$$

- Exp3 upper bound:  $\bar{\mathcal{R}}_T \leq \sqrt{2TK \log(K)}$
- First matching upper bound achieved by INF [Audibert and Bubeck, 2009] (which is Mirror Descent)

# Upgrades

- High Probability bounds: requires a lower-variance estimate of  $\hat{\ell}_t$  or an algorithm that keeps  $p_t(i)$  away from zero
  - Exp3.P [Auer et al., 2002b] uses  $\hat{\ell}_t(i) = \frac{\mathbb{1}_{\{I_t=i\}} \ell_t(I_t) - \beta}{p_t(I_t)}$
  - Exp3-IX [Neu, 2015] uses  $\hat{\ell}_t(i) = \frac{\mathbb{1}_{\{I_t=i\}} \ell_t(I_t)}{p_t(I_t) + \gamma}$
- Experts with bandits; each arm is an expert that recommends actions and you compete with the best expert (Exp4 algorithm) [Auer et al., 2002b]
- Competing with strategies that can switch [Auer, 2002]
- Feedback determined by a graph [Mannor and Shamir, 2011]
- Partial Monitoring [Bartók et al., 2014]
- Combinatorial action spaces...

# Stochastic Bandits

---

# Protocol

## Protocol: Stochastic Bandits Protocol

Given: game length  $T$ , number of arms  $K$

Assume unknown reward distributions  $\nu_1, \dots, \nu_K$

For  $t = 1, 2, \dots, T$ ,

- The learner picks action  $I_t \in \{1, \dots, K\}$
  - The learner observes and receives reward  $X_t \sim \nu_{I_t}$
- 
- Stochastic bandits is an old problem [Thompson, 1933]
  - We will use the following notation
    - Reward of arm  $i$  is sampled from  $\nu_i$  with  $\mu_i := \mathbb{E}_{X \sim \nu_i}[X]$
    - $i^* = \arg \max_j \mu_j$  is the best arm
    - Gaps  $\Delta_i := \mu^* - \mu_i \geq 0$ ,
    - Number of pulls  $N_{i,t} := \sum_{s=1}^t \mathbb{1}_{\{I_s=i\}}$
    - Empirical mean  $\hat{\mu}_{i,t} := \frac{\sum_{s=1}^t X_s \mathbb{1}_{\{I_s=i\}}}{N_{i,t}}$

# Protocol

## Protocol: Stochastic Bandits Protocol

Given: game length  $T$ , number of arms  $K$

Assume unknown reward distributions  $\nu_1, \dots, \nu_K$

For  $t = 1, 2, \dots, T$ ,

- The learner picks action  $I_t \in \{1, \dots, K\}$
  - The learner observes and receives reward  $X_t \sim \nu_{I_t}$
- 
- We still want to minimize the expected regret, which has the useful decomposition

$$\mathbb{E}[\mathcal{R}_T] = T\mu_{i^*} - \sum_{t=1}^T \mathbb{E}[X_t] = \sum_i \Delta_i \mathbb{E}[N_{i,T}]$$

## Assumption: 1-sub-Gaussian reward distributions

For all stochastic bandit problems, we will assume that all arms are 1-sub-Gaussian, i.e.  $\mathbb{E}_{X \sim \mu_i} [e^{\lambda(X - \mu_i)^2 - \lambda^2/2}] \leq 1$ . This implies

$$P(\hat{\mu}_{i,t} - \mu_i \geq \epsilon) \leq e^{-\frac{\epsilon^2 t}{2}}.$$

## Warm-up: Explore-Than-Commit

### Algorithm: Explore-Than-Commit

Given: Game length  $T$ , exploration parameter  $M$

For  $t = 1, 2, \dots, MK$ :

- Choose  $i_t = (t \bmod K)$ , see  $X_t \sim \nu_{i_t}$

Compute empirical means  $\hat{\mu}_{i, MK}$

For  $t = MK + 1, MK + 2, \dots, T$ :

