

# **Stat 260/CS 294-102. Learning in Sequential Decision Problems.**

**Peter Bartlett**

## 1. Linear bandits.

- Full information: mirror descent.
- Bandit information: stochastic mirror descent.

## Full information online prediction games

- Repeated game:

Strategy plays  $a_t \in \mathcal{A}$

Adversary reveals  $\ell_t \in \mathcal{L}$

- Aim to minimize **regret**:

$$R_n = \sum_{t=1}^n \ell_t(a_t) - \min_{a \in \mathcal{A}} \sum_{t=1}^n \ell_t(a).$$

## Online Convex Optimization

- Choosing  $a_t$  to minimize past losses can fail.
- The strategy must avoid overfitting.
- First approach: gradient steps.  
Stay close to previous decisions, but move in a direction of improvement.

# Online Convex Optimization

1. Gradient algorithm.
2. Regularized minimization
  - Bregman divergence
  - Regularized minimization  $\Leftrightarrow$  minimizing latest loss and divergence from previous decision
  - Constrained minimization equivalent to unconstrained plus Bregman projection
  - Linearization
  - Mirror descent
3. Regret bound

## Online Convex Optimization: Gradient Method

$$a_1 \in \mathcal{A},$$
$$a_{t+1} = \Pi_{\mathcal{A}}(a_t - \eta \nabla \ell_t(a_t)),$$

where  $\Pi_{\mathcal{A}}$  is the Euclidean projection on  $\mathcal{A}$ ,

$$\Pi_{\mathcal{A}}(x) = \arg \min_{a \in \mathcal{A}} \|x - a\|.$$

**Theorem:** For  $G = \max_t \|\nabla \ell_t(a_t)\|$  and  $D = \text{diam}(\mathcal{A})$ , the gradient strategy with  $\eta = D/(G\sqrt{n})$  has regret satisfying

$$R_n \leq GD\sqrt{n}.$$

## Online Convex Optimization: Gradient Method

**Example:** (2-ball, 2-ball)

$\mathcal{A} = \{a \in \mathbb{R}^d : \|a\| \leq 1\}$ ,  $\mathcal{L} = \{a \mapsto v \cdot a : \|v\| \leq 1\}$ .  $D = 2$ ,  $G \leq 1$ .

Regret is no more than  $2\sqrt{n}$ .

(And  $O(\sqrt{n})$  is optimal.)

**Example:** (1-ball,  $\infty$ -ball)

$\mathcal{A} = \Delta(k)$ ,  $\mathcal{L} = \{a \mapsto v \cdot a : \|v\|_\infty \leq 1\}$ .

$D = 2$ ,  $G \leq \sqrt{k}$ .

Regret is no more than  $2\sqrt{kn}$ .

Since competing with the whole simplex is equivalent to competing with the vertices (experts) for linear losses, this is worse than exponential weights ( $\sqrt{k}$  versus  $\log k$ ).

## Gradient Method: Proof

$$\begin{aligned}\text{Define} \quad \tilde{a}_{t+1} &= a_t - \eta \nabla \ell_t(a_t), \\ a_{t+1} &= \Pi_{\mathcal{A}}(\tilde{a}_{t+1}).\end{aligned}$$

Fix  $a \in \mathcal{A}$  and consider the measure of progress  $\|a_t - a\|$ .

$$\begin{aligned}\|a_{t+1} - a\|^2 &\leq \|\tilde{a}_{t+1} - a\|^2 \\ &= \|a_t - a\|^2 + \eta^2 \|\nabla \ell_t(a_t)\|^2 - 2\eta \nabla \ell_t(a_t) \cdot (a_t - a).\end{aligned}$$

By convexity,

$$\begin{aligned}\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) &\leq \sum_{t=1}^n \nabla \ell_t(a_t) \cdot (a_t - a) \\ &\leq \frac{\|a_1 - a\|^2 - \|a_{n+1} - a\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^n \|\nabla \ell_t(a_t)\|^2\end{aligned}$$

