

Online Convex Optimization

Tareq Si Salem

Postdoctoral Research Associate at Northeastern University

October, 2023

Outline I

1. Motivating Problem: Prediction from Expert Advice

2. Online Convex Optimization

3. Some Applications

- 3.1 Online Spam Filtering
- 3.2 Online Shortest Paths/Routing
- 3.3 Matrix Completion and Recommendation Systems

4. Online Gradient Descent

- 4.1 Online Gradient Descent
- 4.2 Regret of Online Gradient Descent
- 4.3 Regret of Online Gradient Descent (Proof)

5. Expert Problem (Revisited)

- 5.1 Expert Problem – Applying OGD
- 5.2 Expert Problem – Adapting To Geometry via Online Mirror Descent
- 5.3 Regret of Online Mirror Descent
- 5.4 Simplex/Entropic Setup of OMD

6. Going Beyond Full-Information

- 6.1 Going Beyond Full-Information – Bandit Convex Optimization
- 6.2 Multi-Armed Bandits (MAB)
- 6.3 A Reduction from Limited Information to Full Information
- 6.4 MAB: EXP3 Simultaneous Exploration and Exploitation
- 6.5 Flaxman/Kalai/McMahan (FKM) Algorithm

7. Reduction: Learning \rightarrow OCO

Outline II

7.1 PAC (Probably Approximately Correct) learning and OCO

7.2 Reduction: Learning \rightarrow OCO

8. k -Experts Problem and Switching Costs

8.1 k -Experts Problem

8.2 k -Experts: Casting as a Caching Problem and Dimensionality Reduction

8.3 k -Experts: Recovering the Original Setting

8.4 k -Experts: Switching Costs

8.5 k -Experts: Independent Sampling

8.6 k -Experts: Optimal Transport

8.7 k -Experts: Simpler Approach

8.8 k -Experts: Coupling Schemes a Qualitative Evaluation

Motivating Problem: Prediction from Expert Advice



Motivating Problem: Prediction from Expert Advice

The player has to choose among the advice of n given experts. After making their choice, a loss between in $[0, 1]$ is incurred. A player constructs a “belief” (distribution) over the set of experts

$$\Delta_n \triangleq \{ \mathbf{x} \in [0, 1]^n : \|\mathbf{x}\|_1 = 1 \}. \quad (1)$$

At each time step $t = 1$ to T ,

- Player decides $\mathbf{x}_t \in \Delta_n$ and samples expert $i_t \sim \mathbf{x}_t$
- Adversary picks losses of each expert $\mathbf{l}_t = (l_{t,i})_{i \in [n]}$ and reveals them to the player.
- Player suffers an expected loss $f_t(\mathbf{x}_t) = \mathbb{E}_{i \sim \mathbf{x}_t} [l_{t,i}] = \mathbf{l}_t \cdot \mathbf{x}_t$.

In this setting, the optimization error is ill-defined. Instead, the goal of the player is to minimize *regret*:

$$\text{regret}_T = \sum_{t=1}^T \mathbf{l}_t \cdot \mathbf{x}_t - \min_{i^* \in [n]} \sum_{t=1}^T l_{t,i^*}. \quad (2)$$

The goal of the player is to have sublinear regret ($\text{regret}_T = o(T)$), i.e., on average the losses experienced by the player is **as good as** the best expert in hindsight.

Motivating Problem: Prediction from Expert Advice

A natural strategy. Consider that at timeslot t , the player greedily selects the best expert seen so far, also known as the follow the leader strategy (or fictitious play in game theory).

Theorem 1

FTL strategy incurs $\text{regret}_T = \Omega(T)$.

Proof.

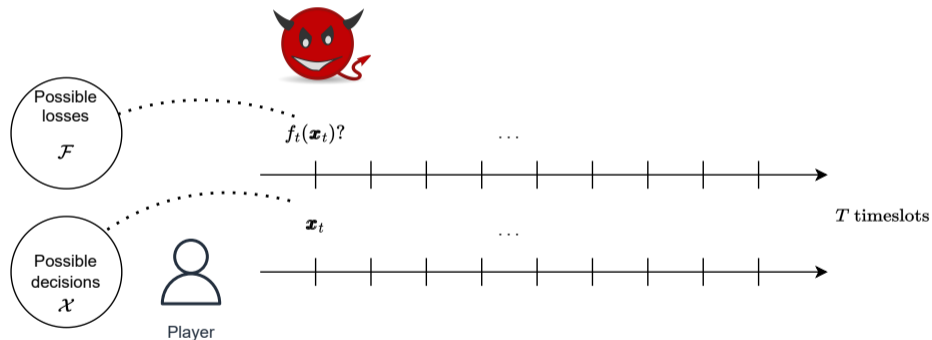
Consider two experts $n = 2$ and a sequence of losses $\mathbf{l}_1 = (1, 0), \mathbf{l}_2 = (0, 1), \dots$. The best expert oscillates between 1 and 2. The player will incur a total loss of T (loss of 1 at every timeslot), whereas selecting a fixed expert incurs a total loss of $T/2$, i.e.,