- Pull arm  $i = \arg \max_j \hat{\mu}_{j, MK}$

- The first strategy you might try
- A proof idea that we will return to: bound regret by first bounding  $\mathbb{E}[N_{i, T}]$ .
- In this simple algorithm,

$$\mathbb{E}[N_{i, T}] = M + (T - MK)P \left( i = \arg \max_j \hat{\mu}_{j, MK} \right)$$

# Explore-than-Commit Upper Bound

Using the sub-Gaussian concentration bound,

$$\begin{aligned} P\left(i = \arg \max_j \hat{\mu}_{j, MK}\right) &\leq P(\hat{\mu}_{i, MK} \geq \hat{\mu}_{i^*, MK}) \\ &= P((\hat{\mu}_{i, MK} - \mu_i) \geq (\hat{\mu}_{i^*, MK} - \mu_{i^*}) + \Delta_i) \\ &\leq e^{-\frac{M\Delta_i^2}{4}} \text{ (the difference is } \sqrt{2/M}\text{-sub-Gaussian)} \end{aligned}$$

## Theorem (Explore-than-Commit upper bound)

$$\mathbb{E}[\mathcal{R}_T] = \sum_i \Delta_i \mathbb{E}[N_{i, T}] \leq \sum_{i=1}^K \Delta_i \left( M + (T - MK)e^{-\frac{M\Delta_i^2}{4}} \right)$$

- If we know  $\Delta$ , then  $m = \frac{4}{\Delta_1^2} \log \frac{T\Delta_1^2}{4}$ , results in  $\mathbb{E}[\mathcal{R}_T] \leq \sum_{i=1}^K \frac{4}{\Delta_1} \log \frac{T\Delta_1^2}{4} + T \frac{4}{T\Delta_1^2} = O\left(\frac{K \log(T)}{\Delta_1}\right)$
- But we don't know  $\Delta$ ...can we be adaptive?



# Algorithm Design Principle: OFU

- OFU: Optimism in the Face of Uncertainty
- We establish some confidence set for the problem instance (e.g. means) to within some confidence set
- We then assume the most favorable instance in the confidence set and act greedily

## Algorithm: UCB1 [Auer et al., 2002a]

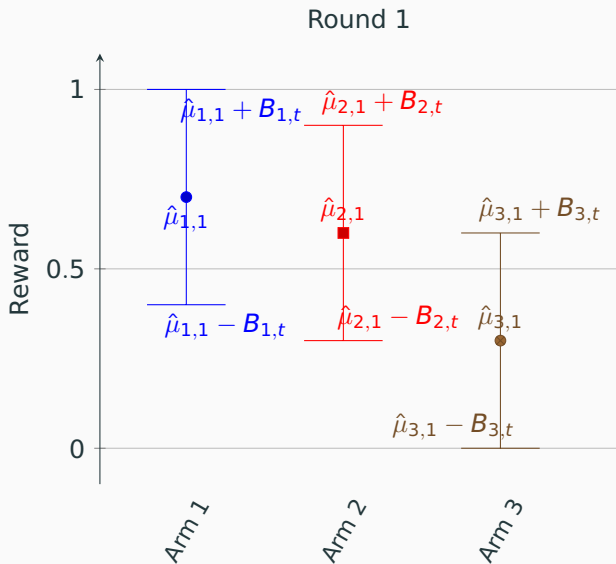
Given: Game length  $T$

Initialize: play every arm once

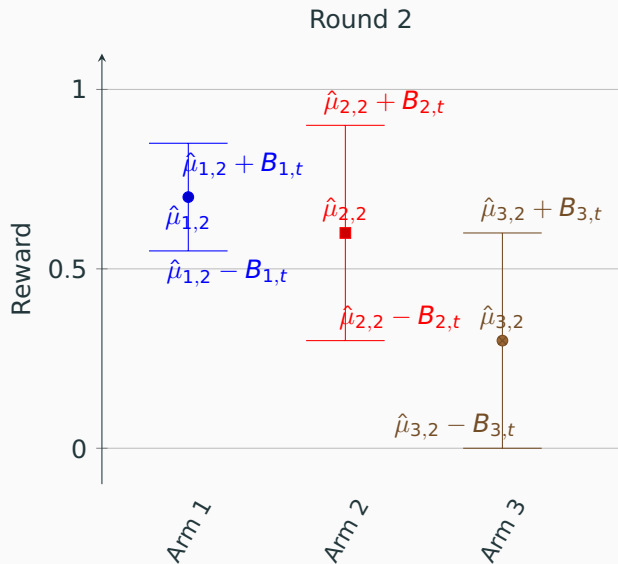
For  $t = K + 1, 2, \dots, T$ :