## Online Convex Optimization: A Regularization Viewpoint

- Suppose  $\ell_t$  is linear:  $\ell_t(a) = g_t \cdot a$ .
- Suppose  $\mathcal{A} = \mathbb{R}^d$ .
- Then minimizing the regularized criterion

$$a_{t+1} = \arg \min_{a \in \mathcal{A}} \left( \eta \sum_{s=1}^t \ell_s(a) + \frac{1}{2} \|a\|^2 \right)$$

corresponds to the gradient step

$$a_{t+1} = a_t - \eta \nabla \ell_t(a_t).$$

## Online Convex Optimization: Regularization

### Regularized minimization

Consider the family of strategies of the form:

$$a_{t+1} = \arg \min_{a \in \mathcal{A}} \left( \eta \sum_{s=1}^t \ell_s(a) + R(a) \right).$$

The regularizer  $R : \mathbb{R}^d \rightarrow \mathbb{R}$  is strictly convex and differentiable.

- $R$  keeps the sequence of  $a_t$ s stable: it diminishes  $\ell_t$ 's influence.
- We can view the choice of  $a_{t+1}$  as trading off two competing forces: making  $\ell_t(a_{t+1})$  small, and keeping  $a_{t+1}$  close to  $a_t$ .
- This is a perspective that motivated many algorithms in the literature.

## Properties of Regularization Methods

In the unconstrained case ( $\mathcal{A} = \mathbb{R}^d$ ), regularized minimization is equivalent to minimizing the latest loss and the distance to the previous decision. The appropriate notion of distance is the **Bregman divergence**

$D_{\Phi_{t-1}}$ :

Define

$$\begin{aligned}\Phi_0 &= R, \\ \Phi_t &= \Phi_{t-1} + \eta \ell_t,\end{aligned}$$

so that

$$\begin{aligned}a_{t+1} &= \arg \min_{a \in \mathcal{A}} \left( \eta \sum_{s=1}^t \ell_s(a) + R(a) \right) \\ &= \arg \min_{a \in \mathcal{A}} \Phi_t(a).\end{aligned}$$

## Bregman Divergence

**Definition:** For a strictly convex, differentiable  $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$ , the Bregman divergence wrt  $\Phi$  is defined, for  $a, b \in \mathbb{R}^d$ , as

$$D_{\Phi}(a, b) = \Phi(a) - (\Phi(b) + \nabla\Phi(b) \cdot (a - b)).$$

$D_{\Phi}(a, b)$  is the difference between  $\Phi(a)$  and the value at  $a$  of the linear approximation of  $\Phi$  about  $b$ . (PICTURE)

## Bregman Divergence

**Example:** For  $a \in \mathbb{R}^d$ , the squared euclidean norm,  $\Phi(a) = \frac{1}{2}\|a\|^2$ , has

$$\begin{aligned} D_{\Phi}(a, b) &= \frac{1}{2}\|a\|^2 - \left( \frac{1}{2}\|b\|^2 + b \cdot (a - b) \right) \\ &= \frac{1}{2}\|a - b\|^2, \end{aligned}$$

the squared euclidean norm.

## Bregman Divergence

**Example:** For  $a \in [0, \infty)^d$ , the unnormalized negative entropy,  $\Phi(a) = \sum_{i=1}^d a_i (\ln a_i - 1)$ , has

$$\begin{aligned} D_{\Phi}(a, b) &= \sum_i (a_i(\ln a_i - 1) - b_i(\ln b_i - 1) - \ln b_i(a_i - b_i)) \\ &= \sum_i \left( a_i \ln \frac{a_i}{b_i} + b_i - a_i \right), \end{aligned}$$

the unnormalized KL divergence.

Thus, for  $a \in \Delta^d$ ,  $\Phi(a) = \sum_i a_i \ln a_i$  has

$$D_{\Phi}(a, b) = \sum_i a_i \ln \frac{a_i}{b_i}.$$

## Bregman Divergence

When the domain of  $\Phi$  is  $\mathcal{S} \subset \mathbb{R}^d$ , in addition to differentiability and strict convexity, we make some more assumptions:

- $\mathcal{S}$  is closed, and its interior is convex.
- For a sequence approaching the boundary of  $\mathcal{S}$ ,  $\|\nabla\Phi(a_n)\| \rightarrow \infty$ .

We say that such a  $\Phi$  is a *Legendre function*.

