

# Learning Methods for Online Prediction Problems

**Peter Bartlett**  
Statistics and EECS  
UC Berkeley

## Course Synopsis

- ▶ A finite comparison class:  $\mathcal{A} = \{1, \dots, m\}$ .
  1. “Prediction with expert advice.”
  2. With perfect predictions:  $\log m$  regret.
  3. Exponential weights strategy:  $\sqrt{n \log m}$  regret.
  4. Refinements and extensions.
  5. Statistical prediction with a finite class.
- ▶ Converting online to batch.
- ▶ Online convex optimization.
- ▶ Log loss.

## Online to Batch Conversion

- ▶ Suppose we have an online strategy that, given observations  $l_1, \dots, l_{t-1}$ , produces  $a_t = A(l_1, \dots, l_{t-1})$ .
- ▶ Can we convert this to a method that is suitable for a probabilistic setting? That is, if the  $l_t$  are chosen i.i.d., can we use  $A$ 's choices  $a_t$  to come up with a  $\hat{a} \in \mathcal{A}$  so that

$$\mathbf{E}l_1(\hat{a}) - \min_{a \in \mathcal{A}} \mathbf{E}l_1(a)$$

is small?

- ▶ Consider the following simple randomized method:
  1. Pick  $T$  uniformly from  $\{0, \dots, n\}$ .
  2. Let  $\hat{a} = A(l_{T+1}, \dots, l_n)$ .

## Online to Batch Conversion

### Theorem

If  $A$  has a regret bound of  $C_{n+1}$  for sequences of length  $n + 1$ , then for any stationary process generating the  $\ell_1, \dots, \ell_{n+1}$ , this method satisfies

$$\mathbf{E} \ell_{n+1}(\hat{a}) - \min_{a \in \mathcal{A}} \mathbf{E} \ell_n(a) \leq \frac{C_{n+1}}{n+1}.$$

(Notice that the expectation averages also over the randomness of the method.)

Proof.

$$\begin{aligned}
 \mathbf{E}l_{n+1}(\hat{a}) &= \mathbf{E}l_{n+1}(\mathbf{A}(l_{T+1}, \dots, l_n)) \\
 &= \mathbf{E} \frac{1}{n+1} \sum_{t=0}^n l_{n+1}(\mathbf{A}(l_{t+1}, \dots, l_n)) \\
 &= \mathbf{E} \frac{1}{n+1} \sum_{t=0}^n l_{n-t+1}(\mathbf{A}(l_1, \dots, l_{n-t})) \\
 &= \mathbf{E} \frac{1}{n+1} \sum_{t=1}^{n+1} l_t(\mathbf{A}(l_1, \dots, l_{t-1})) \\
 &\leq \mathbf{E} \frac{1}{n+1} \left( \min_a \sum_{t=1}^{n+1} l_t(a) + C_{n+1} \right) \\
 &\leq \min_a \mathbf{E}l_t(a) + \frac{C_{n+1}}{n+1}.
 \end{aligned}$$

## Online to Batch Conversion

- ▶ The theorem is for the expectation over the randomness of the method.
- ▶ For a high probability result, we could
  1. Choose  $\hat{a} = \frac{1}{n} \sum_{t=1}^n a_t$ , provided  $\mathcal{A}$  is convex and the  $\ell_t$  are all convex.
  2. Choose

$$\hat{a} = \arg \min_{a_t} \left( \frac{1}{n-t} \sum_{s=t+1}^n \ell_s(a_t) + c \sqrt{\frac{\log(n/\delta)}{n-t}} \right).$$

In both cases, the analysis involves concentration of martingale sequences.

The second (more general) approach does not recover the  $C_n/n$  result: the penalty has the wrong form when  $C_n = o(\sqrt{n})$ .

## Online to Batch Conversion

Key Point:

- ▶ An online strategy with regret bound  $C_n$  can be converted to a batch method.  
The regret per trial in the probabilistic setting is bounded by the regret per trial in the adversarial setting.

## Course Synopsis

- ▶ A finite comparison class:  $\mathcal{A} = \{1, \dots, m\}$ .
- ▶ Converting online to batch.
- ▶ Online convex optimization.
  1. Problem formulation
  2. Empirical minimization fails.
  3. Gradient algorithm.
  4. Regularized minimization
  5. Regret bounds
- ▶ Log loss.