- Compute upper confidence bounds  $B_{i,t-1} = \sqrt{\frac{6 \log(t)}{N_{i,t-1}}}$
- Choose  $I_t = \arg \max_j \hat{\mu}_{j,t-1} + B_{j,t-1}$ , observe  $X_t \sim \nu_{I_t}$
- Update  $N_{i,t} = N_{i,t-1} + \mathbb{1}_{\{I_t=i\}}$  and  $\hat{\mu}_{i,t} = \frac{\sum_{s=1}^t \mathbb{1}_{\{I_s=i\}} X_s}{N_{i,t}}$

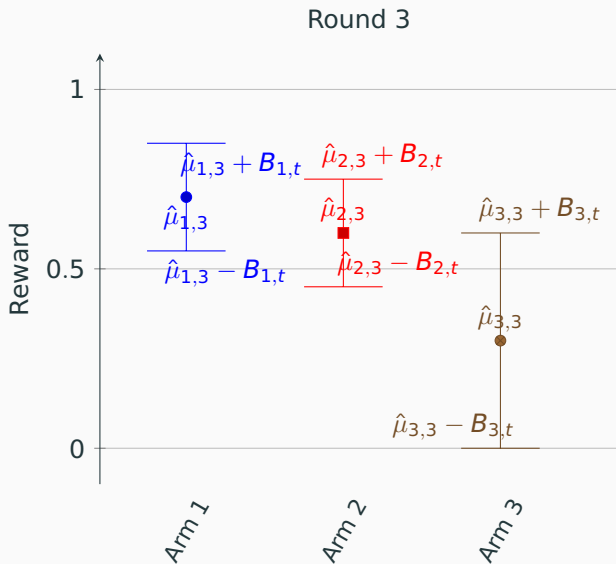
# UCB Illustration



# UCB Illustration



# UCB Illustration



## UCB: Intuition

- Naturally balances exploration and exploitation: an arm has a high UCB if
  - It has a high  $\hat{\mu}_{i,t}$ , or
  - $B_{i,t}$  is large because  $N_{i,t-1}$  is small
- Optimistic because we pretend the rewards are the plausibly best and then do the greedy thing

# UCB: Analysis

- Define  $M_i = \left\lceil \frac{12 \log(n)}{\Delta_i^2} \right\rceil$ , the number of pulls of arm  $i$  such that  $B_{i,t} = \sqrt{\frac{6 \log(t)}{N_{i,t}}} \leq \sqrt{\frac{6 \log(n)}{N_{i,t}}} \leq \frac{\Delta_i}{2}$
- The intuition of the proof is
  1. Since  $\overline{\mathcal{R}}_T = \sum_i \Delta_i \mathbb{E}[N_{i,T}]$ , we bound  $\mathbb{E}[N_{i,T}]$  first.
  2. With high probability, we will never pull arm  $i$  more than  $M_i$  times, so

$$\mathbb{E}[N_{i,T}] = \mathbb{E} \sum_{t=1}^T \mathbb{1}_{\{I_t=i\}} = M_i + \sum_{t=M_i}^T \underbrace{\mathbb{E} \mathbb{1}_{\{I_t=i, N_{i,t} > M_i\}}}_{\text{we will bound this}}$$

3. If  $\{I_t = i, N_{i,t} > M_i\}$  occurs, then the UCB for  $i^*$  or for  $i$  must be wrong (next slide)

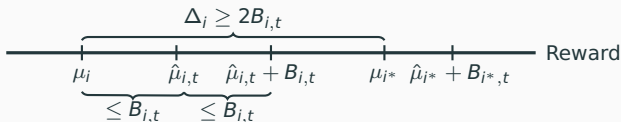
# UCB: Analysis

Claim: if  $\{I_t = i, N_{i,t} > M_i\}$  occurs, then the UCB for  $i^*$  or for  $i$  must be wrong.