## Bregman Divergence Properties

1.  $D_{\Phi} \geq 0$ ,  $D_{\Phi}(a, a) = 0$ .
2.  $D_{A+B} = D_A + D_B$ .
3. For  $\ell$  linear,  $D_{\Phi+\ell} = D_{\Phi}$ .
4. *Bregman projection*,  $\Pi_{\mathcal{A}}^{\Phi}(b) = \arg \min_{a \in \mathcal{A}} D_{\Phi}(a, b)$  is uniquely defined for closed, convex  $\mathcal{A} \subset \mathcal{S}$  (that intersects the interior of  $\mathcal{S}$ ).
5. *Generalized Pythagoras*: for closed, convex  $\mathcal{A}$ ,  $a^* = \Pi_{\mathcal{A}}^{\Phi}(b)$ ,  $a \in \mathcal{A}$ ,  
 $D_{\Phi}(a, b) \geq D_{\Phi}(a, a^*) + D_{\Phi}(a^*, b)$ .
6.  $\nabla_a D_{\Phi}(a, b) = \nabla \Phi(a) - \nabla \Phi(b)$ .
7. For  $\Phi^*$  the Legendre dual of  $\Phi$ ,

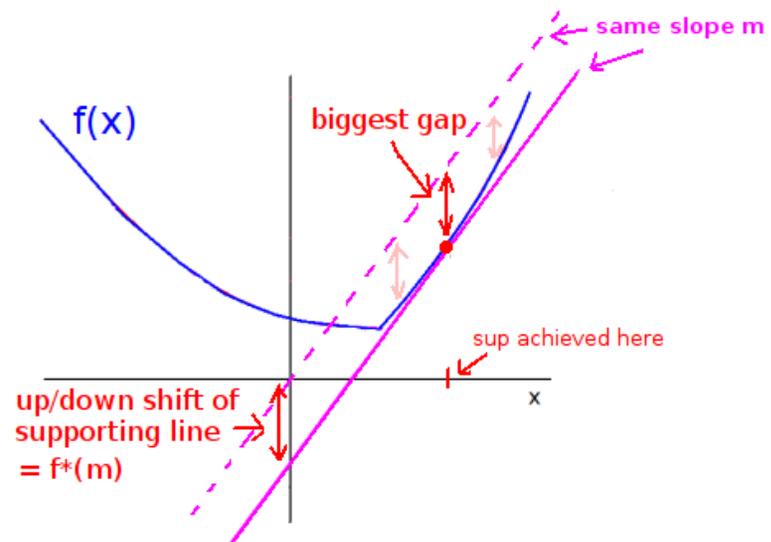
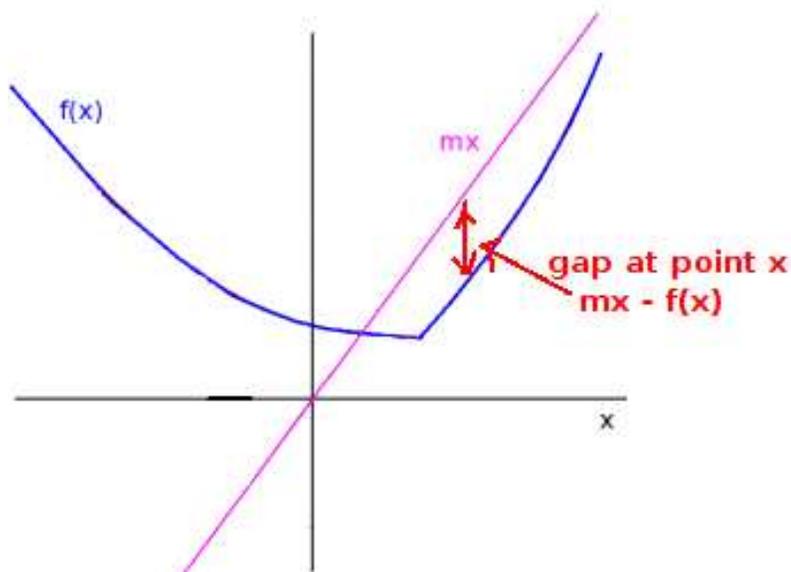
$$\nabla \Phi^* = (\nabla \Phi)^{-1},$$

$$D_{\Phi}(a, b) = D_{\Phi^*}(\nabla \Phi(b), \nabla \Phi(a)).$$

# Legendre Dual

Here, for a Legendre function  $\Phi : \mathcal{S} \rightarrow \mathbb{R}$ , we define the Legendre dual as

$$\Phi^*(u) = \sup_{v \in \mathcal{S}} (u \cdot v - \Phi(v)).$$



(<http://maze5.net/>)