$$\text{regret}_T \geq T - T/2 = \Omega(T). \quad (3)$$

□

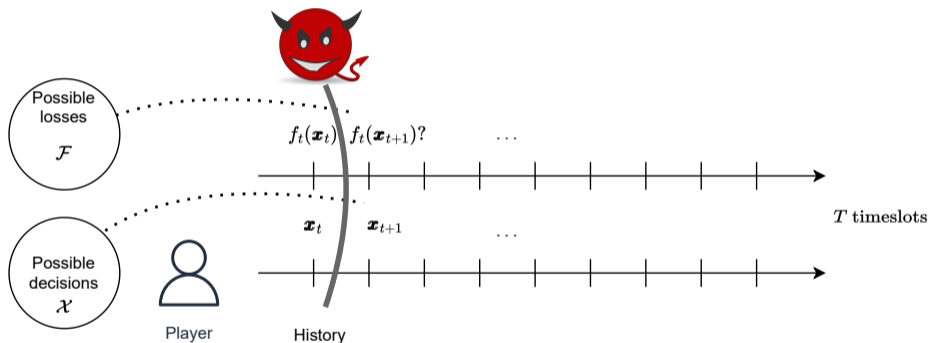
The Online Convex Optimization Setting

A player makes decisions iteratively. At the time of making the decision the outcome or outcomes associated with it is unknown to the player, and can even depend on the action taken by the decision maker.



The Online Convex Optimization Setting

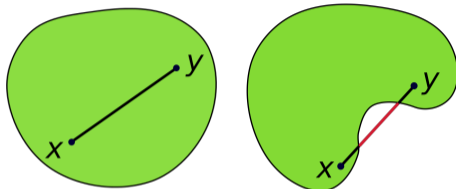
A player makes decisions iteratively. At the time of making the decision the outcome or outcomes associated with it is unknown to the player, and can even depend on the action taken by the decision maker.



- The losses determined by an adversary should not be allowed to be unbounded.
- The decision set must be bounded and/or structured, and possibly infinite.

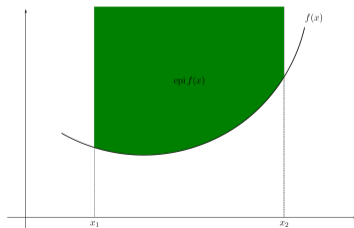
The Online Convex Optimization Setting

OCO models the decision set as a convex set in the Euclidean space $\mathcal{X} \subseteq \mathbb{R}^d$, and the costs as bounded convex functions.



(a) Convex set

(b) Non-convex set



(c) Convex functions

A More Formal Description of OCO

At iteration t , the player picks a new decision \mathbf{x}_{t+1} after incurring the cost f_t , according to a mapping (algorithm) $\mathcal{A}: (\mathcal{X} \times \mathcal{F})^t \rightarrow \mathcal{X}$ given as

$$\mathbf{x}_{t+1} = \mathcal{A}((\mathbf{x}_1, f_1), (\mathbf{x}_2, f_2), \dots, (\mathbf{x}_t, f_t)) \in \mathcal{X}. \quad (4)$$

The regret of the player under policy \mathcal{A} is then

$$\text{regret}_T(\mathcal{A}) \triangleq \sup_{(f_t)_{t \in [T]} \in \mathcal{F}^T} \left\{ \sum_{t=1}^T f_t(\mathbf{x}_t) - \min_{\mathbf{x}^* \in \mathcal{X}} \sum_{t=1}^T f_t(\mathbf{x}^*) \right\}. \quad (5)$$

When the regret is sublinear ($\text{regret}_T = o(T)$), i.e., on average the losses experienced by the player are **as good as** the best decision in hindsight.



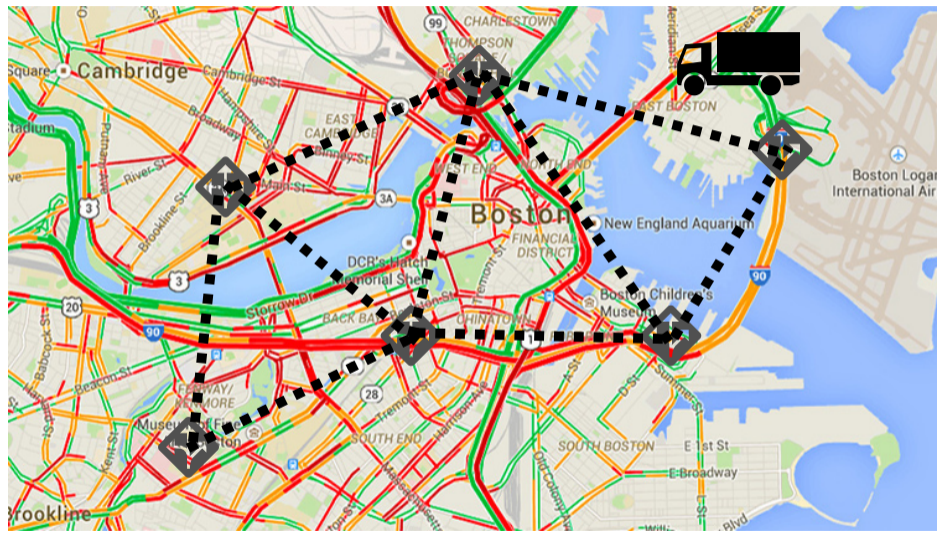
Online Spam Filtering

Consider an online spam-filtering system. Repeatedly, emails arrive in the system and are classified as spam or valid. For simplicity, consider a linear variant of the model:

- Each email is represented as a vector $\mathbf{a} \in \mathbb{R}^d$, where d is the number of words in the dictionary. The email is the aggregation of one-hot encoding of words in the dictionary (“bag-of-words” representation).
- To predict whether an email is spam, we learn a filter, a vector $\mathbf{x} \in \mathbb{R}^d$.
- Given an email \mathbf{a}_t and a label b_t , a loss $l_t(\mathbf{x}) = l_{\mathbf{a}_t, b_t}(\mathbf{x})$ is revealed to the player (e.g., the hinge loss $\text{hinge}(b_t \cdot \mathbf{x}^\top \mathbf{a}_t) = \max\{0, 1 - b_t \cdot \mathbf{x}^\top \mathbf{a}_t\}$).