Suppose that  $N_{i,t} > M_i$ ,  $\hat{\mu}_{i,t} + B_{i,t} \geq \mu_i$ , and  $\hat{\mu}_{i^*,t} + B_{i^*,t} \geq \mu_{i^*}$ . Then

$$\hat{\mu}_{i^*,t} + B_{i^*,t} \geq \mu_{i^*} = \mu_i + \Delta_i \geq \underbrace{\mu_i + 2B_{i,t}}_{\text{by choice of } B_{i,t}} \geq \hat{\mu}_{i,t} + B_{i,t},$$

so the algorithm will not choose  $I_t = i$ . Hence, one (or both) of upper bounds must be wrong.



Hence, we must bound  $P(\hat{\mu}_{i,t} + B_{i,t} \leq \mu_i)$ .

# UCB: Analysis

Using our sub-Gaussian concentration inequality,

$$\begin{aligned} P(\hat{\mu}_{i,t} + B_{i,t} \leq \mu_i) &\leq P\left(\exists s \leq t : \hat{\mu}_{i,s} - \mu_i \leq \sqrt{\frac{6 \log(t)}{s}}\right) \\ &\leq \sum_{s=1}^t P\left(\hat{\mu}_{i,s} - \mu_i \leq \sqrt{\frac{6 \log(t)}{s}}\right) \\ &\leq \sum_{s=1}^t \exp\left\{-\frac{3 \log(t)}{s}\right\} \leq \sum_{s=1}^t t^{-3} \leq t^{-2}. \end{aligned}$$

The same inequality holds for  $i^*$ , so

$$\overline{\mathcal{R}}_T = \sum_i \Delta_i \mathbb{E}[N_{i,T}] \leq \sum_i \Delta_i \left( \frac{12 \log(n)}{\Delta_i^2} + 2 \sum_{t=M_i+1}^T t^{-2} \right).$$



## UCB: Analysis Step 2

### **Theorem (UCB upper bound [Auer, 2002])**

*The UCB1 algorithm on 1-sub-Gaussian data has*

$$\overline{\mathcal{R}}_T \leq \sum_i \frac{12 \log(n)}{\Delta_j} + o(1).$$

## UCB: Analysis Step 2

### Theorem (UCB upper bound [Auer, 2002])

The UCB1 algorithm on 1-sub-Gaussian data has

$$\overline{\mathcal{R}}_T \leq \sum_i \frac{12 \log(n)}{\Delta_j} + o(1).$$

### Theorem (Lower Bound [Lai and Robbins, 1985])

Suppose we have a parametric family  $P_\theta$ , for some  $\theta_1, \dots, \theta_k$

$$\liminf_{T \rightarrow \infty} \frac{\overline{\mathcal{R}}_T}{\log(T)} \geq \sum_{i \neq i^*} \frac{\Delta_i}{KL(P_{\theta_i}, P_{\theta_{i^*}})} \approx O\left(\sum_{i \neq i^*} \frac{1}{\Delta_i}\right)$$

E.g. if  $P_\theta$  is Bernoulli, then  $\frac{(\theta_i - \theta_{i^*})^2}{\theta_{i^*}(1 - \theta_{i^*})} \geq KL(P_{\theta_i}, P_{\theta_{i^*}}) \geq 2(\theta_i - \theta_{i^*})^2$ .