## Legendre Dual

Properties:

- $\Phi^*$  is Legendre.
- $\text{dom}(\Phi^*) = \nabla\Phi(\text{int dom } \Phi)$ .
- $\nabla\Phi^* = (\nabla\Phi)^{-1}$ .
- $D_{\Phi}(a, b) = D_{\Phi^*}(\nabla\Phi(b), \nabla\Phi(a))$ .
- $\Phi^{**} = \Phi$ .

## Properties of Regularization Methods

In the unconstrained case ( $\mathcal{A} = \mathbb{R}^d$ ), regularized minimization is equivalent to minimizing the latest loss and the distance (Bregman divergence) to the previous decision.

**Theorem:** Define  $\tilde{a}_1$  via  $\nabla R(\tilde{a}_1) = 0$ , and set

$$\tilde{a}_{t+1} = \arg \min_{a \in \mathbb{R}^d} (\eta \ell_t(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t)).$$

Then

$$\tilde{a}_{t+1} = \arg \min_{a \in \mathbb{R}^d} \left( \eta \sum_{s=1}^t \ell_s(a) + R(a) \right).$$

## Properties of Regularization Methods

*Proof.* By the definition of  $\Phi_t$ ,

$$\eta\ell_t(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t) = \Phi_t(a) - \Phi_{t-1}(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t).$$

The derivative wrt  $a$  is

$$\begin{aligned} \nabla\Phi_t(a) - \nabla\Phi_{t-1}(a) + \nabla_a D_{\Phi_{t-1}}(a, \tilde{a}_t) \\ = \nabla\Phi_t(a) - \nabla\Phi_{t-1}(a) + \nabla\Phi_{t-1}(a) - \nabla\Phi_{t-1}(\tilde{a}_t) \end{aligned}$$

Setting to zero shows that

$$\nabla\Phi_t(\tilde{a}_{t+1}) = \nabla\Phi_{t-1}(\tilde{a}_t) = \cdots = \nabla\Phi_0(\tilde{a}_1) = \nabla R(\tilde{a}_1) = 0,$$

So  $\tilde{a}_{t+1}$  minimizes  $\Phi_t$ . □

## Properties of Regularization Methods

Constrained minimization is equivalent to unconstrained minimization, followed by Bregman projection:

**Theorem:** For

$$a_{t+1} = \arg \min_{a \in \mathcal{A}} \Phi_t(a),$$

$$\tilde{a}_{t+1} = \arg \min_{a \in \mathbb{R}^d} \Phi_t(a),$$

we have

$$a_{t+1} = \Pi_{\mathcal{A}}^{\Phi_t}(\tilde{a}_{t+1}).$$

## Properties of Regularization Methods

*Proof.* Let  $a'_{t+1}$  denote  $\Pi_{\mathcal{A}}^{\Phi_t}(\tilde{a}_{t+1})$ . First, by definition of  $a_{t+1}$ ,

$$\Phi_t(a_{t+1}) \leq \Phi_t(a'_{t+1}).$$

Conversely,

$$D_{\Phi_t}(a'_{t+1}, \tilde{a}_{t+1}) \leq D_{\Phi_t}(a_{t+1}, \tilde{a}_{t+1}).$$

But  $\nabla \Phi_t(\tilde{a}_{t+1}) = 0$ , so

$$D_{\Phi_t}(a, \tilde{a}_{t+1}) = \Phi_t(a) - \Phi_t(\tilde{a}_{t+1}).$$

Thus,  $\Phi_t(a'_{t+1}) \leq \Phi_t(a_{t+1})$ . □

## Properties of Regularization Methods

**Example:** For **linear**  $\ell_t$ , regularized minimization is equivalent to minimizing the last loss plus the Bregman divergence **wrt**  $R$  to the previous decision:

$$\begin{aligned} & \arg \min_{a \in \mathcal{A}} \left( \eta \sum_{s=1}^t \ell_s(a) + R(a) \right) \\ &= \Pi_{\mathcal{A}}^R \left( \arg \min_{a \in \mathbb{R}^d} (\eta \ell_t(a) + D_R(a, \tilde{a}_t)) \right), \end{aligned}$$

because adding a linear function to  $\Phi$  does not change  $D_{\Phi}$ .

