Topics in Prediction and Learning Lectures 2 and 3: Online Convex Optimization

Peter Bartlett

Computer Science and Statistics University of California at Berkeley

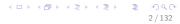
Mathematical Sciences Queensland University of Technology

27 February–9 March, 2017 CREST, ENSAE

イロト 不得下 イヨト イヨト 二日

### A repeated game:

At round t:



# A repeated game:

At round *t*:

**1** Player chooses prediction  $a_t \in A$ .

# A repeated game:

At round *t*:

- **1** Player chooses prediction  $a_t \in A$ .
- 2 Adversary chooses loss  $\ell_t \in \mathcal{L}$ .

・ロト ・回ト ・ヨト ・ヨト

3

2/132

# A repeated game:

At round *t*:

- **1** Player chooses prediction  $a_t \in A$ .
- 2 Adversary chooses loss  $\ell_t \in \mathcal{L}$ .
- 3 Player incurs loss  $\ell_t(a_t)$ .

### A repeated game:

At round *t*:

- **1** Player chooses prediction  $a_t \in A$ .
- 2 Adversary chooses loss  $\ell_t \in \mathcal{L}$ .
- 3 Player incurs loss  $\ell_t(a_t)$ .

# Player's aim:

Minimize regret:

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

# A repeated game:

At round *t*:

- **1** Player chooses prediction  $a_t \in A$ .
- 2 Adversary chooses loss  $\ell_t \in \mathcal{L}$ .
- 3 Player incurs loss  $\ell_t(a_t)$ .

# Player's aim:

Minimize *regret* wrt comparison C:

$$\mathcal{R}_n := \sum_{t=1}^n \ell_t(\mathfrak{a}_t) - \inf_{\hat{\mathfrak{a}} \in \mathcal{C}} \sum_{t=1}^n \ell_t(\hat{\mathfrak{a}}_t).$$

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

# Examples

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

#### Examples

• Quadratic loss:  $\ell_t(a) = ||x_t - a||^2$ .

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

### Examples

- Quadratic loss:  $\ell_t(a) = ||x_t a||^2$ .
- Linear regression:  $\ell_t(a) = (x_t \cdot a y_t)^2$ .

イロト 不得下 イヨト イヨト 二日

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

### Examples

- Quadratic loss:  $\ell_t(a) = ||x_t a||^2$ .
- Linear regression:  $\ell_t(a) = (x_t \cdot a y_t)^2$ .
- Absolute loss linear regression:  $\ell_t(a) = |x_t \cdot a y_t|$ .

イロト 不得下 イヨト イヨト 二日

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

### Examples

- Quadratic loss:  $\ell_t(a) = ||x_t a||^2$ .
- Linear regression:  $\ell_t(a) = (x_t \cdot a y_t)^2$ .
- Absolute loss linear regression:  $\ell_t(a) = |x_t \cdot a y_t|$ .
- Prediction with expert advice:  $\ell_t(a) = w_t^\top a$   $(\mathcal{A} = \Delta^m)$ .

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

#### Examples

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

#### Examples

• Shortest path:  $\ell_t(a) = w_t^\top a$ 

 $(\mathcal{A} = \text{flow}, w_t = \text{edge weights}).$ 

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

#### Examples

• Shortest path:  $\ell_t(a) = w_t^\top a$ 

$$A =$$
flow,  $w_t =$ edge weights).

 $(\mathcal{A} = \Delta^m).$ 

• Portfolio optimization:  $\ell_t(a) = -\log(r_t^{\top}a)$ 

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

#### Examples

- Shortest path:  $\ell_t(a) = w_t^\top a$  ( $\mathcal{A} = \text{flow}, w_t = \text{edge weights}$ ).
- Portfolio optimization:  $\ell_t(a) = -\log(r_t^\top a)$   $(\mathcal{A} = \Delta^m).$
- Collaborative filtering:  $\ell_t(A) = (x_t A_{i_t,j_t})^2$ .

 $(\mathcal{A}=\mathbb{R}^{m\times n}).$ 

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

#### Examples

- Shortest path:  $\ell_t(a) = w_t^\top a$  ( $\mathcal{A} = \text{flow}, w_t = \text{edge weights}$ ).
- Portfolio optimization:  $\ell_t(a) = -\log(r_t^\top a)$   $(\mathcal{A} = \Delta^m).$
- Collaborative filtering:  $\ell_t(A) = (x_t A_{i_t,j_t})^2$ .
- SVM:  $\ell_t(A) = (1 y_t x_t^\top a)_+ + \lambda ||a||^2$ .  $(\mathcal{A} = \mathsf{RKHS})$ .

 $(\mathcal{A}=\mathbb{R}^{m\times n}).$ 

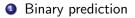
- $\mathcal{A} = \text{ convex subset of } \mathbb{R}^d$ .
- $\mathcal{L} =$  set of convex real functions on  $\mathcal{A}$ .

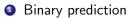
#### Examples

- Shortest path:  $\ell_t(a) = w_t^\top a$  ( $\mathcal{A} = \text{flow}, w_t = \text{edge weights}$ ).
- Portfolio optimization:  $\ell_t(a) = -\log(r_t^\top a)$   $(\mathcal{A} = \Delta^m).$
- Collaborative filtering:  $\ell_t(A) = (x_t A_{i_t,j_t})^2$ .
- SVM:  $\ell_t(A) = (1 y_t x_t^\top a)_+ + \lambda \|a\|^2$ .  $(\mathcal{A} = \mathsf{RKHS})$ .
- Density estimation: ℓ<sub>t</sub>(a) = − log (exp(a'T(y<sub>t</sub>) − A(a))), for exponential family with sufficient statistic T(y).

 $(\mathcal{A}=\mathbb{R}^{m\times n}).$ 

◆□ → < □ → < Ξ → < Ξ → < Ξ → < Ξ → ○ Q ↔ 5/132</p>





② General online convex



② General online convex

5/132

Minimax strategies

• With (perfect) expert advice

- ② General online convex
- Minimax strategies

• With (perfect) expert advice

- Minimax strategy
- ② General online convex
- Minimax strategies

- With (perfect) expert advice
- Minimax strategy
- With imperfect experts: exponential weights

イロト 不得 とくき とくき とうき

- General online convex
- Minimax strategies

• Suppose we are predicting whether it will rain tomorrow.

- Suppose we are predicting whether it will rain tomorrow.
- We have access to a set of *m* experts, who each make a forecast.

- Suppose we are predicting whether it will rain tomorrow.
- We have access to a set of *m* experts, who each make a forecast.
- Can we ensure that we predict almost as well as the best expert?

- Suppose we are predicting whether it will rain tomorrow.
- We have access to a set of *m* experts, who each make a forecast.
- Can we ensure that we predict almost as well as the best expert?
- We'll consider two settings: voting and prediction.

The player votes for a mixture of experts:

The player votes for a mixture of experts: we set  $\mathcal{A} = \Delta^m$ , the probability simplex on  $\{1, \ldots, m\}$ , and the loss function at time t is  $\ell_t(a) = |a^\top f_t - y_t|$ , where  $f_t \in \{0, 1\}^m$  are the forecasts of the experts and  $y_t \in \{0, 1\}$  is the outcome.

The player votes for a mixture of experts: we set  $\mathcal{A} = \Delta^m$ , the probability simplex on  $\{1, \ldots, m\}$ , and the loss function at time t is  $\ell_t(a) = |a^\top f_t - y_t|$ , where  $f_t \in \{0, 1\}^m$  are the forecasts of the experts and  $y_t \in \{0, 1\}$  is the outcome.

#### Prediction

The player votes for a mixture of experts, but the vote can depend on their forecasts:

The player votes for a mixture of experts: we set  $\mathcal{A} = \Delta^m$ , the probability simplex on  $\{1, \ldots, m\}$ , and the loss function at time t is  $\ell_t(a) = |a^\top f_t - y_t|$ , where  $f_t \in \{0, 1\}^m$  are the forecasts of the experts and  $y_t \in \{0, 1\}$  is the outcome.

#### Prediction

The player votes for a mixture of experts, but the vote can depend on their forecasts: we set  $\mathcal{A} = (\Delta^m)^{\{0,1\}^m}$ , and the loss function at time *t* is  $\ell_t(a) = |a(f_t)^\top f_t - y_t|$ .

The player votes for a mixture of experts: we set  $\mathcal{A} = \Delta^m$ , the probability simplex on  $\{1, \ldots, m\}$ , and the loss function at time t is  $\ell_t(a) = |a^\top f_t - y_t|$ , where  $f_t \in \{0, 1\}^m$  are the forecasts of the experts and  $y_t \in \{0, 1\}$  is the outcome.

#### Prediction

The player votes for a mixture of experts, but the vote can depend on their forecasts: we set  $\mathcal{A} = (\Delta^m)^{\{0,1\}^m}$ , and the loss function at time t is  $\ell_t(a) = |a(f_t)^\top f_t - y_t|$ . The comparison class  $\mathcal{C}$  is the set of constant functions. (That is,  $a \in \mathcal{C}$  has  $p \in \Delta^m$  so that for all  $f \in \{0,1\}^m$ , a(f) = p.) Prediction allows the player to see how the experts' predictions compare before making a prediction.

Prediction allows the player to see how the experts' predictions compare before making a prediction.

We write  $\ell_t(e_i) \in \{0, 1\}$  for the loss incurred by expert *i*, where  $e_i \in \Delta^m$  is zero in all but the *i*th coordinate. and  $\ell_t(e_i) \in \{0, 1\}$  is the indicator for expert *i* making an incorrect forecast at time *t*.