1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization
  - ▶ Bregman divergence
  - ▶ Regularized minimization equivalent to minimizing latest loss and divergence from previous decision
  - ▶ Constrained minimization equivalent to unconstrained plus Bregman projection
  - ▶ Linearization
  - ▶ Mirror descent
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Seeing the future
  - ▶ Strong convexity
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions

## Online Convex Optimization

- ▶  $\mathcal{A}$  = convex subset of  $\mathbb{R}^d$ .
- ▶  $\mathcal{L}$  = set of convex real functions on  $\mathcal{A}$ .

For example,

$$\ell_t(\mathbf{a}) = (x_t \cdot \mathbf{a} - y_t)^2.$$

$$\ell_t(\mathbf{a}) = |x_t \cdot \mathbf{a} - y_t|.$$

## Online Convex Optimization: Example

Choosing  $a_t$  to minimize past losses,

$a_t = \arg \min_{a \in \mathcal{A}} \sum_{s=1}^{t-1} \ell_s(a)$ , can fail.

(‘fictitious play,’ ‘follow the leader’)

- ▶ Suppose  $\mathcal{A} = [-1, 1]$ ,  $\mathcal{L} = \{a \mapsto v \cdot a : |v| \leq 1\}$ .
- ▶ Consider the following sequence of losses:

$$\begin{array}{ll} a_1 = 0, & \ell_1(a) = \frac{1}{2}a, \\ a_2 = -1, & \ell_2(a) = -a, \\ a_3 = 1, & \ell_3(a) = a, \\ a_4 = -1, & \ell_4(a) = -a, \\ a_5 = 1, & \ell_5(a) = a, \\ \vdots & \vdots \end{array}$$

- ▶  $a^* = 0$  shows  $L_n^* \leq 0$ , but  $\hat{L}_n = n - 1$ .

## Online Convex Optimization: Example

- ▶ Choosing  $a_t$  to minimize past losses can fail.
- ▶ The strategy must avoid overfitting, just as in probabilistic settings.
- ▶ Similar approaches (regularization; Bayesian inference) are applicable in the online setting.
- ▶ First approach: gradient steps.  
Stay close to previous decisions, but move in a direction of improvement.

# Online Convex Optimization: Gradient Method

$$\begin{aligned} \mathbf{a}_1 &\in \mathcal{A}, \\ \mathbf{a}_{t+1} &= \Pi_{\mathcal{A}}(\mathbf{a}_t - \eta \nabla \ell_t(\mathbf{a}_t)), \end{aligned}$$

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  $G = \max_t \|\nabla \ell_t(\mathbf{a}_t)\|$  and  $D = \text{diam}(\mathcal{A})$ , the gradient strategy with  $\eta = D/(G\sqrt{n})$  has regret satisfying

$$\hat{L}_n - L_n^* \leq GD\sqrt{n}.$$

## Online Convex Optimization: Gradient Method

### Theorem

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

$$\hat{L}_n - L_n^* \leq GD\sqrt{n}.$$

### Example

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

$D = 2$ ,  $G \leq 1$ .

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

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

# Online Convex Optimization: Gradient Method

## 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

$$\hat{L}_n - L_n^* \leq GD\sqrt{n}.$$

## Example

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

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

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

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

Proof.

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

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

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

By convexity,

$$\begin{aligned}\sum_{t=1}^n (\ell_t(\mathbf{a}_t) - \ell_t(\mathbf{a})) &\leq \sum_{t=1}^n \nabla \ell_t(\mathbf{a}_t) \cdot (\mathbf{a}_t - \mathbf{a}) \\ &\leq \frac{\|\mathbf{a}_1 - \mathbf{a}\|^2 - \|\mathbf{a}_{n+1} - \mathbf{a}\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^n \|\nabla \ell_t(\mathbf{a}_t)\|^2\end{aligned}$$



1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization
  - ▶ Bregman divergence
  - ▶ Regularized minimization equivalent to minimizing latest loss and divergence from previous decision
  - ▶ Constrained minimization equivalent to unconstrained plus Bregman projection
  - ▶ Linearization
  - ▶ Mirror descent
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Seeing the future
  - ▶ Strong convexity
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions

## Online Convex Optimization: A Regularization Viewpoint

- ▶ Suppose  $l_t$  is linear:  $l_t(\mathbf{a}) = \mathbf{g}_t \cdot \mathbf{a}$ .
- ▶ Suppose  $\mathcal{A} = \mathbb{R}^d$ .
- ▶ Then minimizing the regularized criterion

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

corresponds to the gradient step

$$\mathbf{a}_{t+1} = \mathbf{a}_t - \eta \nabla l_t(\mathbf{a}_t).$$

## 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.

## Online Convex Optimization: Regularization

### Regularized minimization

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

- ▶  $R$  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. We'll investigate why regularized minimization can be viewed this way.