Such a system has to cope with adversarially generated data and dynamically change with the varying input.

Online Shortest Paths/Routing

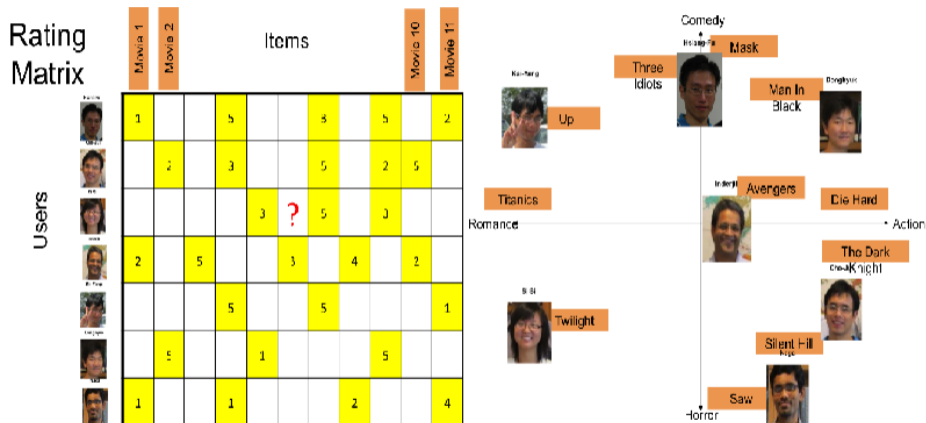


In the online shortest path problem, the decision maker is given a directed graph $G = (V, E)$ and a source-sink pair $u, v \in V$. The set \mathcal{X} of all all distributions over paths (flows) in a graph is a convex set in \mathbb{R}^E , with $\mathcal{O}(m + |V|)$ constraints.

At each time step $t = 1$ to T

- The player chooses a flow $\mathbf{x}_t \in \mathcal{X}$ and samples a path $\mathbf{p}_t \sim \mathbf{x}_t$
- The adversary chooses weights on the edges of the graph $\mathbf{w}_t : E \rightarrow \mathbb{R}$. The player incurs the expected loss $f_t(\mathbf{x}_t) = \mathbf{x}_t \cdot \mathbf{w}_t$.

Matrix Completion and Recommendation Systems



Matrix Completion and Recommendation Systems

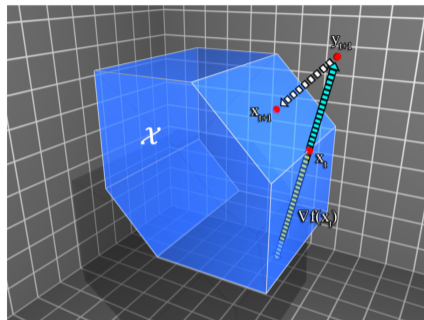
For example, for the case of binary recommendations for music, we have a matrix $\mathbf{X} \in \{0, 1\}^{n \times m}$ where n is the number of persons considered, m is the number of songs in our library, and 0/1 signifies dislike/like respectively.

At each time step $t = 1$ to T

- The player selects a preference matrix $\mathbf{X}_t \in \mathcal{X}$, where $\mathcal{X} \subseteq [0, 1]^{n \times m}$ (with a rank constraint).
- An adversary then chooses a user/song pair (i_t, j_t) along with a “real” preference for this pair $y_t \in \{0, 1\}$. The loss experienced by the decision maker can be described by some convex loss function, $f_t(\mathbf{X}) = (X_{i_t, j_t} - y_t)^2$.

The natural comparator in this scenario is a low-rank matrix, which corresponds to the intuitive assumption that preference is determined by few unknown factors

Online Gradient Descent



At time t , Online Gradient Descent (OGD)'s update rule is

(gradient update)

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \eta_t \nabla f_t(\mathbf{x}_t)$$

(projection step)

$$\mathbf{x}_{t+1} = \Pi_{\mathcal{X}}(\mathbf{y}_{t+1})$$

The operator $\Pi_{\mathcal{X}} : \mathbb{R}^d \rightarrow \mathcal{X}$ is the Euclidean projection $\Pi_{\mathcal{X}}(\mathbf{y}) \triangleq \arg \min_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x} - \mathbf{y}\|_2$.

Theorem 2

Online gradient descent with fixed step sizes η guarantees the following for all $T \geq 1$:

$$\mathbf{regret}_T \leq \frac{D^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t)\|_2^2. \quad (6)$$

For a fixed learning rate $\eta = \frac{D}{G\sqrt{T}}$, the regret is upper bounded by $\mathbf{regret}_T \leq DG\sqrt{T}$.