We can interpret any  $a \in \Delta^m$  equivalently as a prediction,

 $\hat{y}_t = a^{\top} f_t \in [0, 1]$ . And we can view  $\hat{y}_t$  either as the expectation of a random  $\{0, 1\}$ -valued prediction where the loss  $\ell_t(a_t)$  is the probability of a mistake, or as a real-valued prediction, where the loss is the absolute difference between the prediction and the outcome.

The minimax regret is the value of the game:

$$\min_{a_1} \max_{\ell_1} \cdots \min_{a_n} \max_{\ell_n} \left( \sum_{t=1}^n \ell_t(a_t) - \min_{a \in \mathcal{C}} \sum_{t=1}^n \ell_t(a) \right).$$

#### An easier game

Suppose that the adversary is constrained to choose the sequence  $\ell_t$  so that some expert incurs no loss, that is,

$$\min_{a\in\mathcal{C}}\sum_{t=1}^n\ell_t(a)=0.$$

イロト 不得下 イヨト イヨト 二日

9/132

How should we predict?

# Halving Algorithm

[Littlestone, 1988]

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

# Halving Algorithm

• Define the set of experts who have been correct so far:

$$C_t = \{i : \ell_1(e_i) = \cdots = \ell_{t-1}(e_i) = 0\}.$$

[Littlestone, 1988]

10/132

・ロン ・回と ・ヨン ・ヨン

# Halving Algorithm

• Define the set of experts who have been correct so far:

$$C_t = \{i : \ell_1(e_i) = \cdots = \ell_{t-1}(e_i) = 0\}.$$

#### Choose

$$\hat{y}_t = a_t(f_t)^{\top} f_t = \text{majority}\left(\{f_t(j) : j \in C_t\}\right)$$

[Littlestone, 1988]

10/132

・ロン ・回と ・ヨン ・ヨン

# Halving Algorithm

• Define the set of experts who have been correct so far:

$$C_t = \{i : \ell_1(e_i) = \cdots = \ell_{t-1}(e_i) = 0\}.$$

#### Choose

$$\hat{y}_t = a_t(f_t)^{\top} f_t = \text{majority}\left(\{f_t(j) : j \in C_t\}\right)$$

#### Theorem

This strategy has regret no more than  $\log_2 m$ .

[Littlestone, 1988]

イロト 不同下 イヨト イヨト

If the strategy makes a mistake (that is,  $\ell_t(a_t) = 1$ ), then the minority of  $\{f_t(j) : j \in C_t\}$  is correct, so at least half of the experts are eliminated:

 $|C_{t+1}| \leq \frac{|C_t|}{2}.$ 

If the strategy makes a mistake (that is,  $\ell_t(a_t) = 1$ ), then the minority of  $\{f_t(j) : j \in C_t\}$  is correct, so at least half of the experts are eliminated:

$$|C_{t+1}| \leq \frac{|C_t|}{2}.$$

And otherwise  $|C_{t+1}| \leq |C_t|$  (because  $|C_t|$  never increases).

If the strategy makes a mistake (that is,  $\ell_t(a_t) = 1$ ), then the minority of  $\{f_t(j) : j \in C_t\}$  is correct, so at least half of the experts are eliminated:

$$|C_{t+1}| \leq \frac{|C_t|}{2}.$$

And otherwise  $|C_{t+1}| \le |C_t|$  (because  $|C_t|$  never increases). Thus,

$$\sum_{t=1}^{n} \ell_t(a_t) \le \log_2 \frac{|C_1|}{|C_{n+1}|}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

If the strategy makes a mistake (that is,  $\ell_t(a_t) = 1$ ), then the minority of  $\{f_t(j) : j \in C_t\}$  is correct, so at least half of the experts are eliminated:

$$|C_{t+1}| \leq \frac{|C_t|}{2}.$$

And otherwise  $|C_{t+1}| \le |C_t|$  (because  $|C_t|$  never increases). Thus,

$$\sum_{t=1}^{n} \ell_t(a_t) \le \log_2 \frac{|\mathcal{C}_1|}{|\mathcal{C}_{n+1}|} = \log_2 m - \log_2 |\mathcal{C}_{n+1}|$$

・ロ ・ ・ 日 ・ ・ 目 ・ ・ 目 ・ うへで 11/132

If the strategy makes a mistake (that is,  $\ell_t(a_t) = 1$ ), then the minority of  $\{f_t(j) : j \in C_t\}$  is correct, so at least half of the experts are eliminated:

$$|C_{t+1}| \leq \frac{|C_t|}{2}.$$

And otherwise  $|C_{t+1}| \le |C_t|$  (because  $|C_t|$  never increases). Thus,

$$\sum_{t=1}^{n} \ell_t(a_t) \leq \log_2 \frac{|C_1|}{|C_{n+1}|} = \log_2 m - \log_2 |C_{n+1}| \leq \log_2 m.$$

・ロ ・ ・ 日 ・ ・ 目 ・ 日 ・ 日 ・ 日 ・ 11/132

We can do better with a randomized voting strategy.

 $[ \text{Karlin and Peres, 2016} ] \\ < \square \mathrel{\blacktriangleright} < \square \mathrel{\blacktriangleright} < \square \mathrel{\leftarrow} < \square \mathrel{\leftarrow} < \square \mathrel{\leftarrow} < \square \\ 12/132 \\ \end{bmatrix}$ 

We can do better with a randomized voting strategy.

# Random Leader

Choose  $a_t(f_t)$  uniformly on

$$C_t = \{i : \ell_1(e_i) = \cdots = \ell_{t-1}(e_i) = 0\}.$$

We can do better with a randomized voting strategy.

### Random Leader

Choose  $a_t(f_t)$  uniformly on

$$C_t = \{i : \ell_1(e_i) = \cdots = \ell_{t-1}(e_i) = 0\}.$$

#### Theorem

This strategy has regret no more than  $H_m - 1$ , where

$$H_m = \sum_{i=1}^m \frac{1}{i} \in (\ln m, \ln m + 1).$$

We show that, at time t, the strategy can make no more than  $H_{|C_t|} - 1$  mistakes from that time on.

We show that, at time t, the strategy can make no more than  $H_{|C_t|} - 1$  mistakes from that time on.

• This is clearly true when  $|C_t| = 1$ : the strategy never makes another mistake.

We show that, at time t, the strategy can make no more than  $H_{|C_t|} - 1$  mistakes from that time on.

- This is clearly true when  $|C_t| = 1$ : the strategy never makes another mistake.
- Suppose it is true for |C<sub>t+1</sub>| < k, suppose that |C<sub>t</sub>| = k, and suppose that j experts in C<sub>t</sub> make a mistake at time t, where 1 ≤ j ≤ k − 1.

We show that, at time t, the strategy can make no more than  $H_{|C_t|} - 1$  mistakes from that time on.

- This is clearly true when  $|C_t| = 1$ : the strategy never makes another mistake.
- Suppose it is true for  $|C_{t+1}| < k$ , suppose that  $|C_t| = k$ , and suppose that j experts in  $C_t$  make a mistake at time t, where  $1 \le j \le k 1$ . Then the expected number of mistakes made from time t onwards is no more than

$$\frac{J}{k} + H_{k-j} - 1$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

We show that, at time t, the strategy can make no more than  $H_{|C_t|} - 1$  mistakes from that time on.

- This is clearly true when  $|C_t| = 1$ : the strategy never makes another mistake.
- Suppose it is true for  $|C_{t+1}| < k$ , suppose that  $|C_t| = k$ , and suppose that j experts in  $C_t$  make a mistake at time t, where  $1 \le j \le k 1$ . Then the expected number of mistakes made from time t onwards is no more than

$$\frac{J}{k}+H_{k-j}-1\leq H_k-1.$$

◆□ → < □ → < 三 → < 三 → < 三 → < 三 → ○ Q (~ 13/132

# Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

[Karlin and Peres, 2016]

#### Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

[Karlin and Peres, 2016]

Lower bound

Set  $k = \lfloor \log_2 m \rfloor$  so that  $2^k \le m \le 2^{k+1}$ .

#### Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

[Karlin and Peres, 2016]

### Lower bound

Set  $k = \lfloor \log_2 m \rfloor$  so that  $2^k \le m \le 2^{k+1}$ . Consider the following adversary strategy:

#### Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

[Karlin and Peres, 2016]

### Lower bound

Set  $k = \lfloor \log_2 m \rfloor$  so that  $2^k \le m \le 2^{k+1}$ . Consider the following adversary strategy:

• Choose  $C_0$  as the first k experts.

#### Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

[Karlin and Peres, 2016]

## Lower bound

Set  $k = \lfloor \log_2 m \rfloor$  so that  $2^k \le m \le 2^{k+1}$ . Consider the following adversary strategy:

- Choose  $C_0$  as the first k experts.
- At round 1 ≤ t ≤ k, choose C<sub>t+1</sub> ⊂ C<sub>t</sub> uniformly at random from subsets of size |C<sub>t</sub>|/2.

#### Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

[Karlin and Peres, 2016]

## Lower bound

Set  $k = \lfloor \log_2 m \rfloor$  so that  $2^k \le m \le 2^{k+1}$ . Consider the following adversary strategy:

- Choose  $C_0$  as the first k experts.
- At round 1 ≤ t ≤ k, choose C<sub>t+1</sub> ⊂ C<sub>t</sub> uniformly at random from subsets of size |C<sub>t</sub>|/2.
- Choose  $y_t \in \{0, 1\}$  uniformly at random.

#### Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

[Karlin and Peres, 2016]

## Lower bound

Set  $k = \lfloor \log_2 m \rfloor$  so that  $2^k \le m \le 2^{k+1}$ . Consider the following adversary strategy:

- Choose  $C_0$  as the first k experts.
- At round 1 ≤ t ≤ k, choose C<sub>t+1</sub> ⊂ C<sub>t</sub> uniformly at random from subsets of size |C<sub>t</sub>|/2.
- Choose  $y_t \in \{0, 1\}$  uniformly at random.

Set

$$f_t^i = egin{cases} y_t & ext{ for } i \in C_{t+1}, \ 1-y_t & ext{ otherwise}. \end{cases}$$

Clearly, after k rounds there is still a perfect expert.

Clearly, after k rounds there is still a perfect expert. The expected number of mistakes of any player strategy is

 $\frac{k}{2}$ 

Clearly, after k rounds there is still a perfect expert. The expected number of mistakes of any player strategy is

・ロト ・回ト ・ヨト ・ヨト

15/132

$$\frac{k}{2} = \frac{\lfloor \log_2 m \rfloor}{2}$$

Clearly, after k rounds there is still a perfect expert. The expected number of mistakes of any player strategy is

$$\frac{k}{2} = \frac{\lfloor \log_2 m \rfloor}{2} \ge \lfloor \log_4 m \rfloor.$$

Set  $a_t(f_t)^{\top} f_t = \phi(p_t)\hat{y} + (1 - \phi(p_t))(1 - \hat{y})$ , where

Set  $a_t(f_t)^{\top} f_t = \phi(p_t)\hat{y} + (1 - \phi(p_t))(1 - \hat{y})$ , where

 $\hat{y} = \text{majority}\left(\{f_t(j) : j \in C_t\}\right),$ 

Set  $a_t(f_t)^{\top} f_t = \phi(p_t)\hat{y} + (1 - \phi(p_t))(1 - \hat{y})$ , where

$$\hat{y} = \text{majority}\left(\left\{f_t(j) : j \in C_t\right\}\right),\$$

$$p_t = \frac{1}{|C_t|} \left|\left\{i \in C_t : f_t(i) = \text{majority}\left(\left\{f_t(j) : j \in C_t\right\}\right)\right\}\right|$$

Set  $a_t(f_t)^{\top} f_t = \phi(p_t)\hat{y} + (1 - \phi(p_t))(1 - \hat{y})$ , where

$$\begin{split} \hat{y} &= \mathsf{majority}\left(\{f_t(j) : j \in C_t\}\right),\\ p_t &= \frac{1}{|C_t|} \left| \left\{ i \in C_t : f_t(i) = \mathsf{majority}\left(\{f_t(j) : j \in C_t\}\right) \right\} \right|,\\ \phi(p) &= 1 + \log_4 p. \end{split}$$

Set  $a_t(f_t)^{\top} f_t = \phi(p_t)\hat{y} + (1 - \phi(p_t))(1 - \hat{y})$ , where

$$\begin{split} \hat{y} &= \mathsf{majority}\left(\{f_t(j) : j \in C_t\}\right), \\ p_t &= \frac{1}{|C_t|} \left| \left\{ i \in C_t : f_t(i) = \mathsf{majority}\left(\{f_t(j) : j \in C_t\}\right) \right\} \right| \\ \phi(p) &= 1 + \log_4 p. \end{split}$$

16/132

That is, follow the majority with probability  $\phi(p_t)$ .

Set  $a_t(f_t)^{\top} f_t = \phi(p_t)\hat{y} + (1 - \phi(p_t))(1 - \hat{y})$ , where

$$\begin{split} \hat{y} &= \mathsf{majority}\left(\left\{f_t(j) : j \in C_t\right\}\right), \\ p_t &= \frac{1}{|C_t|} \left| \left\{i \in C_t : f_t(i) = \mathsf{majority}\left(\left\{f_t(j) : j \in C_t\right\}\right)\right\} \right| \\ \phi(p) &= 1 + \log_4 p. \end{split}$$

That is, follow the majority with probability  $\phi(p_t)$ .

(NB:  $\phi(p) = 1$  corresponds to the halving algorithm.  $\phi(p) = p$  corresponds to voting uniformly on  $C_t$ .)

We'd like an upper bound on the expected number of mistakes of the form  $\log_a m$ .

We'd like an upper bound on the expected number of mistakes of the form  $\log_a m$ . To make the inductive proof of this bound work, we need to consider two cases.

We'd like an upper bound on the expected number of mistakes of the form  $\log_a m$ . To make the inductive proof of this bound work, we need to consider two cases. First, if the majority is correct  $(y_t = \hat{y})$ , then we need

 $\log_a(p_t m) + (1 - \phi(p_t)) \le \log_a m.$ 

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

We'd like an upper bound on the expected number of mistakes of the form  $\log_a m$ . To make the inductive proof of this bound work, we need to consider two cases. First, if the majority is correct  $(y_t = \hat{y})$ , then we need

 $\log_a(p_t m) + (1 - \phi(p_t)) \le \log_a m.$ 

Second, if the minority is correct, then we need

 $\log_a((1-p_t)m) + \phi(p_t) \le \log_a m.$ 

イロン イロン イヨン イヨン 三日

17 / 132

#### Proof

 $\log_a(p_t m) + (1 - \phi(p_t)) \le \log_a m,$  $\log_a((1 - p_t)m) + \phi(p_t) \le \log_a m.$ 

#### Proof

$$\log_a(p_t m) + (1 - \phi(p_t)) \le \log_a m,$$
  
$$\log_a((1 - p_t)m) + \phi(p_t) \le \log_a m.$$

Rearranging and combining, we need

 $1 + \log_a p_t \le \phi(p_t) \le - \log_a (1 - p_t)$ 

#### Proof

$$\log_a(p_t m) + (1 - \phi(p_t)) \le \log_a m,$$
  
$$\log_a((1 - p_t)m) + \phi(p_t) \le \log_a m.$$

Rearranging and combining, we need

$$egin{aligned} &1+\log_a p_t \leq \phi(p_t) \leq -\log_a(1-p_t)\ &\Leftrightarrow &\log_a(ap_t) \leq \log_a\left(rac{1}{1-p_t}
ight). \end{aligned}$$

・ロ ・ ・ 日 ・ ・ 目 ・ ・ 目 ・ うへで 18/132

#### Proof

$$\log_a(p_t m) + (1 - \phi(p_t)) \le \log_a m,$$
  
$$\log_a((1 - p_t)m) + \phi(p_t) \le \log_a m.$$

Rearranging and combining, we need

$$egin{aligned} &1+\log_a p_t \leq \phi(p_t) \leq -\log_a(1-p_t)\ &\Leftrightarrow &\log_a(ap_t) \leq \log_a\left(rac{1}{1-p_t}
ight). \end{aligned}$$

The largest *a* satisfying  $ap_t(1-p_t) \le 1$  is a = 4.

$$\log_a(p_t m) + (1 - \phi(p_t)) \le \log_a m,$$
  
$$\log_a((1 - p_t)m) + \phi(p_t) \le \log_a m.$$

Rearranging and combining, we need

$$egin{aligned} &1+\log_a p_t \leq \phi(p_t) \leq -\log_a(1-p_t)\ &\Leftrightarrow &\log_a(ap_t) \leq \log_a\left(rac{1}{1-p_t}
ight). \end{aligned}$$

The largest a satisfying  $ap_t(1-p_t) \le 1$  is a = 4. So any  $\phi(p_t)$  between  $1 + \log_4 p_t$  and  $-\log_4(1-p_t)$  will suffice.

#### Theorem

The minimax regret is between  $\lfloor \log_4 m \rfloor$  and  $\log_4 m$ .

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

#### Binary prediction

- With (perfect) expert advice
- Minimax strategy
- With imperfect experts: exponential weights

イロト 不得下 イヨト イヨト 二日

- eneral online convex
- Minimax strategies

We return to the voting setting, and allow even the best expert to make mistakes.

We return to the voting setting, and allow even the best expert to make mistakes.

#### Voting

The player votes for a mixture of experts: we set  $\mathcal{A} = \Delta^m$ , the probability simplex on  $\{1, \ldots, m\}$ , and the loss function at time t is  $\ell_t(a) = \sum_{i=1}^m a_i \ell_t(e_i)$ , where  $e_i \in \Delta^m$  is zero in all but the *i*th coordinate, and  $\ell_t(e_i) \in \{0, 1\}$  is the indicator for the *i*th expert making an incorrect forecast at time t.

#### **Exponential Weights**

[Littlestone and Warmuth, 1994]

◆□ → ◆□ → ◆ 三 → ◆ 三 → ○ へ ○
22/132

#### **Exponential Weights**

• Maintain a set of (unnormalized) weights over experts:

#### Exponential Weights

• Maintain a set of (unnormalized) weights over experts:

$$w_1^i = 1,$$

#### Exponential Weights

• Maintain a set of (unnormalized) weights over experts:

$$w_1^i = 1,$$
  
$$w_{t+1}^i = w_t^i \exp\left(-\eta \ell_t(e_i)\right).$$

#### Exponential Weights

• Maintain a set of (unnormalized) weights over experts:

$$w_1^i = 1,$$
  
$$w_{t+1}^i = w_t^i \exp\left(-\eta \ell_t(e_i)\right).$$

• Here,  $\eta > 0$  is a parameter of the algorithm.