# Lower Bound Reasoning

- Fix a strategy and consider two problem instances:
  1.  $\nu_1, \nu_2, \dots, \nu_k$ ; with  $P$  as the joint distribution over  $(I_t, r_{i,t})$
  2.  $\nu_1, \nu'_2, \dots, \nu_k$ ; with  $P'$  as the joint distribution over  $(I_t, r_{i,t})$
  3. The optimal arm is different:  $\mu'_2 \geq \mu_1 \geq \mu_2 \geq \mu_3 \geq \dots$
  4. The data from  $P$  and  $P'$  will look very similar
- An algorithm that does well on  $P$  must not pull arm 2 too many times; hence, it will not do well on  $P'$
- “Similar” is quantified by a change-of-measure identity; e.g.  
 $P'(A) = e^{-k\widehat{kl}_{N_2, T}} P(A)$ , where  $\widehat{kl}_t = \sum_{s=1}^t \log \frac{d\nu_2}{d\nu'_2}(X_{2,s})$
- Hence, an algorithm cannot tell if it is  $P$  or  $P'$  and must get high regret under  $P'$ , mistakenly believing it is playing in  $P$

# Algorithm design principle: probability matching

- We put a  $\pi$  over  $\mu_i$  and a likelihood  $\nu_i = P(\cdot|\mu_i)$  over arm  $i$
- We choose  $P(I_t = i) = P(\mu_i = \mu_{i^*} | \text{history})$
- We usually pick conjugate models (e.g.  $\mu_i \sim N(0, 1)$ ,  $X_t \sim N(\mu_i, 1)$ )

## Algorithm: Thompson Sampling

Given: game length  $T$ , prior  $\pi(\mu)$ , likelihoods  $p(\cdot|\mu)$

Initialize posteriors  $p_{i,0}(\mu) = \pi(\mu)$

For  $t = 1, 2, \dots, T$ :

- Draw  $\theta_{i,t} \sim p_{i,t-1}$  for all  $i$
- Choose  $I_t = \arg \max_i \theta_{i,t}$
- Receive and observe  $X_t \sim \nu_{I_t}$
- Update the posterior  $p_{I_t,t}(\mu) = p(X_t|\mu)p_{I_t,t-1}(\mu)$

- We put a  $\pi$  over  $\mu_i$  and a likelihood  $\nu_i = P(\cdot|\mu_i)$  over arm  $i$
- We choose  $P(I_t = i) = P(\mu_i = \mu_{i^*} | \text{history})$
- We usually pick conjugate models (e.g.  $\mu_i \sim N(0, 1)$ ,  $X_t \sim N(\mu_i, 1)$ )

# Thompson Sampling: Overview

- *Not* Bayesian; uses Bayesian techniques, but the guarantees are frequentist
- Arms with small  $N_{i,t}$  implies a wide posterior, hence a good probability of being selected
- Generally performs empirically better than UCB
- Arms with small  $N_{i,t}$  implies a wide posterior, hence a good probability of being selected
- Analysis is difficult

# Thompson Sampling: Upper Bound

## Theorem (Agrawal and Goyal [2013])

For binary rewards, Gamma-Beta Thompson sampling has

$$\mathbb{E}[R_T] \leq (1 + \epsilon) \sum_{i \neq i^*} \Delta_i \frac{\log(T)}{kl(\mu_i, \mu_{i^*})} + O\left(\frac{N}{\epsilon^2}\right).$$

- The proof is much more technical than UCB's
- We cannot rely on the upper bounds being correct w.h.p.
- For some to-be-tuned  $\mu_j \leq x_j \leq y_j \leq \mu_{j^*}$ , we have

$$\begin{aligned} \mathbb{E}[N_{i,T}] &\leq \sum_{t=1}^T P(I_t = i) \\ &\leq \sum_{t=1}^T P(I_t = i, \hat{\mu}_{i,t-1} \leq x_i, \theta_{i,t} \geq y_i) \quad \left(O\left(\frac{\log(T)}{kl(x_j, y_j)}\right)\right) \\ &\quad + \sum_{t=1}^T P(I_t = i, \hat{\mu}_{i,t-1} \leq x_i, \theta_{i,t} \leq y_i) \quad (\text{the tricky case}) \\ &\quad + \sum_{t=1}^T P(I_t = i, \hat{\mu}_{i,t-1} \geq x_i) \quad (\text{Small by concentration}) \end{aligned}$$