Regret of Online Gradient Descent (Proof)

A fundamental property of projections into convex bodies is that for an arbitrary $\mathbf{x}' \in \mathbb{R}^d$, we have for all $\mathbf{x} \in \mathcal{X}$:

$$\|\Pi_{\mathcal{X}}(\mathbf{x}') - \mathbf{x}\|_2^2 \leq \|\mathbf{x}' - \mathbf{x}\|_2^2$$

Applying the above,

$$\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2 = \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\Pi_{\mathcal{X}}(\mathbf{x}_t - \eta \nabla f_t(\mathbf{x}_t)) - \mathbf{x}^*\|^2 \quad (7)$$

$$\geq \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_t - \eta \nabla f_t(\mathbf{x}_t) - \mathbf{x}^*\|^2 \quad (8)$$

$$= 2\eta \nabla f_t(\mathbf{x}_t) \cdot (\mathbf{x}_t - \mathbf{x}^*) - \eta^2 \|\nabla f_t(\mathbf{x}_t)\|_2^2 \quad (9)$$

And so,

$$\nabla f_t(\mathbf{x}_t) \cdot (\mathbf{x}_t - \mathbf{x}^*) \leq \frac{1}{2\eta} \left(\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2^2 \right) + \frac{\eta}{2} \|\nabla f_t(\mathbf{x}_t)\|_2^2 \quad (10)$$

Summing over t ,

$$\sum_{t=1}^T \nabla f_t(\mathbf{x}_t) \cdot (\mathbf{x}_t - \mathbf{x}^*) \leq \frac{1}{2\eta} \left(\|\mathbf{x}_1 - \mathbf{x}^*\|_2^2 - \|\mathbf{x}_{T+1} - \mathbf{x}^*\|_2^2 \right) + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t)\|_2^2 \quad (11)$$

$$\leq \frac{D^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t)\|_2^2 \quad (12)$$

Expert Problem – Applying OGD

We can apply OGD to the expert problem. Recall the the decision set \mathcal{X} is the probability simplex Δ_n and the loss of picking an expert is in $[0, 1]$. The diameter of the set is $\|\mathbf{x} - \mathbf{x}'\|_2 \leq D = 1$ for $\mathbf{x}, \mathbf{x}' \in \Delta_n$. The gradients are bounded by $\|\nabla f_t(\mathbf{x}_t)\|_2^2 \leq G = \sqrt{n}$. The regret is sublinear with the rate \sqrt{nT} . **The dependency on the number of experts can be improved further!**

We will describe a gradient-based scheme that can be tailored to the geometry of the problem.

Expert Problem – Adapting To Geometry via Online Mirror Descent

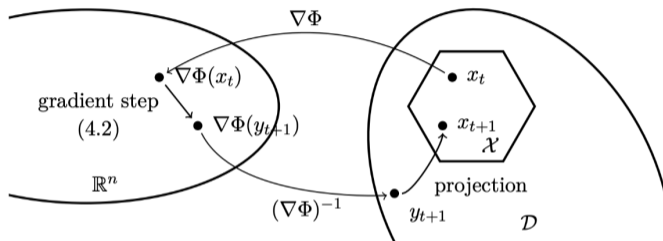


Figure: Illustration of online mirror descent

The policy is parameterized by (1) a learning rate η , and (2) a mirror map $\Phi : \mathcal{D} \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$. The projection step generalizes Euclidean projection and it is given by

$$\Pi_{\mathcal{X}}^{\Phi}(\mathbf{y}) \triangleq \arg \min_{\mathbf{x} \in \mathcal{D} \cap \mathcal{X}} D_{\Phi}(\mathbf{x}, \mathbf{y}) = \arg \min_{\mathbf{x} \in \mathcal{D} \cap \mathcal{X}} \Phi(\mathbf{x}) - \Phi(\mathbf{y}) - \nabla\Phi(\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y}) \quad (\text{Bregman projection}). \quad (13)$$

Theorem 3

Let Φ be a mirror map ρ -strongly convex on $\mathcal{X} \cap \mathcal{D}$ w.r.t. $\|\cdot\|$. Let $R^2 = \sup_{\mathbf{x} \in \mathcal{X} \cap \mathcal{D}} \Phi(\mathbf{x}) - \Phi(\mathbf{x}_1)$, and f_t be convex and G -Lipschitz w.r.t. $\|\cdot\|$. Then online mirror descent with $\eta = \frac{R}{G} \sqrt{\frac{2\rho}{T}}$ satisfies

$$\sum_{t=1}^T f_t(\mathbf{x}_t) - f_t(\mathbf{x}) \leq \frac{D_{\Phi}(\mathbf{x}, \mathbf{x}_1)}{\eta} + \eta \frac{\sum_{t=1}^T \|\nabla f_t(\mathbf{x}_t)\|_*^2}{2\rho} = RL \sqrt{\frac{2}{\rho} T}. \quad (14)$$

- A differentiable function f is α -strongly-convex with respect to a norm $\|\cdot\|$ if $\forall \mathbf{x} \in \mathcal{X}$

$$f(\mathbf{x}) - f(\mathbf{y}) \leq \nabla f(\mathbf{x}) \cdot (\mathbf{x} - \mathbf{y}) - \frac{\alpha}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{y} \in \mathcal{X}. \quad (15)$$

- Let $\|\cdot\|$ be a norm on \mathbb{R}^d . The dual norm $\|\cdot\|_*$ is defined as $\|\mathbf{u}\|_* = \sup_{\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\| \leq 1} \mathbf{x} \cdot \mathbf{u}$.