#### Exponential Weights

• Maintain a set of (unnormalized) weights over experts:

$$w_1^i = 1,$$
  
$$w_{t+1}^i = w_t^i \exp\left(-\eta \ell_t(e_i)\right).$$

- Here,  $\eta > 0$  is a parameter of the algorithm.
- Choose a<sub>t</sub> as the normalized vector,

$$a_t = \frac{1}{\sum_{i=1}^m w_t^i} w_t.$$

#### Theorem

The exponential weights strategy with parameter

$$\eta = \sqrt{\frac{8\ln m}{n}}$$

has regret satisfying

$$R_n \leq \sqrt{\frac{n \ln m}{2}}.$$

[Cesa-Bianchi, Freund, Haussler, Helmbold, Schapire, and Warmuth, 1997]

### Proof Idea

We use a measure of progress:

$$W_t = \sum_{i=1}^m w_t^i.$$

#### Proof Idea

We use a measure of progress:

$$\mathcal{N}_t = \sum_{i=1}^m w_t^i.$$

•  $W_n$  grows at least as

$$\exp\left(-\eta\min_{i}\sum_{t=1}^{n}\ell_{t}(e_{i})
ight).$$

#### Proof Idea

We use a measure of progress:

$$N_t = \sum_{i=1}^m w_t^i.$$

• 
$$W_n$$
 grows at least as

$$\exp\left(-\eta\min_{i}\sum_{t=1}^{n}\ell_{t}(e_{i})
ight).$$

2  $W_n$  grows no faster than

$$\exp\left(-\eta\sum_{t=1}^n\ell_t(a_t)\right).$$

$$\ln \frac{W_{n+1}}{W_1}$$

$$\ln \frac{W_{n+1}}{W_1} = \ln \left( \sum_{i=1}^m w_{n+1}^i \right) - \ln m$$

$$\ln \frac{W_{n+1}}{W_1} = \ln \left( \sum_{i=1}^m w_{n+1}^i \right) - \ln m$$
$$= \ln \left( \sum_{i=1}^m \exp \left( -\eta \sum_t \ell_t(e_i) \right) \right) - \ln m$$

$$n \frac{W_{n+1}}{W_1} = \ln\left(\sum_{i=1}^m w_{n+1}^i\right) - \ln m$$
$$= \ln\left(\sum_{i=1}^m \exp\left(-\eta \sum_t \ell_t(e_i)\right)\right) - \ln m$$
$$\geq \ln\left(\max_i \exp\left(-\eta \sum_t \ell_t(e_i)\right)\right) - \ln m$$

$$\ln \frac{W_{n+1}}{W_1} = \ln \left( \sum_{i=1}^m w_{n+1}^i \right) - \ln m$$
$$= \ln \left( \sum_{i=1}^m \exp\left(-\eta \sum_t \ell_t(e_i)\right) \right) - \ln m$$
$$\geq \ln \left( \max_i \exp\left(-\eta \sum_t \ell_t(e_i)\right) \right) - \ln m$$
$$= -\eta \min_i \left( \sum_t \ell_t(e_i) \right) - \ln m$$

$$\ln \frac{W_{n+1}}{W_1} = \ln \left( \sum_{i=1}^m w_{n+1}^i \right) - \ln m$$
$$= \ln \left( \sum_{i=1}^m \exp\left(-\eta \sum_t \ell_t(e_i)\right) \right) - \ln m$$
$$\ge \ln \left( \max_i \exp\left(-\eta \sum_t \ell_t(e_i)\right) \right) - \ln m$$
$$= -\eta \min_i \left( \sum_t \ell_t(e_i) \right) - \ln m$$
$$= -\eta \inf_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) - \ln m.$$

### Proof idea:

$$\ln rac{W_{t+1}}{W_t}$$

26/132

$$\ln \frac{W_{t+1}}{W_t} = \ln \left( \frac{\sum_{i=1}^m \exp(-\eta \ell_t(e_i)) w_t^i}{\sum_i w_t^i} \right)$$

$$\ln \frac{W_{t+1}}{W_t} = \ln \left( \frac{\sum_{i=1}^m \exp(-\eta \ell_t(e_i)) w_t^i}{\sum_i w_t^i} \right)$$
$$\leq -\eta \frac{\sum_i \ell_t(e_i) w_t^i}{\sum_i w_t^i} + \frac{\eta^2}{8}$$

#### Proof idea:

$$\ln \frac{W_{t+1}}{W_t} = \ln \left( \frac{\sum_{i=1}^m \exp(-\eta \ell_t(e_i)) w_t^i}{\sum_i w_t^i} \right)$$
$$\leq -\eta \frac{\sum_i \ell_t(e_i) w_t^i}{\sum_i w_t^i} + \frac{\eta^2}{8}$$

where we have used Hoeffding's inequality: for a random variable  $X \in [a, b]$  and  $\lambda \in \mathbb{R}$ ,

$$\ln\left(\mathbb{E}e^{\lambda X}\right) \leq \lambda \mathbb{E}X + \frac{\lambda^2(b-a)^2}{8}$$

(人間) ト イヨト イヨト

#### Proof idea:

$$\ln \frac{W_{t+1}}{W_t} = \ln \left( \frac{\sum_{i=1}^m \exp(-\eta \ell_t(e_i)) w_t^i}{\sum_i w_t^i} \right)$$
$$\leq -\eta \frac{\sum_i \ell_t(e_i) w_t^i}{\sum_i w_t^i} + \frac{\eta^2}{8}$$
$$= -\eta \ell_t(a_t) + \frac{\eta^2}{8},$$

where we have used Hoeffding's inequality: for a random variable  $X \in [a, b]$  and  $\lambda \in \mathbb{R}$ ,

$$\ln\left(\mathbb{E}e^{\lambda X}\right) \leq \lambda \mathbb{E}X + \frac{\lambda^2(b-a)^2}{8}$$

26/132

(人間) トイヨト イヨー

$$-\eta \inf_{a\in\mathbb{R}^d}\sum_{t=1}^n \ell_t(a) - \ln m \leq \ln \frac{W_{n+1}}{W_1} \leq -\eta \sum_{t=1}^n \ell_t(a_t) + \frac{n\eta^2}{8}.$$

$$-\eta \inf_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) - \ln m \le \ln \frac{W_{n+1}}{W_1} \le -\eta \sum_{t=1}^n \ell_t(a_t) + \frac{n\eta^2}{8}.$$
  
Thus,  
$$R_n \le \frac{\ln m}{\eta} + \frac{\eta n}{8}.$$

#### Proof idea:

$$-\eta \inf_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) - \ln m \le \ln \frac{W_{n+1}}{W_1} \le -\eta \sum_{t=1}^n \ell_t(a_t) + \frac{n\eta^2}{8}.$$
  
Thus,  
$$R_n \le \frac{\ln m}{n} + \frac{\eta n}{8}.$$

 $\eta$ 

Choosing the optimal  $\eta$  gives the result:

# Prediction with Expert Advice

### Proof idea:

Thus,

$$-\eta \inf_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) - \ln m \le \ln \frac{W_{n+1}}{W_1} \le -\eta \sum_{t=1}^n \ell_t(a_t) + \frac{n\eta^2}{8}.$$

$$R_n \leq rac{\ln m}{\eta} + rac{\eta n}{8}.$$

Choosing the optimal  $\eta$  gives the result:

#### Theorem

The exponential weights strategy with parameter  $\eta = \sqrt{8 \ln m/n}$  has regret no more than  $\sqrt{\frac{n \ln m}{2}}$ .

## Key Points

For a finite set of actions (experts):

### Key Points

For a finite set of actions (experts):

• If one action is perfect (i.e., has zero loss), the minimax strategy gives per round regret of

 $\frac{\log_4 m}{n}.$ 

### Key Points

For a finite set of actions (experts):

• Exponential weights gives per round regret of

 $\sqrt{\frac{\ln m}{2n}}.$ 

・ロ ・ ・ 日 ・ ・ 三 ・ ・ 三 ・ 三 ・ つ Q (\* 28 / 132)

**(**) boundedness:  $\ell_t(e_i) \in [0, 1]$  (for Hoeffding's inequality), and

**(**) boundedness:  $\ell_t(e_i) \in [0, 1]$  (for Hoeffding's inequality), and

linearity,

$$\ell_t(a_t) = \sum_i \ell_t(e_i) w_t^i.$$

**(**) boundedness:  $\ell_t(e_i) \in [0, 1]$  (for Hoeffding's inequality), and

linearity,

$$\ell_t(a_t) = \sum_i \ell_t(e_i) w_t^i.$$

For linearity, an inequality would have sufficed,

$$\ell_t(a_t) \leq \sum_i \ell_t(e_i) w_t^i,$$

**(**) boundedness:  $\ell_t(e_i) \in [0, 1]$  (for Hoeffding's inequality), and

linearity,

$$\ell_t(a_t) = \sum_i \ell_t(e_i) w_t^i.$$

For linearity, an inequality would have sufficed,

$$\ell_t(a_t) \leq \sum_i \ell_t(e_i) w_t^i,$$

which corresponds to convexity of  $\ell_t$ .