# Thompson Sampling: Proof Outline

- The tricky case is  $\sum_{t=1}^T P(I_t = i, \hat{\mu}_{i,t-1} \leq x_i, \theta_{i,t} \leq y_i)$
- This happens when we have enough samples of  $i$  but not many of  $i^*$
- A key lemma argues that, on  $\hat{\mu}_{i,t-1} \leq x_i, \theta_{i,t} \leq y_i$ , the probability of picking  $i$  is a constant less than of picking  $i^*$ :

$$\begin{aligned} & \sum_{t=1}^T P(I_t = i, \hat{\mu}_{i,t-1} \leq x_i, \theta_{i,t} \leq y_i) \\ & \leq \sum_{t=1}^T \underbrace{\frac{P(\theta_{i^*,t} \leq y_i)}{P(\theta_{i^*,t} > y_j)}}_{\text{exponentially small}} P(I_t = i^*, \hat{\mu}_{i,t-1} \leq x_i, \theta_{i,t} \leq y_i) = O(1) \end{aligned}$$

- Hence, we will quickly get enough samples of  $i^*$

# Best of Both Worlds

- The stochastic and adversarial algorithms are quite different
- A natural question: is there an algorithm that
  - gets  $\mathcal{R}_T = O(\sqrt{TK})$  regret for adversarial
  - gets  $\mathcal{R}_t = O(\sum_i \log(T)/\Delta_i)$  regret for stochastic
  - without knowing the setting?
- Bubeck and Slivkins [2012] propose an algorithm that assumes stochastic but falls back to UCB once adversarial data is detected
- Zimmert and Seldin [2019] show that (for pseudo-regret), it is possible
  - Their algorithm: online mirror descent with  $\frac{1}{2}$ -Tsallis entropy
  - $\Psi(w) = -\sum_i 4(\sqrt{w_i} - \frac{1}{2}w_i)$



# Pure Exploration

---

## A new problem

- What if we only wanted to identify the best arm  $i^*$  without caring about loss along the way?
- Intuitively, we would explore more; we are happy to accrue less reward if we get more useful samples.
- More similar to hypothesis testing; useful for selecting treatments
- Known as “Best Arm Identification” or “Pure Exploration”

## Two settings

### Protocol: Best-arm identification protocol

Given: number of arms  $K$

For  $t = 1, 2, \dots$ ,

- The learner picks arm  $I_t \in \{1, \dots, K\}$
- The learner observes  $X_t \sim \nu_{I_t}$
- The learner decides whether to stop

The learner returns arm  $A$

Two settings:

	<i>fixed-confidence</i>	<i>fixed-budget</i>
Input	$\delta > 0, \epsilon > 0$	$T$
Goal	arm $A$ is $(\epsilon, \delta)$ -PAC	maximize $P(A = i^*)$
Stopping	once learner is confident	after $T$ rounds

Arm  $A$  is  $(\epsilon, \delta)$ -PAC if  $P(\mu_A \geq \mu_{j^*} - \epsilon) \geq 1 - \delta$ .

- Standard stochastic bandit algorithms under explore (they fail to meet lower bounds on this problem)
- Many can be adapted
  - LUCB [?]
  - Top-Two Thompson Sampling [?]
- Instead, we will describe a new algorithm design principle

# Algorithm Design Principle: Action Elimination

## Algorithm: Successive Elimination

Given: confidence  $\delta > 0$

Initialize plausibly-best set  $S = \{1, \dots, K\}$

For  $t = 1, 2, \dots$ :

- Pull all arms in  $S$  and update  $\hat{\mu}_{i,t}$
- Calculate  $B_t = \sqrt{t^{-1} \log(nt^2/\delta)}$
- Remove  $i$  from  $S$  if  $\underbrace{\max_{j \in S} \hat{\mu}_{j,t} - B_t}_{\text{Lowest } \mu_i^* \text{ could be}} \geq \underbrace{\hat{\mu}_{i,t} + B_t}_{\text{highest } \mu_i \text{ could be}}$
- If  $|S| = 1$ , stop and return  $A = S$ .