## 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  $b$ .

## Bregman Divergence

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

### 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

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

### 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 range of  $\Phi$  is  $\mathcal{A} \subset \mathbb{R}^d$ , in addition to differentiability and strict convexity, we make two more assumptions:

- ▶ The interior of  $\mathcal{A}$  is convex,
- ▶ For a sequence approaching the boundary of  $\mathcal{A}$ ,  
 $\|\nabla\Phi(\mathbf{a}_n)\| \rightarrow \infty$ .

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. *Bregman projection*,  $\Pi_{\mathcal{A}}^\Phi(b) = \arg \min_{a \in \mathcal{A}} D_\Phi(a, b)$  is uniquely defined for closed, convex  $\mathcal{A}$ .
4. *Generalized Pythagoras*: for closed, convex  $\mathcal{A}$ ,  $b^* = \Pi_{\mathcal{A}}^\Phi(b)$ , and  $a \in \mathcal{A}$ ,

$$D_\Phi(a, b) \geq D_\Phi(a, a^*) + D_\Phi(a^*, b).$$

5.  $\nabla_a D_\Phi(a, b) = \nabla \Phi(a) - \nabla \Phi(b)$ .
6. For  $\ell$  linear,  $D_{\Phi+\ell} = D_\Phi$ .
7. For  $\Phi^*$  the Legendre dual of  $\Phi$ ,

$$\begin{aligned}\nabla \Phi^* &= (\nabla \Phi)^{-1}, \\ D_\Phi(a, b) &= D_{\Phi^*}(\nabla \phi(b), \nabla \phi(a)).\end{aligned}$$

## Legendre Dual

For a Legendre function  $\Phi : \mathcal{A} \rightarrow \mathbb{R}$ , the Legendre dual is

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

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

## Legendre Dual

### Example

For  $\Phi = \frac{1}{2} \| \cdot \|_p^2$ , the Legendre dual is  $\Phi^* = \frac{1}{2} \| \cdot \|_q^2$ , where  $1/p + 1/q = 1$ .

### Example

For  $\Phi(a) = \sum_{i=1}^d e^{a_i}$ ,

$$\nabla \Phi(a) = (e^{a_1}, \dots, e^{a_d})'$$

so

$$(\nabla \Phi)^{-1}(u) = \nabla \Phi^*(u) = (\ln u_1, \dots, \ln u_d)'$$

and  $\Phi^*(u) = \sum_i u_i (\ln u_i - 1)$ .

# Online Convex Optimization

1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization
  - ▶ Bregman divergence
  - ▶ Regularized minimization equivalent to minimizing latest loss and divergence from previous decision
  - ▶ Constrained minimization equivalent to unconstrained plus Bregman projection
  - ▶ Linearization
  - ▶ Mirror descent
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Seeing the future
  - ▶ Strong convexity
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions

## 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

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

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:

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

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

1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization
  - ▶ Bregman divergence
  - ▶ Regularized minimization equivalent and Bregman divergence from previous
  - ▶ Constrained minimization equivalent to unconstrained plus Bregman projection
  - ▶ **Linearization**
  - ▶ Mirror descent
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Seeing the future
  - ▶ Strong convexity
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions

## Properties of Regularization Methods: Linear Loss

We can replace  $\ell_t$  by  $\nabla \ell_t(\mathbf{a}_t)$ , and this leads to an upper bound on regret.