# Online convex optimization

### Binary prediction

- With (perfect) expert advice
- Minimax strategy
- With imperfect experts: exponential weights

#### ② General online convex

- Empirical minimization fails
- Gradient algorithm.
- A regularization viewpoint
- Bregman divergence
- Properties of regularization
- Linearization
- Mirror descent
- Regret bounds
- Strongly convex losses
- Adaptive regularization
- Minimax strategies

### The problem

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

### The problem

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

## Minimax regret

$$\min_{a_1} \max_{\ell_1} \cdots \min_{a_n} \max_{\ell_n} \left( \sum_{t=1}^n \ell_t(a_t) - \min_{a \in \mathcal{A}} \sum_{t=1}^n \ell_t(a) \right).$$

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  $\mathcal{A} = [-1, 1]$ ,  $\mathcal{L} = \{a \mapsto v \cdot a : |v| \leq 1\}$ .

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

$$a_1=0,$$

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

$$\ell_1(a) = rac{1}{2}a$$
 $a_1 = 0,$ 

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

$$\ell_1(a) = rac{1}{2}a, \ a_1 = 0, \qquad a_2 = -1,$$

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

$$\ell_1(a) = rac{1}{2}a, \ \ \ell_2(a) = -a, \ a_1 = 0, \ \ \ a_2 = -1,$$

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')

イロト 不得下 イヨト イヨト 二日

32 / 132

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

$$\ell_1(a) = rac{1}{2}a, \ \ \ell_2(a) = -a, \ a_1 = 0, \ \ \ a_2 = -1, \ \ \ \ a_3 = 1,$$

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

$$\ell_1(a) = rac{1}{2}a, \ \ \ell_2(a) = -a, \ \ \ell_3(a) = a, \ \ a_1 = 0, \ \ \ a_2 = -1, \ \ \ a_3 = 1,$$

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:

$$\ell_1(a) = \frac{1}{2}a, \quad \ell_2(a) = -a, \quad \ell_3(a) = a,$$
  
 $a_1 = 0, \qquad a_2 = -1, \qquad a_3 = 1, \qquad a_4 = -1,$ 

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

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

$$\ell_1(a) = rac{1}{2}a, \ \ell_2(a) = -a, \ \ell_3(a) = a, \ \ell_4(a) = -a, \ a_1 = 0, \ a_2 = -1, \ a_3 = 1, \ a_4 = -1,$$

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

$$\ell_1(a) = \frac{1}{2}a, \quad \ell_2(a) = -a, \quad \ell_3(a) = a, \quad \ell_4(a) = -a,$$
  
 $a_1 = 0, \qquad a_2 = -1, \qquad a_3 = 1, \qquad a_4 = -1, \qquad a_5 = 1,$ 

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

$$\ell_1(a) = \frac{1}{2}a, \quad \ell_2(a) = -a, \quad \ell_3(a) = a, \quad \ell_4(a) = -a, \quad \ell_5(a) = a, \\ a_1 = 0, \qquad a_2 = -1, \qquad a_3 = 1, \qquad a_4 = -1, \qquad a_5 = 1,$$

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:

$$\ell_1(a) = rac{1}{2}a, \ \ell_2(a) = -a, \ \ell_3(a) = a, \ \ell_4(a) = -a, \ \ell_5(a) = a, \ a_1 = 0, \ a_2 = -1, \ a_3 = 1, \ a_4 = -1, \ a_5 = 1,$$

•  $a^* = 0$  shows  $\min_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) \le 0$ , but  $\sum_{t=1}^n \ell_t(a_t) = n-1$ .

• Choosing  $a_t$  to minimize past losses can fail.

- Choosing  $a_t$  to minimize past losses can fail.
- The strategy must avoid overfitting, just as in probabilistic settings.

- Choosing *a*<sub>t</sub> 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.

- Choosing *a*<sub>t</sub> 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

### Binary prediction

- ② General online convex
  - Empirical minimization fails
  - Gradient algorithm.
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
- Minimax strategies

 $a_1 \in \mathcal{A},$ 

 < □ > < □ > < □ > < Ξ > < [Zinkevich, 2003] 35 / 132

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

 < □ > < □ > < □ > < Ξ > < [Zinkevich, 2003] 35 / 132

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

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



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

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

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

$$\mathsf{a}_1 \in \mathcal{A}, \qquad \qquad \mathsf{a}_{t+1} = \mathsf{\Pi}_\mathcal{A} \left( \mathsf{a}_t - \eta 
abla \ell_t(\mathsf{a}_t) 
ight),$$

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

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

#### Theorem

the gradient strategy with  $\eta = D/(G\sqrt{n})$ 

$$\mathsf{a}_1 \in \mathcal{A}, \qquad \qquad \mathsf{a}_{t+1} = \mathsf{\Pi}_\mathcal{A} \left( \mathsf{a}_t - \eta 
abla \ell_t(\mathsf{a}_t) 
ight),$$

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  $G = \max_t \|\nabla \ell_t(a_t)\|$ the gradient strategy with  $\eta = D/(G\sqrt{n})$ 

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

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  $D = \operatorname{diam}(\mathcal{A})$ , the gradient strategy with  $\eta = D/(G\sqrt{n})$ 

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

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  $D = \operatorname{diam}(\mathcal{A})$ , the gradient strategy with  $\eta = D/(G\sqrt{n})$  has regret satisfying

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

 < □ > < □ > < □ > < □ > < [Zinkevich, 2003] 35 / 132

#### Example

 $\mathcal{A} = \{ \mathbf{a} \in \mathbb{R}^d : \|\mathbf{a}\| \leq 1 \},$ 

#### Example

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

(ロ)、(型)、(注)、(注)、(注)、(注)、(2)、(36/132)

#### Example

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

#### Example

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

#### Example

$$\begin{aligned} \mathcal{A} &= \{ a \in \mathbb{R}^d : \|a\| \leq 1 \}, \ \mathcal{L} &= \{ a \mapsto v \cdot a : \|v\| \leq 1 \}. \\ D &= 2, \ G \leq 1. \\ \text{Regret is no more than } 2\sqrt{n}. \end{aligned}$$

### Example

$$\begin{aligned} \mathcal{A} &= \{ a \in \mathbb{R}^d : \|a\| \leq 1 \}, \ \mathcal{L} &= \{ a \mapsto v \cdot a : \|v\| \leq 1 \}. \\ D &= 2, \ G \leq 1. \\ \text{Regret is no more than } 2\sqrt{n}. \end{aligned}$$

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



### Example

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

### Example

$$\mathcal{A} = \Delta^m, \ \mathcal{L} = \{ a \mapsto v \cdot a : \|v\|_{\infty} \le 1 \}.$$
  
 
$$D = \sqrt{2},$$

### Example

$$\mathcal{A} = \Delta^m, \ \mathcal{L} = \{ a \mapsto v \cdot a : \|v\|_{\infty} \le 1 \}.$$
  
 
$$D = \sqrt{2}, \ G \le \sqrt{m}.$$

37 / 132

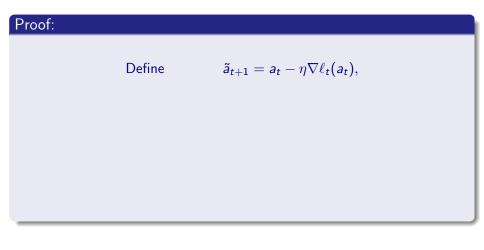
#### Example

 $\begin{array}{l} \mathcal{A} = \Delta^m, \ \mathcal{L} = \{a \mapsto v \cdot a : \|v\|_{\infty} \leq 1\}. \\ D = \sqrt{2}, \ G \leq \sqrt{m}. \\ \text{Regret is no more than } \sqrt{2mn}. \end{array}$ 

#### Example

 $\begin{array}{l} \mathcal{A} = \Delta^m, \ \mathcal{L} = \{a \mapsto v \cdot a : \|v\|_{\infty} \leq 1\}, \\ D = \sqrt{2}, \ G \leq \sqrt{m}. \\ \text{Regret is no more than } \sqrt{2mn}. \end{array}$ 

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



### Proof:

Define	$egin{aligned} & ilde{a}_{t+1} = a_t - \eta  abla \ell_t(a_t), \ & ilde{a}_{t+1} = \Pi_\mathcal{A}( ilde{a}_{t+1}). \end{aligned}$

Define	$\tilde{a}_{t+1} = a_t - \eta \nabla \ell_t(a_t),$
	$a_{t+1} = \Pi_\mathcal{A}(\widetilde{a}_{t+1}).$

Fix a comparator  $a \in \mathcal{A}$ 

◆□ → < □ → < 三 → < 三 → < 三 → < 三 → ○ < ○ 38/132

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

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

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

Fix a comparator  $a \in \mathcal{A}$  and consider the measure of progress  $\|a_t - a\|$ .  $\|a_{t+1} - a\|^2$ 

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

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

$$\|a_{t+1} - a\|^2 \le \|\tilde{a}_{t+1} - a\|^2$$

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

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

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

◆□ → ◆□ → ◆ ■ → ◆ ■ → ● ● 今 Q ↔
38/132

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

By convexity,

$$\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) \leq \sum_{t=1}^n 
abla \ell_t(a_t) \cdot (a_t - a)$$

By convexity,

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

By convexity,

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

◆□ → < □ → < 亘 → < 亘 → < 亘 → < 亘 → ○ Q (~ 39/132

# **Online Convex Optimization**

### Binary prediction

- ② General online convex
  - Empirical minimization fails
  - Gradient algorithm
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
- Minimax strategies

An observation: gradient algorithm is regularized minimization.