Ball setup. When $\Phi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2$, the dual and primal spaces are identical ($\nabla \Phi(\mathbf{x}) = \mathbf{x}$ is the identity mapping) and the Bregman projection is simply the Euclidean projection $\Pi_{\mathcal{X}}$. OGD is an instance of OMD.

Simplex/Entropic Setup of OMD

Simplex/Entropic setup. A more interesting choice of a mirror map is given by the negative entropy

$$\Phi(\mathbf{x}) = \sum_{i=1}^d x_i \log(x_i). \quad (16)$$

The Bregman divergence of this mirror map is given by $D_{\Phi}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d x_i \log(\frac{x_i}{y_i})$. The policy configured for simplex decision sets Δ_n amounts to

$$\text{(gradient update)} \quad \mathbf{x}_{t+1,i} = x_{t,i} e^{-\eta l_{t,i}}, i \in [d] \quad (17)$$

$$\text{(projection step)} \quad \mathbf{x}_{t+1} = \frac{\mathbf{y}_{t+1}}{\|\mathbf{y}_{t+1}\|_1}. \quad (18)$$

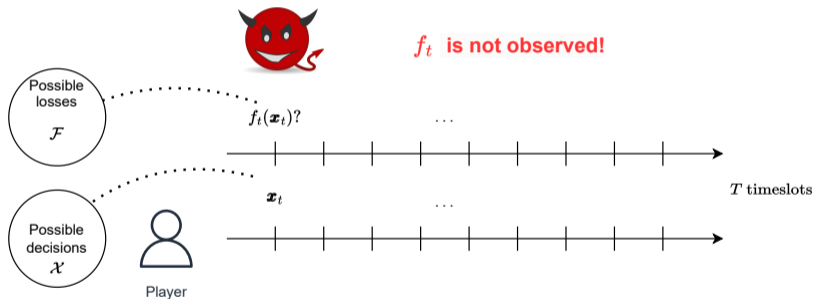
For $\mathbf{x}_1 = (1/n, 1/n, \dots, 1/n)$, one has $R^2 = \log(n)$ and $\|\nabla f_t(\mathbf{y})\|_{\infty} \leq G = 1$. (The norm $\|\cdot\|_{\infty}$ is the dual norm of $\|\cdot\|_1$.) The map Φ is 1-strongly convex w.r.t. $\|\cdot\|_1$ over the simplex Δ_n (Pinsker's inequality). The regret under this configuration has rate of $\sqrt{\log(n)T}$ instead of \sqrt{nT} for OGD. This setup corresponds to the well known **Hedge/Multiplicative Weights Update** algorithm.

Going Beyond Full-Information



Going Beyond Full-Information – Bandit Convex Optimization

As opposed to the OCO model, in which the decision maker has access to a gradient oracle for f_t over \mathcal{X} , in BCO the loss $f_t(\mathbf{x}_t)$ **is the only feedback available** to the online player at iteration t .



Multi-Armed Bandits (MAB)

A classical model for decision making under uncertainty is the multi-armed bandit (MAB) model. This setting is identical to the setting of prediction from expert advice, the only difference being the feedback available to the decision maker. **Only the loss of the selected expert is revealed!**

The MAB problem exhibits an exploration-exploitation tradeoff:

- An efficient (low regret) algorithm has to explore the value of the different actions in order to make the best decision.
- On the other hand, having gained sufficient information about the environment, a reasonable algorithm needs to exploit this action by picking the best action.

A Reduction from Limited Information to Full Information

- 1: Input: convex set $\mathcal{K} \subset \mathbb{R}^n$, first order online algorithm \mathcal{A} .
- 2: Let $\mathbf{x}_1 = \mathcal{A}(\emptyset)$.
- 3: **for** $t = 1$ to T **do**
- 4: Generate distribution \mathcal{D}_t , sample $\mathbf{y}_t \sim \mathcal{D}_t$ with $\mathbf{E}[\mathbf{y}_t] = \mathbf{x}_t$.
- 5: Play \mathbf{y}_t .
- 6: Observe $f_t(\mathbf{y}_t)$, generate \mathbf{g}_t with $\mathbf{E}[\mathbf{g}_t] = \nabla f_t(\mathbf{x}_t)$.
- 7: Let $\mathbf{x}_{t+1} = \mathcal{A}(\mathbf{g}_1, \dots, \mathbf{g}_t)$.
- 8: **end for**

Lemma 4

Let $(f'_t(\mathbf{x}) = \mathbf{g}_t \cdot \mathbf{x})_{t \in [T]}$ be a sequence of linear cost functions, and $(f_t)_{t \in [T]}$ be a sequence of differentiable function, and \mathcal{X} be the convex decision set. Let \mathcal{A} be an OCO policy with regret bound $\text{regret}_T(\mathcal{A}) \leq B_{\mathcal{A}}(\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_T)$ in the full-information setting against losses f'_t . Define the points (\mathbf{x}_t) as: $\mathbf{x}_1 \in \mathcal{X}$, $\mathbf{x}_t \leftarrow \mathcal{A}_t((\mathbf{x}_s, f'_s)_{s \in [t-1]})$ where each \mathbf{g}_t is a vector valued random variable such that:

$$\mathbb{E}[\mathbf{g}_t | \mathbf{x}_1, f_1, \dots, \mathbf{x}_t, f_t] = \nabla f_t(\mathbf{x}_t), \quad (19)$$