### Theorem

*Any strategy for online linear optimization, with regret satisfying*

$$\sum_{t=1}^n g_t \cdot \mathbf{a}_t - \min_{\mathbf{a} \in \mathcal{A}} \sum_{t=1}^n g_t \cdot \mathbf{a} \leq C_n(g_1, \dots, g_n)$$

*can be used to construct a strategy for online convex optimization, with regret*

$$\sum_{t=1}^n \ell_t(\mathbf{a}_t) - \min_{\mathbf{a} \in \mathcal{A}} \sum_{t=1}^n \ell_t(\mathbf{a}) \leq C_n(\nabla \ell_1(\mathbf{a}_1), \dots, \nabla \ell_n(\mathbf{a}_n)).$$

### Proof.

Convexity implies  $\ell_t(\mathbf{a}_t) - \ell_t(\mathbf{a}) \leq \nabla \ell_t(\mathbf{a}_t) \cdot (\mathbf{a}_t - \mathbf{a})$ .

## Properties of Regularization Methods: Linear Loss

### Key Point:

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

Thus, 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

$$\begin{aligned}\nabla R(\tilde{\mathbf{a}}_{t+1}) &= -\eta \sum_{s=1}^t \mathbf{g}_s, \\ \nabla R(\tilde{\mathbf{a}}_t) &= -\eta \sum_{s=1}^{t-1} \mathbf{g}_s,\end{aligned}$$

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

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



# Online Convex Optimization

1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization and Bregman divergences
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Seeing the future
  - ▶ Strong convexity
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions



## Online Convex Optimization: Regularization

Regularized minimization

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

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

## Regularization Methods: Regret

### Theorem

For  $\mathcal{A} = \mathbb{R}^d$ , regularized minimization suffers regret against any  $a \in \mathcal{A}$  of

$$\sum_{t=1}^n \ell_t(\mathbf{a}_t) - \sum_{t=1}^n \ell_t(\mathbf{a}) = \frac{D_R(\mathbf{a}, \mathbf{a}_1) - D_{\Phi_n}(\mathbf{a}, \mathbf{a}_{n+1})}{\eta} + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(\mathbf{a}_t, \mathbf{a}_{t+1}),$$

and thus

$$\hat{L}_n \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(\mathbf{a}) + \frac{D_R(\mathbf{a}, \mathbf{a}_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(\mathbf{a}_t, \mathbf{a}_{t+1}).$$

So the sizes of the steps  $D_{\Phi_t}(\mathbf{a}_t, \mathbf{a}_{t+1})$  determine the regret bound.

## Regularization Methods: Regret

### Theorem

For  $\mathcal{A} = \mathbb{R}^d$ , regularized minimization suffers regret

$$\hat{L}_n \leq \inf_{\mathbf{a} \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(\mathbf{a}) + \frac{D_R(\mathbf{a}, \mathbf{a}_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(\mathbf{a}_t, \mathbf{a}_{t+1}).$$

Notice that we can write

$$\begin{aligned} D_{\Phi_t}(\mathbf{a}_t, \mathbf{a}_{t+1}) &= D_{\Phi_t^*}(\nabla \Phi_t(\mathbf{a}_{t+1}), \nabla \Phi_t(\mathbf{a}_t)) \\ &= D_{\Phi_t^*}(0, \nabla \Phi_{t-1}(\mathbf{a}_t) + \eta \nabla \ell_t(\mathbf{a}_t)) \\ &= D_{\Phi_t^*}(0, \eta \nabla \ell_t(\mathbf{a}_t)). \end{aligned}$$

So it is the size of the gradient steps,  $D_{\Phi_t^*}(0, \eta \nabla \ell_t(\mathbf{a}_t))$ , that determines the regret.

## Regularization Methods: Regret Bounds

### Example

Suppose  $R = \frac{1}{2} \|\cdot\|^2$ . Then we have

$$\hat{L}_n \leq L_n^* + \frac{\|a^* - a_1\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^n \|g_t\|^2.$$

And if  $\|g_t\| \leq G$  and  $\|a^* - a_1\| \leq D$ , choosing  $\eta$  appropriately gives  $\hat{L}_n \leq L_n^* \leq DG\sqrt{n}$ .

# Online Convex Optimization

1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization and Bregman divergences
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Seeing the future
  - ▶ Strong convexity
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions

## Regularization Methods: Regret Bounds

Seeing the future gives small regret:

Theorem

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

$$\sum_{t=1}^n \ell_t(\mathbf{a}_{t+1}) - \sum_{t=1}^n \ell_t(\mathbf{a}) \leq \frac{1}{\eta} (R(\mathbf{a}) - R(\mathbf{a}_1)).$$

# Regularization Methods: Regret Bounds

Proof.