An observation: gradient algorithm is regularized minimization.

• Suppose  $\ell_t$  is linear:  $\ell_t(a) = g_t \cdot a$ .

An observation: gradient algorithm is regularized minimization.

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

An observation: gradient algorithm is regularized minimization.

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

An observation: gradient algorithm is regularized minimization.

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

$$egin{split} egin{aligned} egin{aligned} egin{aligned} eta_{t+1} &= rg\min_{oldsymbol{a}\in\mathcal{A}}\left(\eta\sum_{s=1}^t\ell_s(oldsymbol{a}) + rac{1}{2}\|oldsymbol{a}\|^2
ight) \end{split}$$

corresponds to the gradient step  $a_{t+1} = a_t - \eta \nabla \ell_t(a_t)$ .

• Indeed, setting the derivative to zero gives

An observation: gradient algorithm is regularized minimization.

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

$$m{a}_{t+1} = rg\min_{m{a}\in\mathcal{A}}\left(\eta\sum_{s=1}^t \ell_s(m{a}) + rac{1}{2}\|m{a}\|^2
ight)$$

41 / 132

corresponds to the gradient step  $a_{t+1} = a_t - \eta \nabla \ell_t(a_t)$ .

Indeed, setting the derivative to zero gives

$$a_{t+1} = -\eta \sum_{s=1}^t \nabla \ell_s$$

An observation: gradient algorithm is regularized minimization.

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

$$m{a}_{t+1} = rg\min_{m{a}\in\mathcal{A}}\left(\eta\sum_{s=1}^t \ell_s(m{a}) + rac{1}{2}\|m{a}\|^2
ight)$$

corresponds to the gradient step  $a_{t+1} = a_t - \eta \nabla \ell_t(a_t)$ .

Indeed, setting the derivative to zero gives

$$a_{t+1} = -\eta \sum_{s=1}^{t} \nabla \ell_s$$
$$a_t = -\eta \sum_{s=1}^{t-1} \nabla \ell_s.$$

An observation: gradient algorithm is regularized minimization.

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

$$m{a}_{t+1} = rg\min_{m{a}\in\mathcal{A}}\left(\eta\sum_{s=1}^t \ell_s(m{a}) + rac{1}{2}\|m{a}\|^2
ight)$$

corresponds to the gradient step  $a_{t+1} = a_t - \eta \nabla \ell_t(a_t)$ .

• Indeed, setting the derivative to zero gives

$$a_{t+1} = -\eta \sum_{s=1}^{t} \nabla \ell_s = a_t - \eta \nabla \ell_t.$$
$$a_t = -\eta \sum_{s=1}^{t-1} \nabla \ell_s.$$

41/132

### Definition: Regularized minimization

Consider the family of strategies of the form:

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

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

# Online Convex Optimization: Regularization

### Regularized minimization

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

#### 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.

#### 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<sub>t+1</sub> as trading off two competing forces: making l<sub>t</sub>(a<sub>t+1</sub>) small, and keeping a<sub>t+1</sub> close to a<sub>t</sub>.

#### 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<sub>t+1</sub> as trading off two competing forces: making l<sub>t</sub>(a<sub>t+1</sub>) small, and keeping a<sub>t+1</sub> close to a<sub>t</sub>.
- This is a perspective that motivated many algorithms in the literature. We'll investigate why regularized minimization can be viewed this way.

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.

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

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

#### Definition

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

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

#### Definition

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

So

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

・・・・
 ・・・
 ・・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・
 ・・

#### Definition: Bregman Divergence

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

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

[Bregman, 1967]

45 / 132

・ロット 全部 マイヨット キョン

#### Definition: Bregman Divergence

For a strictly convex, differentiable  $\Phi : \mathbb{R}^d \to \mathbb{R}$ , 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, 1967]

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

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

#### Example

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

$$D_{\Phi}(a,b) = rac{1}{2} \|a\|^2 - \left(rac{1}{2} \|b\|^2 + b \cdot (a-b)
ight)$$

<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

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

#### Example

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

$$egin{aligned} D_{\Phi}(a,b) &= rac{1}{2} \|a\|^2 - \left(rac{1}{2} \|b\|^2 + b \cdot (a-b)
ight) \ &= rac{1}{2} \|a-b\|^2, \end{aligned}$$

・ロ ・ ・ 一 ・ ・ 注 ト ・ 注 ・ 注 ・ つ へ (\*) 46 / 132

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

#### Example

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

$$egin{aligned} D_{\Phi}(a,b) &= rac{1}{2} \|a\|^2 - \left(rac{1}{2} \|b\|^2 + b \cdot (a-b)
ight) \ &= rac{1}{2} \|a-b\|^2, \end{aligned}$$

the squared euclidean norm.

・ロ ・ ・ 日 ・ ・ 目 ・ 日 ・ 日 ・ 日 ・ の へ ペ 46 / 132

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

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

### Example

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

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

### Example

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

$$D_{\Phi}(a,b) = \sum_{i} (a_i(\ln a_i - 1) - b_i(\ln b_i - 1) - \ln b_i(a_i - b_i))$$

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

#### Example

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

$$egin{aligned} D_{\Phi}(a,b) &= \sum_i \left( a_i (\ln a_i - 1) - b_i (\ln b_i - 1) - \ln b_i (a_i - b_i) 
ight) \ &= \sum_i \left( a_i \ln rac{a_i}{b_i} + b_i - a_i 
ight), \end{aligned}$$

the unnormalized KL 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

$$egin{aligned} D_{\Phi}(a,b) &= \sum_i \left(a_i(\ln a_i-1) - b_i(\ln b_i-1) - \ln b_i(a_i-b_i)
ight) \ &= \sum_i \left(a_i\ln rac{a_i}{b_i} + b_i - a_i
ight), \end{aligned}$$

the unnormalized KL divergence.

Thus, for  $a \in \Delta^d$ ,  $\Phi(a) = \sum_i a_i \ln a_i$  has  $D_{\phi}(a, b) = \sum_i a_i \ln \frac{a_i}{b_i}$ .

• The interior of  $\mathcal{A}$  is convex,

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

• The interior of  $\mathcal{A}$  is convex,

• For a sequence approaching the boundary of  $\mathcal{A}$ ,  $\|\nabla \Phi(a_n)\| \to \infty$ . We say that such a  $\Phi$  is a *Legendre function*.

## Properties



**1**  $D_{\Phi} \ge 0, \ D_{\Phi}(a, a) = 0.$ 

#### Properties

- **1**  $D_{\Phi} \geq 0$ ,  $D_{\Phi}(a, a) = 0$ .
- $D_{A+B} = D_A + D_B.$

- $D_{\Phi} \ge 0, \ D_{\Phi}(a, a) = 0.$
- $D_{A+B} = D_A + D_B.$
- Segman projection, ⊓<sup>Φ</sup><sub>A</sub>(b) = arg min<sub>a∈A</sub> D<sub>Φ</sub>(a, b) is uniquely defined for closed, convex A.

- **1**  $D_{\Phi} \geq 0, \ D_{\Phi}(a, a) = 0.$
- $D_{A+B} = D_A + D_B.$
- Segman projection, Π<sup>Φ</sup><sub>A</sub>(b) = arg min<sub>a∈A</sub> D<sub>Φ</sub>(a, b) is uniquely defined for closed, convex A.
- Generalized Pythagorus: for closed, convex  $\mathcal{A}$ ,  $a^* = \Pi^{\Phi}_{\mathcal{A}}(b)$ ,  $a \in \mathcal{A}$ :

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

- **1**  $D_{\Phi} \geq 0, \ D_{\Phi}(a, a) = 0.$
- $D_{A+B} = D_A + D_B.$
- Bregman projection, ∏<sup>Φ</sup><sub>A</sub>(b) = arg min<sub>a∈A</sub> D<sub>Φ</sub>(a, b) is uniquely defined for closed, convex A.
- **(**) Generalized Pythagorus: for closed, convex  $\mathcal{A}$ ,  $a^* = \prod_{\mathcal{A}}^{\Phi}(b)$ ,  $a \in \mathcal{A}$ :

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

< □ > < □ > < □ > < ⊇ > < ⊇ > < ⊇ > < ⊇ > ○ Q (C 50 / 132

• For  $\ell$  affine,  $D_{\Phi+\ell} = D_{\Phi}$ .

- For  $\ell$  affine,  $D_{\Phi+\ell} = D_{\Phi}$ .
- **③** For  $\Phi^*$  the Legendre dual of  $\Phi$ ,

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

- For  $\ell$  affine,  $D_{\Phi+\ell} = D_{\Phi}$ .
- **③** For  $\Phi^*$  the Legendre dual of  $\Phi$ ,

$$abla \Phi^* = (
abla \Phi)^{-1},$$
  
 $D_{\Phi}(a, b) = D_{\Phi^*}(
abla \phi(b), 
abla \phi(a)).$ 

## Definition: Legendre Dual

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

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

### Definition: Legendre Dual

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

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

•  $\Phi^*$  is Legendre.

### Definition: Legendre Dual

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

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

- Φ<sup>\*</sup> is Legendre.
- dom( $\Phi^*$ ) =  $\nabla \Phi(\text{int dom } \Phi)$ .

#### Definition: Legendre Dual

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

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

- Φ<sup>\*</sup> is Legendre.
- dom( $\Phi^*$ ) =  $\nabla \Phi(\text{int dom } \Phi)$ .
- $\nabla \Phi^* = (\nabla \Phi)^{-1}$ .