Then the following holds for all $\mathbf{x}^* \in \mathcal{X}$:

$$\mathbb{E} \left[\sum_{t=1}^T f_t(\mathbf{x}_t) \right] - \sum_{t=1}^T f_t(\mathbf{x}^*) \leq \mathbb{E} [B_{\mathcal{A}}(\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_T)]. \quad (20)$$

MAB: EXP3 Simultaneous Exploration and Exploitation

- 1: Input: parameter $\varepsilon > 0$. Set $\mathbf{x}_1 = (1/n)\mathbf{1}$.
- 2: **for** $t \in \{1, 2, \dots, T\}$ **do**
- 3: Choose $i_t \sim \mathbf{x}_t$ and play i_t .
- 4: Let

$$\hat{\ell}_t(i) = \begin{cases} \frac{1}{\mathbf{x}_t(i_t)} \cdot \ell_t(i_t), & i = i_t \\ 0, & \text{otherwise} \end{cases}$$

- 5: Update $\mathbf{y}_{t+1}(i) = \mathbf{x}_t(i)e^{-\varepsilon\hat{\ell}_t(i)}$, $\mathbf{x}_{t+1} = \frac{\mathbf{y}_{t+1}}{\|\mathbf{y}_{t+1}\|_1}$
- 6: **end for**

Theorem 5

EXP3 with non-negative losses and $\varepsilon = \sqrt{\frac{\log(n)}{Tn}}$ guarantees the following regret bound:

$$\mathbb{E} \left[\sum_{t=1}^T l_{t,i_t} \right] - \min_{i^* \in [n]} \sum_{t=1}^T l_{t,i^*} \leq 2\sqrt{Tn \log(n)}. \quad (21)$$

Flaxman/Kalai/McMahan (FKM) Algorithm

- 1: Input: decision set \mathcal{K} containing $\mathbf{0}$, set $\mathbf{x}_1 = \mathbf{0}$, parameters δ, η .
- 2: **for** $t = 1$ to T **do**
- 3: Draw $\mathbf{u}_t \in \mathbb{S}_1$ uniformly at random, set $\mathbf{y}_t = \mathbf{x}_t + \delta \mathbf{u}_t$.
- 4: Play \mathbf{y}_t , observe and incur loss $f_t(\mathbf{y}_t)$. Let $\mathbf{g}_t = \frac{\eta}{\delta} f_t(\mathbf{y}_t) \mathbf{u}_t$.
- 5: Update $\mathbf{x}_{t+1} = \Pi_{\mathcal{K}_\delta} [\mathbf{x}_t - \eta \mathbf{g}_t]$.
- 6: **end for**

Theorem 6

FKM algorithm with parameters $\frac{D}{nT^{3/4}}$ and $\delta = \frac{1}{T^{1/4}}$ guarantees the following expected regret bound

$$\mathbb{E} \left[\sum_{t=1}^T f_t(\mathbf{x}_t) \right] - \sum_{t=1}^T f_t(\mathbf{x}^*) \leq 9nDGT^{3/4}. \quad (22)$$

Proof.

(Proof Sketch) The algorithm constructs an unbiased estimator of a convex function that is δ -far from the true cost function f_t . The estimator's variance scales with $\frac{1}{\delta^2}$. Selecting $\delta = \frac{1}{T^{1/4}}$ with a slightly less aggressive OGD, configured with stepsize $\eta = \frac{D}{nT^{3/4}}$; combined with the reduction in Lemma 4 ensures the above bound. \square

PAC (Probably Approximately Correct) learning and OCO

A learning algorithm has access to samples from an unknown distribution

$$(\mathbf{x}, y) \sim \mathcal{D}, \mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}. \quad (23)$$

The goal is to be able to predict y as a function of \mathbf{x} , i.e., to learn a mapping (hypothesis) $h: \mathcal{X} \rightarrow \mathcal{Y}$, that minimizes a prediction error according to a (usually convex) loss function $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ given as

$$\text{error}(h) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [l(h(\mathbf{x}), y)]. \quad (24)$$

Definition 7

(Agnostic PAC learning) The hypothesis class \mathcal{H} is agnostically PAC learnable with respect to loss function $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ if the following holds. There exists an algorithm \mathcal{A} that accepts $S_T = \{(x_t, y_t), t \in [T]\}$ and returns hypothesis $h_{\mathcal{A}(S_T)} \in \mathcal{H}$ that satisfies: for any $\epsilon, \delta > 0$ there exists a sufficiently large natural number $T = T(\epsilon, \delta)$ such that for any distribution \mathcal{D} over pairs (\mathbf{x}, y) and $T(\epsilon, \delta)$ samples from this distribution, it holds that with probability at least $1 - \delta$

$$\text{error}(h_{\mathcal{A}(S_T)}) \leq \min_{h \in \mathcal{H}} \text{error}(h) + \epsilon \quad (25)$$

Reduction: Learning \rightarrow OCO

- 1: Input: OCO algorithm \mathcal{A} , convex hypothesis class $\mathcal{H} \subseteq \mathbb{R}^d$, convex loss function ℓ .
- 2: Let $h_1 \leftarrow \mathcal{A}(\emptyset)$.
- 3: **for** $t = 1$ to T **do**
- 4: Draw labeled example $(\mathbf{x}_t, y_t) \sim \mathcal{D}$.
- 5: Let $f_t(h) = \ell(h(\mathbf{x}_t), y_t)$.
- 6: Update
$$h_{t+1} = \mathcal{A}(f_1, \dots, f_t).$$
- 7: **end for**
- 8: Return $\bar{h} = \frac{1}{T} \sum_{t=1}^T h_t$.