Since  $\mathbf{a}_{t+1}$  minimizes  $\Phi_t$ ,

$$\begin{aligned}\eta \sum_{s=1}^t \ell_s(\mathbf{a}) + R(\mathbf{a}) &\geq \eta \sum_{s=1}^t \ell_s(\mathbf{a}_{t+1}) + R(\mathbf{a}_{t+1}) \\ &= \eta \ell_t(\mathbf{a}_{t+1}) + \eta \sum_{s=1}^{t-1} \ell_s(\mathbf{a}_{t+1}) + R(\mathbf{a}_{t+1}) \\ &\geq \eta \ell_t(\mathbf{a}_{t+1}) + \eta \sum_{s=1}^{t-1} \ell_s(\mathbf{a}_t) + R(\mathbf{a}_t) \\ &\vdots \\ &\geq \eta \sum_{s=1}^t \ell_s(\mathbf{a}_{s+1}) + R(\mathbf{a}_1).\end{aligned}$$

## Regularization Methods: Regret Bounds

### Theorem

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

$$\sum_{t=1}^n \ell_t(\mathbf{a}_{t+1}) - \sum_{t=1}^n \ell_t(\mathbf{a}) \leq \frac{1}{\eta} (R(\mathbf{a}) - R(\mathbf{a}_1)).$$

Thus, if  $\mathbf{a}_t$  and  $\mathbf{a}_{t+1}$  are close, then regret is small:

### Corollary

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

$$\sum_{t=1}^n (\ell_t(\mathbf{a}_t) - \ell_t(\mathbf{a})) \leq \sum_{t=1}^n (\ell_t(\mathbf{a}_t) - \ell_t(\mathbf{a}_{t+1})) + \frac{1}{\eta} (R(\mathbf{a}) - R(\mathbf{a}_1)).$$

So how can we control the increments  $\ell_t(\mathbf{a}_t) - \ell_t(\mathbf{a}_{t+1})$ ?



1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization
  - ▶ Bregman divergence
  - ▶ Regularized minimization equivalent and Bregman divergence from previous
  - ▶ Constrained minimization equivalent to unconstrained plus Bregman projection
  - ▶ Linearization
  - ▶ Mirror descent
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Seeing the future
  - ▶ **Strong convexity**
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions

# Regularization Methods: Regret Bounds

## Definition

We say  $R$  is strongly convex wrt a norm  $\|\cdot\|$  if, for all  $a, b$ ,

$$R(a) \geq R(b) + \nabla R(b) \cdot (a - b) + \frac{1}{2} \|a - b\|^2.$$

For linear losses and strongly convex regularizers, the dual norm of the gradient is small:

## Theorem

If  $R$  is strongly convex wrt a norm  $\|\cdot\|$ , and  $\ell_t(a) = g_t \cdot a$ , then

$$\|a_t - a_{t+1}\| \leq \eta \|g_t\|_*,$$

where  $\|\cdot\|_*$  is the dual norm to  $\|\cdot\|$ :

$$\|v\|_* = \sup\{|v \cdot a| : a \in \mathcal{A}, \|a\| \leq 1\}.$$

## Regularization Methods: Regret Bounds

Proof.

$$R(a_t) \geq R(a_{t+1}) + \nabla R(a_{t+1}) \cdot (a_t - a_{t+1}) + \frac{1}{2} \|a_t - a_{t+1}\|^2,$$

$$R(a_{t+1}) \geq R(a_t) + \nabla R(a_t) \cdot (a_{t+1} - a_t) + \frac{1}{2} \|a_t - a_{t+1}\|^2.$$

Combining,

$$\|a_t - a_{t+1}\|^2 \leq (\nabla R(a_t) - \nabla R(a_{t+1})) \cdot (a_t - a_{t+1})$$

Hence,

$$\|a_t - a_{t+1}\| \leq \|\nabla R(a_t) - \nabla R(a_{t+1})\|_* = \|\eta g_t\|_*.$$



## Regularization Methods: Regret Bounds

This leads to the regret bound:

### Corollary

*For linear losses, if  $R$  is strongly convex wrt  $\|\cdot\|$ , then for all  $a \in \mathcal{A}$ ,*

$$\sum_{t=1}^n (\ell_t(\mathbf{a}_t) - \ell_t(\mathbf{a})) \leq \eta \sum_{t=1}^n \|\mathbf{g}_t\|_*^2 + \frac{1}{\eta} (R(\mathbf{a}) - R(\mathbf{a}_1)).$$

Thus, for  $\|\mathbf{g}_t\|_* \leq G$  and  $R(\mathbf{a}) - R(\mathbf{a}_1) \leq D^2$ , choosing  $\eta$  appropriately gives regret no more than  $2GD\sqrt{n}$ .