#### Definition: Legendre Dual

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

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

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

#### Definition: Legendre Dual

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

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

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

# 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.

イロン イヨン イヨン イヨン 三日

52/132

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

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

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

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

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**

#### Binary prediction

- ② General online convex
  - Empirical minimization fails
  - Gradient algorithm
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
- Minimax strategies

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

#### Theorem

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

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

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

#### Theorem

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

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

Then

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

・ロ ・ ・ 一 ・ ・ 三 ・ ・ 三 ・ 三 ・ つ へ ()
54/132

By the definition of  $\Phi_t$ ,

## $\left(\Phi_t(a) := \eta \sum_{s=1}^t \ell_s(a) + R(a)\right)$

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

By the definition of  $\Phi_t$ ,

$$\left(\Phi_t(a) := \eta \sum_{s=1}^t \ell_s(a) + R(a)\right)$$

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

The derivative wrt *a* is

$$egin{aligned} 
abla \Phi_t(a) &- 
abla \Phi_{t-1}(a) + 
abla_a D_{\Phi_{t-1}}(a, \tilde{a}_t) \ &= 
abla \Phi_t(a) - 
abla \Phi_{t-1}(a) + 
abla \Phi_{t-1}(a) - 
abla \Phi_{t-1}(\tilde{a}_t) \end{aligned}$$

By the definition of  $\Phi_t$ ,  $(\Phi_t(a) := \eta \sum_{s=1}^t \ell_s(a) + R(a))$  $\eta \ell_t(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t) = \Phi_t(a) - \Phi_{t-1}(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t).$ The derivative wrt a is

$$egin{aligned} & 
abla \Phi_{t}(a) - 
abla \Phi_{t-1}(a) + 
abla_{\Phi_{t-1}}(a, a_t) \ &= 
abla \Phi_t(a) - 
abla \Phi_{t-1}(a) + 
abla \Phi_{t-1}(a) - 
abla \Phi_{t-1}( ilde{a}_t) \end{aligned}$$

Setting to zero shows that

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

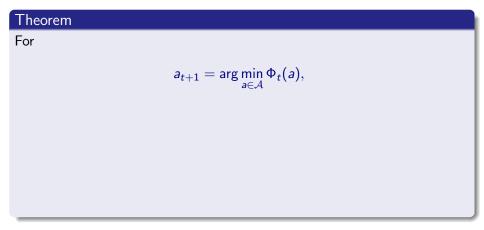
By the definition of  $\Phi_t$ ,  $(\Phi_t(a) := \eta \sum_{s=1}^t \ell_s(a) + R(a))$   $\eta \ell_t(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t) = \Phi_t(a) - \Phi_{t-1}(a) + D_{\Phi_{t-1}}(a, \tilde{a}_t).$ The derivative wrt *a* is  $\nabla \Phi_t(a) - \nabla \Phi_{t-1}(a) + \nabla_s D_{\Phi_{t-1}}(a, \tilde{a}_t)$ 

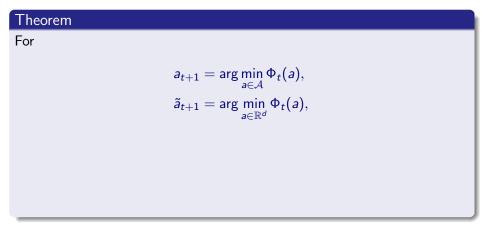
$$= \nabla \Phi_t(a) - \nabla \Phi_{t-1}(a) + \nabla \Phi_{t-1}(a) - \nabla \Phi_{t-1}(\tilde{a}_t)$$

Setting to zero shows that

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

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





Theorem	
For	
	$(x) \rightarrow (x)$
	$a_{t+1} = \arg\min_{a \in \mathcal{A}} \Phi_t(a),$
	$\widetilde{a}_{t+1} = rg\min_{a\in \mathbb{R}^d} \Phi_t(a),$
we have	
	$a_{t+1} = \Pi^{igoplus_t}_A(\widetilde{a}_{t+1}).$

Let  $a'_{t+1}$  denote  $\Pi^{\Phi_t}_{\mathcal{A}}(\tilde{a}_{t+1})$ .



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

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,

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

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,

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

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

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

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

Conversely,

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

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

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

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

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

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

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

イロト 不同下 イヨト イヨト

58 / 132

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

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

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

イロト 不得下 イヨト イヨト 二日

58/132

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

(e.g., *R* squared Euclidean norm)

# **Online Convex Optimization**

#### Binary prediction

- ② General online convex
  - Empirical minimization fails
  - Gradient algorithm
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
- Minimax strategies

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

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

#### Theorem

Any strategy for online linear optimization, with regret satisfying

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

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

#### Theorem

Any strategy for online linear optimization, with regret satisfying

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

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

$$\sum_{t=1}^n \ell_t(a_t) - \min_{a \in \mathcal{A}} \sum_{t=1}^n \ell_t(a) \leq C_n(\nabla \ell_1(a_1), \dots, \nabla \ell_n(a_n)).$$

◆□ → ◆□ → ◆ ■ → ◆ ■ → ● ● ○ へ ○
60/132

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

#### Theorem

Any strategy for online linear optimization, with regret satisfying

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

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

$$\sum_{t=1}^n \ell_t(a_t) - \min_{a \in \mathcal{A}} \sum_{t=1}^n \ell_t(a) \leq C_n(\nabla \ell_1(a_1), \ldots, \nabla \ell_n(a_n)).$$

#### Proof:

Convexity implies  $\ell_t(a_t) - \ell_t(a) \leq \nabla \ell_t(a_t) \cdot (a_t - a)$ .

#### Key Point:

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

#### 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  $\ell_t$ .

# Online convex optimization

#### Binary prediction

- ② General online convex
  - Empirical minimization fails
  - Gradient algorithm
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
- Minimax strategies

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

see [Nemirovsky and Yudin, 1983]

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

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

see [Nemirovsky and Yudin, 1983]

63 / 132

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

see [Nemirovsky and Yudin, 1983]

63 / 132

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  $\tilde{a}_{t+1}$ .

see [Nemirovsky and Yudin, 1983]

For the unconstrained minimization, we have

For the unconstrained minimization, we have

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

For the unconstrained minimization, we have

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

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

For the unconstrained minimization, we have

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

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

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

・ロ ・ ・ 日 ・ ・ 注 ・ く 注 ・ 注 ・ う Q (\*) 64/132

### **Online Convex Optimization**

### Binary prediction

- ② General online convex
  - Empirical minimization fails
  - Gradient algorithm
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
- Minimax strategies

### Recall: 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  $R : \mathbb{R}^d \to \mathbb{R}$  is strictly convex and differentiable.

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

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a)$$
  
=  $\frac{D_R(a, a_1) - D_{\Phi_n}(a, a_{n+1})}{\eta} + \frac{1}{\eta} \sum_{t=1}^{n} D_{\Phi_t}(a_t, a_{t+1}),$ 

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

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a)$$
  
=  $\frac{D_R(a, a_1) - D_{\Phi_n}(a, a_{n+1})}{\eta} + \frac{1}{\eta} \sum_{t=1}^{n} D_{\Phi_t}(a_t, a_{t+1}),$ 

and thus

$$\sum_{t=1}^n \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(a) + \frac{D_R(a, a_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(a_t, a_{t+1}).$$

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

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a)$$
  
=  $\frac{D_R(a, a_1) - D_{\Phi_n}(a, a_{n+1})}{\eta} + \frac{1}{\eta} \sum_{t=1}^{n} D_{\Phi_t}(a_t, a_{t+1}),$ 

and thus

$$\sum_{t=1}^n \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(a) + \frac{D_R(a, a_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(a_t, a_{t+1}).$$

So the sizes of the steps  $D_{\Phi_t}(a_t, a_{t+1})$  determine the regret bound.

67/132

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

$$\sum_{t=1}^n \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(a) + \frac{D_R(a, a_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(a_t, a_{t+1}).$$

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

$$\sum_{t=1}^n \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(a) + \frac{D_R(a, a_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(a_t, a_{t+1}).$$

Notice that, because  $a_{t+1}$  is the unconstrained minimizer of  $\Phi_t$ ,

$$D_{\Phi_t}(a_t, a_{t+1}) = D_{\Phi_t^*}(\nabla \Phi_t(a_{t+1}), \nabla \Phi_t(a_t))$$

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

$$\sum_{t=1}^n \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(a) + \frac{D_R(a, a_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(a_t, a_{t+1}).$$

Notice that, because  $a_{t+1}$  is the unconstrained minimizer of  $\Phi_t$ ,

$$egin{aligned} D_{\Phi_t}(a_t,a_{t+1}) &= D_{\Phi_t^*}(
abla \Phi_t(a_{t+1}),
abla \Phi_t(a_t)) \ &= D_{\Phi_t^*}(0,
abla \Phi_{t-1}(a_t) + \eta 
abla \ell_t(a_t)) \end{aligned}$$

68 / 132

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

$$\sum_{t=1}^n \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(a) + \frac{D_R(a, a_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(a_t, a_{t+1}).$$

Notice that, because  $a_{t+1}$  is the unconstrained minimizer of  $\Phi_t$ ,

$$egin{aligned} D_{\Phi_t}(a_t,a_{t+1}) &= D_{\Phi_t^*}(
abla \Phi_t(a_{t+1}),
abla \Phi_t(a_t)) \ &= D_{\Phi_t^*}(0,
abla \Phi_{t-1}(a_t) + \eta 
abla \ell_t(a_t)) \ &= D_{\Phi_t^*}(0,\eta 
abla \ell_t(a_t)). \end{aligned}$$