Theorem 8

Let \mathcal{A} be an OCO algorithm whose regret after T iterations is guaranteed to be bounded by $\text{regret}_T(\mathcal{A})$. Then for any $\delta > 0$, with probability at least $1 - \delta$, it holds that

$$\text{error}(\bar{h}) \leq \min_{h^* \in \mathcal{H}} \text{error}(h^*) + \frac{\text{regret}_T(\mathcal{A})}{T} + \sqrt{\frac{8 \log(2/\delta)}{T}}. \quad (26)$$

For $T = \mathcal{O}\left(\frac{1}{\epsilon^2} \log(1/\delta) + T_\epsilon(\mathcal{A})\right)$, where $T_\epsilon(\mathcal{A})$ is the integer T such that $\text{regret}_T(\mathcal{A})/T \leq \epsilon$, we have

$$\text{error}(\bar{h}) \leq \min_{h^* \in \mathcal{H}} \text{error}(h^*) + \epsilon. \quad (27)$$

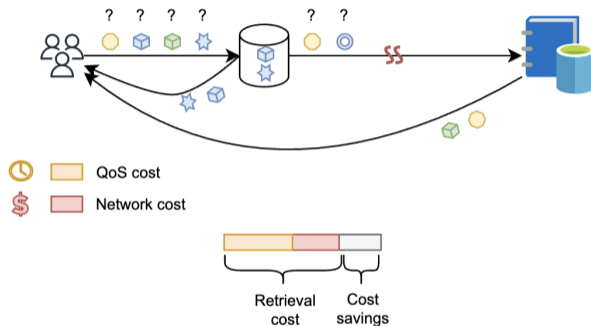
Let us consider a generalization of the expert problem.¹ The player has to choose k experts among n given experts. The number of configurations under this setting is $\binom{n}{k}$, so constructing a distribution $\Delta_{\binom{n}{k}}$ already prohibits having any efficient algorithms. Instead, we can keep track of the marginal probabilities instead, i.e., the decisions in

$$\Delta_{n,k} \triangleq \{ \mathbf{x} \in [0, 1]^n : \|\mathbf{x}\|_1 = k \}. \quad (28)$$

¹No-Regret Caching via Online Mirror Descent. ACM Transactions on Modeling and Performance Evaluation of Computing Systems. T. Si Salem, S. Ioannidis, G. Neglia.

k -Experts: Casting as a Caching Problem and Dimensionality Reduction

For $\mathbf{x} \in \Delta_{n,k}$, the component x_i corresponds to the probability of selecting file $i \in [n]$.

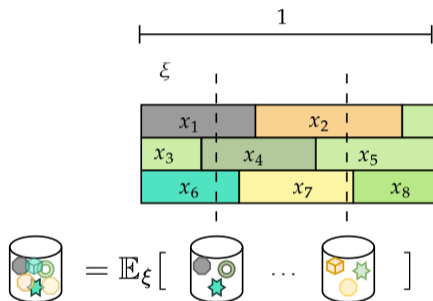


When a request batch \mathbf{r}_t arrives, the cache incurs the following cost:

$$f_{\mathbf{r}_t}(\mathbf{x}_t) = \sum_{i=1}^N w_i r_{t,i} (1 - x_{t,i}). \quad (29)$$

k -Experts: Recovering the Original Setting

We can construct a sampling scheme Ξ that outputs $z_t \in \Delta_{n,k} \cap \{0, 1\}^n$ such that $\mathbb{E}[z_t] = \mathbf{x}_t$.



This scheme is called Madow's sampling. Under Ξ the expected cost of z_t is identical to the cost of \mathbf{x}_t , i.e.,

$$\mathbb{E}_\Xi [f_{r_t}(z_t)] = f_{r_t}(\mathbf{x}_t). \quad (30)$$

The expected regret $\text{regret}_T(\mathcal{A}, \Xi)$ is the same as the regret of \mathcal{A} !

In practice, there is a cost associated with switching the decision from z_t to z_{t+1} , denoted by $\text{UC}_{r_t}(z_t, z_{t+1})$. In the fractional setup, gradient-based algorithms adapt their states proportionally to the learning rate, i.e., $\|\mathbf{x}_{t+1} - \mathbf{x}_t\| = \mathcal{O}(\eta)$. The total update cost measured in some norm $\|\cdot\|$ is given by

$$\sum_{t=1}^T \|\mathbf{x}_{t+1} - \mathbf{x}_t\| = \mathcal{O}(\eta T) = \mathcal{O}(\sqrt{T}). \quad (31)$$

The last equality is obtained for $\eta = \Theta\left(\frac{1}{\sqrt{T}}\right)$.