## Regularization Methods: Regret Bounds

### Example

Consider  $R(a) = \frac{1}{2}\|a\|^2$ ,  $a_1 = 0$ , and  $\mathcal{A}$  contained in a Euclidean ball of diameter  $D$ .

Then  $R$  is strongly convex wrt  $\|\cdot\|$  and  $\|\cdot\|_* = \|\cdot\|$ . And the mapping between primal and dual spaces is the identity.

So if  $\sup_{a \in \mathcal{A}} \|\nabla \ell_t(a)\| \leq G$ , then regret is no more than  $2GD\sqrt{n}$ .

## Regularization Methods: Regret Bounds

### Example

Consider  $\mathcal{A} = \Delta^m$ ,  $R(a) = \sum_i a_i \ln a_i$ . Then the mapping between primal and dual spaces is  $\nabla R(a) = \ln(a)$  (component-wise). And the divergence is the KL divergence,

$$D_R(a, b) = \sum_i a_i \ln(a_i/b_i).$$

And  $R$  is strongly convex wrt  $\|\cdot\|_1$  (check!).

Suppose that  $\|g_t\|_\infty \leq 1$ . Also,  $R(a) - R(a_1) \leq \ln m$ , so the regret is no more than  $2\sqrt{n \ln m}$ .

## Regularization Methods: Regret Bounds

### Example

$\mathcal{A} = \Delta^m$ ,  $R(a) = \sum_j a_j \ln a_j$ .

What are the updates?

$$\begin{aligned} a_{t+1} &= \Pi_{\mathcal{A}}^R(\tilde{a}_{t+1}) \\ &= \Pi_{\mathcal{A}}^R(\nabla R^*(\nabla R(\tilde{a}_t) - \eta g_t)) \\ &= \Pi_{\mathcal{A}}^R(\nabla R^*(\ln(\tilde{a}_t \exp(-\eta g_t)))) \\ &= \Pi_{\mathcal{A}}^R(\tilde{a}_t \exp(-\eta g_t)), \end{aligned}$$

where the  $\ln$  and  $\exp$  functions are applied component-wise. This is **exponentiated gradient**: mirror descent with  $\nabla R = \ln$ . It is easy to check that the projection corresponds to normalization,  $\Pi_{\mathcal{A}}^R(\tilde{a}) = \tilde{a} / \|\tilde{a}\|_1$ .

## Regularization Methods: Regret Bounds

Notice that when the losses are linear, exponentiated gradient is exactly the **exponential weights strategy** we discussed for a finite comparison class.

Compare  $R(a) = \sum_i a_i \ln a_i$  with  $R(a) = \frac{1}{2} \|a\|^2$ ,  
for  $\|g_t\|_\infty \leq 1$ ,  $\mathcal{A} = \Delta^m$ :

$O(\sqrt{n \ln m})$  versus  $O(\sqrt{mn})$ .



# Online Convex Optimization

1. Problem formulation
2. Empirical minimization fails.
3. Gradient algorithm.
4. Regularized minimization
  - ▶ Bregman divergence
  - ▶ Regularized minimization equivalent and Bregman divergence from previous
  - ▶ Constrained minimization equivalent to unconstrained plus Bregman projection
  - ▶ Linearization
  - ▶ Mirror descent
5. Regret bounds
  - ▶ Unconstrained minimization
  - ▶ Strong convexity
  - ▶ Examples (gradient, exponentiated gradient)
  - ▶ Extensions

## Regularization Methods: Extensions

- ▶ Instead of

$$a_{t+1} = \arg \min_{a \in \mathcal{A}} (\eta \ell_t(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t)) ,$$

we can use

$$a_{t+1} = \arg \min_{a \in \mathcal{A}} (\eta \ell_t(a) + D_{\Phi_{t-1}}(a, a_t)) .$$

And analogous results apply. For instance, this is the approach used by the first gradient method we considered.

- ▶ We can get faster rates with stronger assumptions on the losses...