- $S$  is a list of plausibly-best arms
- Each epoch, all arms that cannot be the best (if the bounds hold) are removed

# Successive Elimination Analysis

- Define the “bad event”  $\mathcal{E} = \bigcup_{i,t} \{|\hat{\mu}_{i,t} - \mu_i| \leq B_t(\delta)\}$ : we have

$$\begin{aligned} P(\mathcal{E}) &\leq \sum_{i,t} P\left(|\hat{\mu}_{i,t} - \mu_i| \leq \sqrt{t^{-1} \log(Kt^2/\delta)}\right) \leq \sum_{i,t} 2e^{-\frac{1}{2} \log\left(\frac{Kt^2}{\delta}\right)} \\ &\leq \sum_{i,t} \frac{2e^{-2\delta}}{Kt^2} = \frac{2e^{-2\pi^2}}{6} \delta \leq \delta \end{aligned}$$

- (Correctness) If  $\mathcal{E}$  does not happen,
  - $|\hat{\mu}_{i^*} - \mu_{i^*}| \leq B_t$  and  $|\mu_j - \hat{\mu}_j| \leq B_t$  for all  $j$ . Thus,  
 $\hat{\mu}_j - \hat{\mu}_{i^*} \leq (\mu_{i^*} - \hat{\mu}_{i^*}) + (\mu_{i^*} - \mu_j) + (-\mu_j - \hat{\mu}_{i^*}) \leq 2B_t$
  - $i$  is removed if  $\max_{j \in S} \hat{\mu}_{j,t} - \hat{\mu}_{i,t} \geq 2B_t \Rightarrow i^*$  is never removed
  - $\lim_{t \rightarrow \infty} B_t(\delta) \rightarrow 0$ : every arm will eventually be removed
  - Hence, Successive Elimination is  $(0, \delta)$ -PAC
- (Sample Complexity): arm  $i$  will be eliminated once  $\Delta_i \leq 2B_t$ 
  - We can verify that  $N_i = O(\Delta_i^{-2} \log(K/\delta\Delta_i))$  is sufficient
  - Total sample complexity of  $\sum_i \Delta_i^{-2} \log(K/\delta\Delta_i)$
- Can convert to a  $(\epsilon, \delta)$ -PAC algorithm by stopping early

# Linear Stochastic Bandits

---

## Bonus: Linear Contextual Bandits

### Protocol: Contextual Linear Bandit Protocol

Given: game length  $T$ , number of arms  $K$

For  $t = 1, 2, \dots, T$ ,

- The learner sees one context per arm  $c_{1,t}, \dots, c_{K,t}$
- The learner picks action  $I_t \in \{1, \dots, K\}$
- The learner observes and receives reward  $X_t = \langle c_{I_t,t}, \theta^* \rangle + \xi_t$

Regret is defined w.r.t. an agent that knows the true  $\theta$ :

$$\mathcal{R}_T = \sum_{t=1}^T \max_i x_{i,t}^\top \theta^* - \sum_{t=1}^T X_t$$



# Algorithm Design Principle: Optimism

## Algorithm: OFUL [Abbasi-Yadkori et al., 2011]

Initialize  $\hat{\theta}_0 = 0$ ,  $B_0 = \mathbb{R}^d$

For  $t = 1, 2, \dots, T$ :

- Receive contexts  $c_{1,t}, \dots, c_{K,t}$
  - Choose  $(I_t, \tilde{\theta}_t) = \arg \max_{j \in \{1, \dots, K\}, \theta \in B_{t-1}} \theta^\top c_{j,t}$  (optimism)
  - Observe  $X_t = c_{I_t,t}^\top \theta^* + \xi_t$
  - Calculate  $V_t = \sum_{s=1}^t c_s c_s^\top + \lambda I$  and  $r_t = \sqrt{\log \frac{\det(V_t)}{\delta^2 \lambda^d}} + \sqrt{\lambda} \|\theta^*\|$
  - Calculate  $\hat{\theta}_t = V_t^{-1} \left( \sum_{s=1}^t c_s X_s \right)$  (ridge)
  - Update  $B_t = \{ \theta : (\theta - \hat{\theta}_t)^\top V_t (\theta - \hat{\theta}_t) \leq r_t \}$
- 
- If  $\xi_t$  is 1-sub-Gaussian,  $B_t$  is a confidence sequence with  $P(\forall t > 0 : \theta^* \in B_t) \geq 1 - \delta$  (more examples in [de la Peña et al., 2009, Howard et al., 2020])