Randomness prevents this property to be transferred to the integral actions z_t .

k -Experts: Independent Sampling

When considering the extended regret ($E\text{-Regret}_T(\mathcal{A}, \Xi)$) we lose immediately the regret guarantee:

Theorem

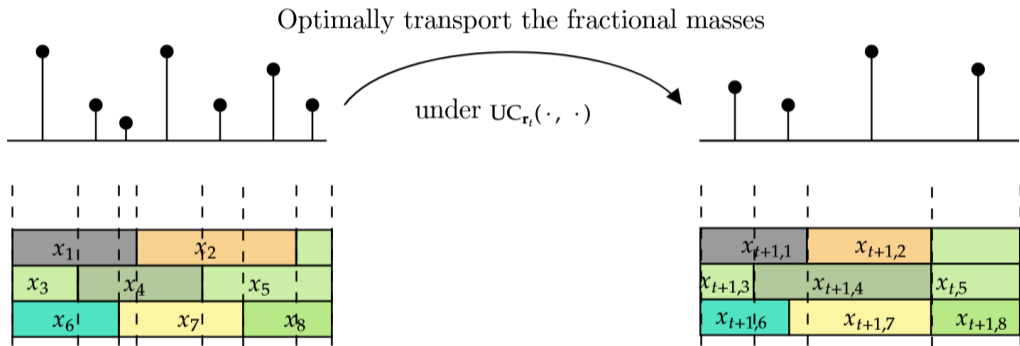
Any randomized caching policy constructed by an online policy \mathcal{A} combined with online independent rounding as Ξ leads to $\Omega(T)$ $E\text{-Regret}(\mathcal{A}, \Xi)$.

Imposing dependence (coupling) between the two consecutive random states may significantly reduce the expected update cost.

Optimizing for all possible coupling is a **min-cost flow problem**, or when establishing flows between distributions it is an **optimal transport problem**.

$$\begin{aligned} \mathbf{f} = & \arg \min_{[f_{i,j}]_{(i,j) \in [|\mathbf{p}_t|] \times [|\mathbf{p}_{t+1}|]}} \mathbb{E} [\text{UC}(\mathbf{z}_t, \mathbf{z}_{t+1})] = \sum_{i=1}^{|\mathbf{p}_t|} \sum_{j=1}^{|\mathbf{p}_{t+1}|} \text{UC}_{r_t}(\zeta_t^i, \zeta_{t+1}^j) f_{i,j} \\ \text{s.t.} \quad & \sum_{j=1}^{|\mathbf{p}_{t+1}|} f_{i,j} = p_{t,i}, \quad \sum_{i=1}^{|\mathbf{p}_t|} f_{i,j} = p_{t+1,j}, \quad f_{i,j} \in [0, 1], \forall (i, j) \in [|\mathbf{p}_t|] \times [|\mathbf{p}_{t+1}|]. \end{aligned}$$

k -Experts: Optimal Transport

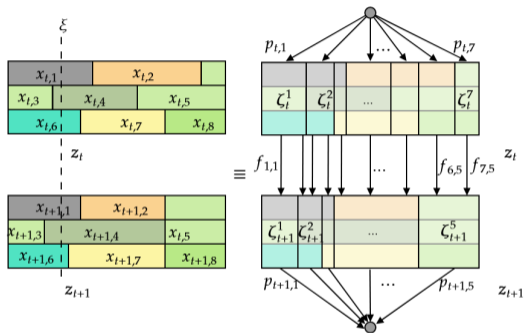


Remark

We prove that this scheme selected as Ξ , coupled with a no-regret policy \mathcal{A} , has sublinear extended regret guarantee. **However, it has a time-complexity $\mathcal{O}(N^3)$.**

k -Experts: Simpler Approach

Online Coupled Rounding

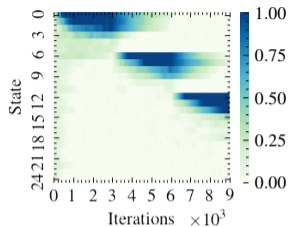


Theorem

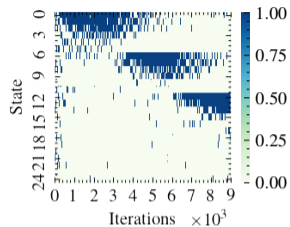
A no-regret policy \mathcal{A} combined with online coupled rounding Ξ has $\mathcal{O}(\sqrt{T})$ E-Regret $_T(\mathcal{A}, \Xi)$.

Online Coupled Rounding has linear time complexity.

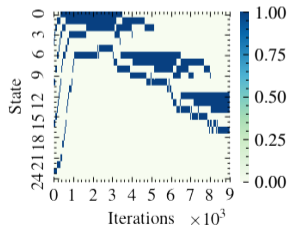
k -Experts: Coupling Schemes a Qualitative Evaluation



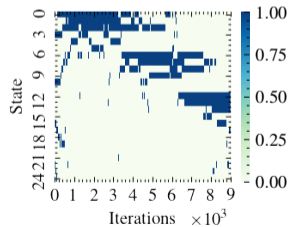
(a) Fractional cache states



(b) Online independent rounding



(c) Online coupled rounding



(d) Online optimally-coupled rounding