## Linear Loss

We can replace  $\ell_t$  by  $\nabla \ell_t(a_t)$ , and this leads to an upper bound on regret.

Thus, for convex losses, we can work with **linear**  $\ell_t$ .

## Regularization Methods: Mirror Descent

Regularized minimization for linear losses can be viewed as **mirror descent**—taking a gradient step in a dual space:

**Theorem:** The decisions

$$\tilde{a}_{t+1} = \arg \min_{a \in \mathbb{R}^d} \left( \eta \sum_{s=1}^t g_s \cdot a + R(a) \right)$$

can be written

$$\tilde{a}_{t+1} = (\nabla R)^{-1} (\nabla R(\tilde{a}_t) - \eta g_t).$$

This corresponds to first mapping from  $\tilde{a}_t$  through  $\nabla R$ , then taking a step in the direction  $-g_t$ , then mapping back through  $(\nabla R)^{-1} = \nabla R^*$  to  $\tilde{a}_{t+1}$ .

## Regularization Methods: Mirror Descent

*Proof.* For the unconstrained minimization, we have

$$\nabla R(\tilde{a}_{t+1}) = -\eta \sum_{s=1}^t g_s,$$

$$\nabla R(\tilde{a}_t) = -\eta \sum_{s=1}^{t-1} g_s,$$

so  $\nabla R(\tilde{a}_{t+1}) = \nabla R(\tilde{a}_t) - \eta g_t$ , which can be written

$$\tilde{a}_{t+1} = \nabla R^{-1} (\nabla R(\tilde{a}_t) - \eta g_t).$$

□

## Mirror Descent

Given:

compact, convex  $\mathcal{A} \subseteq \mathbb{R}^d$ , closed, convex  $\mathcal{S} \supset \mathcal{A}$ ,  $\eta > 0$ ,  $\mathcal{S} \supset \mathcal{A}$ ,  
Legendre  $R : \mathcal{S} \rightarrow \mathbb{R}$ . Set  $a_1 \in \arg \min_{a \in \mathcal{A}} R(a)$ .

For round  $t$ :

1. Play  $a_t$ ; observe  $\ell_t \in \mathbb{R}^d$ .
2.  $w_{t+1} = \nabla R^* (\nabla R(a_t) - \eta \nabla \ell_t(a_t))$ .
3.  $a_{t+1} = \arg \min_{a \in \mathcal{A}} D_R(a, w_{t+1})$ .

## Exponential weights as mirror descent

For  $\mathcal{A} = \Delta(k)$  and  $R(a) = \sum_{i=1}^k (a_i \log a_i - a_i)$ , this reduces to exponential weights:

$$\nabla R(u)_i = \log a_i,$$

$$R^*(u) = \sum_i e^{u_i},$$

$$\nabla R^*(u)_i = \exp(u_i),$$

$$\nabla R(w_{t+1})_i = \log(w_{t+1,i}) = \log a_{t,i} - \eta \nabla \ell_t(a_t)_i,$$

$$w_{t+1,i} = a_{t,i} \exp(-\eta \nabla \ell_t(a_t)_i),$$

$$D_R(a, b) = \sum_i \left( a_i \log \frac{a_i}{b_i} + b_i - a_i \right),$$

$$a_{t+1,i} \propto w_{t+1,i}.$$

## Mirror descent regret

**Theorem:** Suppose that, for all  $a \in \mathcal{A} \cap \text{int}(\mathcal{S})$ ,  $\ell \in \mathcal{L}$ ,  $\nabla R(a) - \eta \nabla \ell(a) \in \nabla R(\text{int}(\mathcal{S}))$ . For any  $a \in \mathcal{A}$ ,

$$\begin{aligned} & \sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) \\ & \leq \frac{1}{\eta} \left( R(a) - R(a_1) + \sum_{t=1}^n D_{R^*} \left( \nabla R(a_t) - \eta \nabla \ell_t(a_t), \nabla R(a_t) \right) \right). \end{aligned}$$