# Analysis

- Regret decomposes over rounds:
- Recall that  $(I_t, \tilde{\theta}_t) = \arg \max_{i \in \{1, \dots, K\}, \theta \in B_{t-1}} \theta^\top C_{i,t}$

$$\begin{aligned} \mathcal{R}_t - \mathcal{R}_{t-1} &= c_{i_t^*}^\top \theta^* - c_{I_t}^\top \theta^* \\ &\leq c_{I_t}^\top \tilde{\theta}_t - c_{I_t}^\top \theta^* && \text{(by optimism)} \\ &\leq c_{I_t}^\top (\tilde{\theta}_t - \hat{\theta}_{t-1}) + c_{I_t}^\top (\hat{\theta}_{t-1} - \theta^*) \\ &\leq \|c_{I_t}\|_{V_t} \underbrace{\|\tilde{\theta}_t - \hat{\theta}_{t-1}\|_{V_t}}_{\leq r_t} + \|c_{I_t}\|_{V_t} \underbrace{\|\hat{\theta}_{t-1} - \theta^*\|_{V_t}}_{\leq r_t} \end{aligned}$$

- After some algebra, we can show, with probability  $\geq 1 - \delta$ , that

$$\mathcal{R}_T = O\left(\frac{d \log(1/\delta)}{\Delta}\right)$$

- The shared structure lets us learn a lot!

# Review

- Setting: adversarial bandits
  - Exp3 (exponential weights)
- Setting: stochastic bandits
  - UCB (optimism)
  - Thompson Sampling (probability matching)
- Setting: pure exploration
  - Successive Elimination (action-elimination)
- Setting: linear contextual bandits
  - OFUL (optimism)

Thanks!

# References

---

- Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- Shipra Agrawal and Navin Goyal. Further optimal regret bounds for thompson sampling. In *Artificial intelligence and statistics*, pages 99–107, 2013.
- Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits. 2009.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov): 397–422, 2002.
- Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002a.

Peter Auer, Nicolo Cesa-Bianchi, Yoav Freund, and Robert E Schapire. The nonstochastic multiarmed bandit problem. *SIAM journal on computing*, (1):48–77, 2002b.

Gábor Bartók, Dean P Foster, Dávid Pál, Alexander Rakhlin, and Csaba Szepesvári. Partial monitoring – classification, regret bounds, and algorithms. *Mathematics of Operations Research*, (4):967–997, 2014.

Sébastien Bubeck and Aleksandrs Slivkins. The best of both worlds: Stochastic and adversarial bandits. In *Conference on Learning Theory*, pages 42–1, 2012.

Victor H. de la Peña, Michael J. Klass, and Tze Leung Lai. Theory and applications of multivariate self-normalized processes. *Stochastic Processes and their Applications*, 119(12): 4210–4227, December 2009. ISSN 0304-4149.

Steven R. Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Time-uniform chernoff bounds via nonnegative supermartingales. *Probab. Surveys*, 17:257–317, 2020. doi: 10.1214/18-PS321. URL <https://doi.org/10.1214/18-PS321>.

Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in Applied Mathematics*, (1):4–22, 1985.

Shie Mannor and Ohad Shamir. From bandits to experts: On the value of side-observations. In *Advances in Neural Information Processing Systems*, pages 684–692, 2011.

Gergely Neu. Explore no more: Improved high-probability regret bounds for non-stochastic bandits. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3168–3176, 2015.

William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.

Julian Zimmert and Yevgeny Seldin. An optimal algorithm for stochastic and adversarial bandits. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 467–475, 2019.