## Theorem

Define

$$\mathbf{a}_{t+1} = \arg \min_{\mathbf{a} \in \mathbb{R}^d} \left( \sum_{t=1}^n \eta_t \ell_t(\mathbf{a}) + R(\mathbf{a}) \right).$$

For any  $\mathbf{a} \in \mathbb{R}^d$ ,

$$\hat{L}_n - \sum_{t=1}^n \ell_t(\mathbf{a}) \leq \sum_{t=1}^n \frac{1}{\eta_t} (D_{\Phi_t}(\mathbf{a}_t, \mathbf{a}_{t+1}) + D_{\Phi_{t-1}}(\mathbf{a}, \mathbf{a}_t) - D_{\Phi_t}(\mathbf{a}, \mathbf{a}_{t+1})).$$

If we linearize the  $\ell_t$ , we have

$$\hat{L}_n - \sum_{t=1}^n \ell_t(\mathbf{a}) \leq \sum_{t=1}^n \frac{1}{\eta_t} (D_R(\mathbf{a}_t, \mathbf{a}_{t+1}) + D_R(\mathbf{a}, \mathbf{a}_t) - D_R(\mathbf{a}, \mathbf{a}_{t+1})).$$

But what if  $\ell_t$  are strongly convex?

## Regularization Methods: Strongly Convex Losses

### Theorem

If  $\ell_t$  is  $\sigma$ -strongly convex wrt  $R$ , that is, for all  $a, b \in \mathbb{R}^d$ ,

$$\ell_t(a) \geq \ell_t(b) + \nabla \ell_t(b) \cdot (a - b) + \frac{\sigma}{2} D_R(a, b),$$

then for any  $a \in \mathbb{R}^d$ , this strategy with  $\eta_t = \frac{2}{t\sigma}$  has regret

$$\hat{L}_n - \sum_{t=1}^n \ell_t(a) \leq \sum_{t=1}^n \frac{1}{\eta_t} D_R(a_t, a_{t+1}).$$

## Strongly Convex Losses: Proof idea

$$\begin{aligned} & \sum_{t=1}^n (\ell_t(\mathbf{a}_t) - \ell_t(\mathbf{a})) \\ & \leq \sum_{t=1}^n \left( \nabla \ell_t(\mathbf{a}_t) \cdot (\mathbf{a}_t - \mathbf{a}) - \frac{\sigma}{2} D_R(\mathbf{a}, \mathbf{a}_t) \right) \\ & \leq \sum_{t=1}^n \frac{1}{\eta_t} \left( D_R(\mathbf{a}_t, \mathbf{a}_{t+1}) + D_R(\mathbf{a}, \mathbf{a}_t) - D_R(\mathbf{a}, \mathbf{a}_{t+1}) - \frac{\eta_t \sigma}{2} D_R(\mathbf{a}, \mathbf{a}_t) \right) \\ & \leq \sum_{t=1}^n \frac{1}{\eta_t} D_R(\mathbf{a}_t, \mathbf{a}_{t+1}) + \sum_{t=2}^n \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - \frac{\sigma}{2} \right) D_R(\mathbf{a}, \mathbf{a}_t) \\ & \quad + \left( \frac{1}{\eta_1} - \frac{\sigma}{2} \right) D_R(\mathbf{a}, \mathbf{a}_1). \end{aligned}$$

And choosing  $\eta_t$  appropriately eliminates the second and third terms.

## Strongly Convex Losses

### Example

For  $R(a) = \frac{1}{2}\|a\|^2$ , we have

$$\hat{L}_n - L_n^* \leq \frac{1}{2} \sum_{t=1}^n \frac{1}{\eta_t} \|\eta_t \nabla \ell_t\|^2 \leq \sum_{t=1}^n \frac{G^2}{t\sigma} = O\left(\frac{G^2}{\sigma} \log n\right).$$

## Strongly Convex Losses

**Key Point:** When the loss is strongly convex wrt the regularizer, the regret rate can be faster; in the case of quadratic  $R$  (and  $\ell_t$ ), it is  $O(\log n)$ , versus  $O(\sqrt{n})$ .

## Course Synopsis

- ▶ A finite comparison class:  $\mathcal{A} = \{1, \dots, m\}$ .
- ▶ Converting online to batch.
- ▶ Online convex optimization.
- ▶ Log loss.