*Proof:* Fix  $a \in \mathcal{A}$ . Since the  $\ell_t$  are convex,

$$\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) \leq \sum_{t=1}^n \nabla \ell_t(a_t)^T (a_t - a).$$

## Mirror descent regret: proof

The choice of  $w_{t+1}$  and the fact that  $\nabla R^{-1} = \nabla R^*$  show that

$$\nabla R(w_{t+1}) = \nabla R(a_t) - \eta \nabla \ell_t(a_t).$$

Hence,

$$\begin{aligned} \eta \nabla \ell_t(a_t)^T (a_t - a) &= (a - a_t)^T (\nabla R(w_{t+1}) - \nabla R(a_t)) \\ &= D_R(a, a_t) + D_R(a_t, w_{t+1}) - D_R(a, w_{t+1}). \end{aligned}$$

Generalized Pythagoras' inequality shows that the projection  $a_{t+1}$  satisfies

$$D_R(a, w_{t+1}) \geq D_R(a, a_{t+1}) + D_R(a_{t+1}, w_{t+1}).$$

## Mirror descent regret: proof

$$\begin{aligned} & \eta \sum_{t=1}^n \nabla \ell_t(a_t)^T (a_t - a) \\ & \leq \sum_{t=1}^n \left( D_R(a, a_t) + D_R(a_t, w_{t+1}) - D_R(a, w_{t+1}) \right. \\ & \quad \left. - D_R(a, a_{t+1}) - D_R(a_{t+1}, w_{t+1}) \right) \\ & = D_R(a, a_1) - D_R(a, a_{n+1}) + \sum_{t=1}^n (D_R(a_t, w_{t+1}) - D_R(a_{t+1}, w_{t+1})) \\ & \leq D_R(a, a_1) + \sum_{t=1}^n D_R(a_t, w_{t+1}). \end{aligned}$$

## Mirror descent regret: proof

$$\begin{aligned} &= D_R(a, a_1) + \sum_{t=1}^n D_{R^*}(\nabla R(w_{t+1}), \nabla R(a_t)) \\ &= D_R(a, a_1) + \sum_{t=1}^n D_{R^*}(\nabla R(a_t) - \eta \nabla \ell_t(a_t), \nabla R(a_t)) \\ &= R(a) - R(a_1) + \sum_{t=1}^n D_{R^*}(\nabla R(a_t) - \eta \nabla \ell_t(a_t), \nabla R(a_t)). \end{aligned}$$

## Linear bandit setting

- See only  $\ell_t(a_t)$ ;  $\nabla \ell_t(a_t)$  is unseen.
- Instead of  $a_t$ , strategy plays a noisy version,  $x_t$ .
- Strategy uses  $\ell_t(x_t)$  to give an unbiased estimate of  $\nabla \ell_t(a_t)$ .

## Stochastic mirror descent

Given:

compact, convex  $\mathcal{A} \subseteq \mathbb{R}^d$ ,  $\eta > 0$ ,  $\mathcal{S} \supset \mathcal{A}$ , Legendre  $R : \mathcal{S} \rightarrow \mathbb{R}$ .

Set  $a_1 \in \arg \min_{a \in \mathcal{A}} R(a)$ .

For round  $t$ :

1. Play **noisy version**  $x_t$  of  $a_t$ ; observe  $\ell_t(x_t)$ .
2. Compute estimate  $\tilde{g}_t$  of  $\nabla \ell_t(a_t)$ .
3.  $w_{t+1} = \nabla R^* (\nabla R(a_t) - \eta \tilde{g}_t)$ .
4.  $a_{t+1} = \arg \min_{a \in \mathcal{A}} D_R(a, w_{t+1})$ .

## Regret of stochastic mirror descent

**Theorem:** Suppose that, for all  $a \in \mathcal{A} \cap \text{int}(\mathcal{S})$  and linear  $\ell \in \mathcal{L}$ ,  $\mathbb{E}[\tilde{g}_t | a_t] = \nabla \ell_t(a_t)$  and  $\nabla R(a) - \eta \tilde{g}_t(a) \in \nabla R(\text{int}(\mathcal{S}))$ .

For any  $a \in \mathcal{A}$ ,

$$\begin{aligned} & \sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) \\ & \leq \frac{1}{\eta} \left( R(a) - R(a_1) + \sum_{t=1}^n \mathbb{E} D_{R^*} \left( \nabla R(a_t) - \eta \tilde{g}_t, \nabla R(a_t) \right) \right) \\ & \quad + \sum_{t=1}^n \mathbb{E} [\| \|a_t - \mathbb{E}[x_t | a_t]\| \| \tilde{g}_t \|_*]. \end{aligned}$$

## Regret: proof

$$\begin{aligned} & \mathbb{E} \sum_{t=1}^n (\ell_t(x_t) - \ell_t(a)) \\ &= \mathbb{E} \sum_{t=1}^n (\ell_t(x_t) - \ell_t(a_t) + \ell_t(a_t) - \ell_t(a)) \\ &= \mathbb{E} \sum_{t=1}^n (\mathbb{E} [\ell_t^T(x_t - a_t) | a_t] + \ell_t(a_t) - \ell_t(a)) \\ &\leq \mathbb{E} \sum_{t=1}^n \|a_t - \mathbb{E}[x_t | a_t]\| \|\tilde{g}_t\|_* + \mathbb{E} \sum_{t=1}^n \nabla \ell_t(a_t)^T (a_t - a) \\ &= \mathbb{E} \sum_{t=1}^n \|a_t - \mathbb{E}[x_t | a_t]\| \|\tilde{g}_t\|_* + \mathbb{E} \sum_{t=1}^n \tilde{g}_t^T (a_t - a). \end{aligned}$$

## Regret: proof

Applying the regret bound for the (random) linear losses  $a \mapsto \tilde{g}_t^T a$  gives

$$\begin{aligned} &\leq \mathbb{E} \sum_{t=1}^n \|a_t - \mathbb{E}[x_t | a_t]\| \|\tilde{g}_t\|_* \\ &\quad + \frac{1}{\eta} \left( R(a) - R(a_1) + \sum_{t=1}^n \mathbb{E} D_{R^*}(\nabla R(a_t) - \eta \tilde{g}_t, \nabla R(a_t)) \right). \end{aligned}$$

## Regret: Euclidean ball

Consider  $B = \{a \in \mathbb{R}^d : \|a\| \leq 1\}$  (with the Euclidean norm).

Ingredients:

1. Distribution of  $x_t$ , given  $a_t$ :

$$x_t = \xi_t \frac{a_t}{\|a_t\|} + (1 - \xi_t) \epsilon_t e_{I_t},$$

where  $\xi_t$  is Bernoulli( $\|a_t\|$ ),  $\epsilon_t$  is uniform  $\pm 1$ , and  $I_t$  is uniform on  $\{1, \dots, d\}$ , so  $\mathbb{E}[x_t | a_t] = a_t$ .

2. Estimate  $\tilde{\ell}_t$  of loss  $\ell_t$ :

$$\tilde{\ell}_t = d \frac{1 - \xi_t}{1 - \|a_t\|} x_t^T \ell_t x_t,$$

so  $\mathbb{E}[\tilde{\ell}_t | a_t] = \ell_t$ .

## Regret: Euclidean ball

**Theorem:** Consider stochastic mirror descent on  $\mathcal{A} = (1 - \gamma)B$ , with these choices and  $R(a) = -\log(1 - \|a\|) - \|a\|$ . Then for  $\eta d \leq 1/2$ ,

$$\bar{R}_n \leq \gamma n + \frac{\log(1/\gamma)}{\eta} + \eta \sum_{t=1}^n \mathbb{E} \left[ (1 - \|a_t\|) \|\tilde{\ell}_t\|^2 \right].$$

For  $\gamma = 1/\sqrt{n}$  and  $\eta = \sqrt{\log n / (2nd)}$ ,

$$\bar{R}_n \leq 3\sqrt{dn \log n}.$$