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

$$\sum_{t=1}^n \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \left( \sum_{t=1}^n \ell_t(a) + \frac{D_R(a, a_1)}{\eta} \right) + \frac{1}{\eta} \sum_{t=1}^n D_{\Phi_t}(a_t, a_{t+1}).$$

Notice that, because  $a_{t+1}$  is the unconstrained minimizer of  $\Phi_t$ ,

$$D_{\Phi_t}(a_t, a_{t+1}) = D_{\Phi_t^*}(\nabla \Phi_t(a_{t+1}), \nabla \Phi_t(a_t))$$
  
=  $D_{\Phi_t^*}(0, \nabla \Phi_{t-1}(a_t) + \eta \nabla \ell_t(a_t))$   
=  $D_{\Phi_t^*}(0, \eta \nabla \ell_t(a_t)).$ 

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

### Example

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

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

### Example

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

$$\sum_{t=1}^{n} \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \sum_{t=1}^{n} \ell_t(a) + \frac{\|a^* - a_1\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{n} \|g_t\|^2.$$

#### Example

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

$$\sum_{t=1}^{n} \ell_t(a_t) \leq \inf_{a \in \mathbb{R}^d} \sum_{t=1}^{n} \ell_t(a) + \frac{\|a^* - a_1\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^{n} \|g_t\|^2.$$

And if  $||g_t|| \le G$  and  $||a^* - a_1|| \le D$ , choosing  $\eta$  appropriately gives regret  $\le DG\sqrt{n}$ .

### Seeing the future gives small regret

For regularized minimization, that is,

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

see also [Kalai and Vempala, 2005]

### Seeing the future gives small regret

For regularized minimization, that is,

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

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

$$\sum_{t=1}^{n} \ell_t(a_{t+1}) - \sum_{t=1}^{n} \ell_t(a) \leq \frac{1}{\eta}(R(a) - R(a_1)).$$

see also [Kalai and Vempala, 2005]

### Seeing the future gives small regret

For regularized minimization, that is,

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

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

$$\sum_{t=1}^{n} \ell_t(a_{t+1}) - \sum_{t=1}^{n} \ell_t(a) \leq \frac{1}{\eta}(R(a) - R(a_1)).$$

(NB: This is cheating!)

see also [Kalai and Vempala, 2005]



### Proof:

Since  $a_{t+1}$  minimizes  $\Phi_t$ ,

### Proof:

Since  $a_{t+1}$  minimizes  $\Phi_t$ ,

$$\eta \sum_{s=1}^{t} \ell_s(a) + R(a) \ge \eta \sum_{s=1}^{t} \ell_s(a_{t+1}) + R(a_{t+1})$$

### Proof:

Since  $a_{t+1}$  minimizes  $\Phi_t$ ,

$$\eta \sum_{s=1}^{t} \ell_s(a) + R(a) \ge \eta \sum_{s=1}^{t} \ell_s(a_{t+1}) + R(a_{t+1})$$
$$= \eta \ell_t(a_{t+1}) + \eta \sum_{s=1}^{t-1} \ell_s(a_{t+1}) + R(a_{t+1})$$

### Proof:

Since  $a_{t+1}$  minimizes  $\Phi_t$ , and  $a_t$  minimizes  $\Phi_{t-1}$ ,

$$\eta \sum_{s=1}^{t} \ell_s(a) + R(a) \ge \eta \sum_{s=1}^{t} \ell_s(a_{t+1}) + R(a_{t+1})$$
$$= \eta \ell_t(a_{t+1}) + \eta \sum_{s=1}^{t-1} \ell_s(a_{t+1}) + R(a_{t+1})$$
$$\ge \eta \ell_t(a_{t+1}) + \eta \sum_{s=1}^{t-1} \ell_s(a_t) + R(a_t)$$

### Proof:

Since  $a_{t+1}$  minimizes  $\Phi_t$ , and  $a_t$  minimizes  $\Phi_{t-1}$ ,

$$\eta \sum_{s=1}^{t} \ell_{s}(a) + R(a) \ge \eta \sum_{s=1}^{t} \ell_{s}(a_{t+1}) + R(a_{t+1})$$
$$= \eta \ell_{t}(a_{t+1}) + \eta \sum_{s=1}^{t-1} \ell_{s}(a_{t+1}) + R(a_{t+1})$$
$$\ge \eta \ell_{t}(a_{t+1}) + \eta \sum_{s=1}^{t-1} \ell_{s}(a_{t}) + R(a_{t})$$
$$\vdots$$
$$\ge \eta \sum_{s=1}^{t} \ell_{s}(a_{s+1}) + R(a_{1}).$$

### Theorem

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

$$\sum_{t=1}^{n} \ell_t(a_{t+1}) - \sum_{t=1}^{n} \ell_t(a) \leq \frac{1}{\eta}(R(a) - R(a_1)).$$

<ロ > < 部 > < 言 > < 言 > 言 の < C 72/132

### Theorem

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

$$\sum_{t=1}^n \ell_t(a_{t+1}) - \sum_{t=1}^n \ell_t(a) \le rac{1}{\eta}(R(a) - R(a_1)).$$

Thus, if  $a_t$  and  $a_{t+1}$  are close, then regret is small:

### Corollary

For all  $a \in A$ ,

$$\sum_{t=1}^n \left(\ell_t(\mathsf{a}_t)-\ell_t(\mathsf{a})
ight) \leq \sum_{t=1}^n \left(\ell_t(\mathsf{a}_t)-\ell_t(\mathsf{a}_{t+1})
ight) + rac{1}{\eta}\left(R(\mathsf{a})-R(\mathsf{a}_1)
ight).$$

#### Theorem

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

$$\sum_{t=1}^n \ell_t(a_{t+1}) - \sum_{t=1}^n \ell_t(a) \le rac{1}{\eta}(R(a) - R(a_1)).$$

Thus, if  $a_t$  and  $a_{t+1}$  are close, then regret is small:

### Corollary

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

$$\sum_{t=1}^n \left(\ell_t(\mathsf{a}_t)-\ell_t(\mathsf{a})
ight) \leq \sum_{t=1}^n \left(\ell_t(\mathsf{a}_t)-\ell_t(\mathsf{a}_{t+1})
ight) + rac{1}{\eta}\left(\mathsf{R}(\mathsf{a})-\mathsf{R}(\mathsf{a}_1)
ight).$$

So how can we control the increments  $\ell_t(a_t) - \ell_t(a_{t+1})$ ?

### 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.

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

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

where  $a_{t+1}$  minimizes  $\Phi_t$  and  $\|\cdot\|_*$  is the dual norm to  $\|\cdot\|_:$ 

 $\|v\|_* = \sup\{v \cdot a : \|a\| \le 1\}.$ 

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

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

where  $a_{t+1}$  minimizes  $\Phi_t$  and  $\|\cdot\|_*$  is the dual norm to  $\|\cdot\|$ :

 $\|v\|_* = \sup\{v \cdot a : \|a\| \le 1\}.$ 

Note that the definition implies a generalization of the Cauchy-Schwarz inequality: for ||a|| > 0,

$$\mathbf{v}\cdot\frac{\mathbf{a}}{\|\mathbf{a}\|}\leq \|\mathbf{v}\|_*.$$

### Proof:

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

# Proof:

$$egin{aligned} &R(a_t) \geq R(a_{t+1}) + 
abla R(a_{t+1}) \cdot (a_t - a_{t+1}) + rac{1}{2} \|a_t - a_{t+1}\|^2, \ &R(a_{t+1}) \geq R(a_t) + 
abla R(a_t) \cdot (a_{t+1} - a_t) + rac{1}{2} \|a_t - a_{t+1}\|^2. \end{aligned}$$

## Proof:

$$egin{aligned} &R(a_t) \geq R(a_{t+1}) + 
abla R(a_{t+1}) \cdot (a_t - a_{t+1}) + rac{1}{2} \|a_t - a_{t+1}\|^2, \ &R(a_{t+1}) \geq R(a_t) + 
abla R(a_t) \cdot (a_{t+1} - a_t) + rac{1}{2} \|a_t - a_{t+1}\|^2. \end{aligned}$$

Combining,

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

## Proof:

$$egin{aligned} &R(a_t) \geq R(a_{t+1}) + 
abla R(a_{t+1}) \cdot (a_t - a_{t+1}) + rac{1}{2} \|a_t - a_{t+1}\|^2, \ &R(a_{t+1}) \geq R(a_t) + 
abla R(a_t) \cdot (a_{t+1} - a_t) + rac{1}{2} \|a_t - a_{t+1}\|^2. \end{aligned}$$

Combining,

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

Hence,

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

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

This leads to the regret bound:

## Corollary

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

$$\sum_{t=1}^{n} (\ell_t(a_t) - \ell_t(a)) \leq \eta \sum_{t=1}^{n} \|g_t\|_*^2 + \frac{1}{\eta} (R(a) - R(a_1)).$$

This leads to the regret bound:

### Corollary

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

$$\sum_{t=1}^{n} (\ell_t(a_t) - \ell_t(a)) \leq \eta \sum_{t=1}^{n} \|g_t\|_*^2 + \frac{1}{\eta} (R(a) - R(a_1)).$$

Thus, for  $||g_t||_* \leq G$  and  $R(a) - R(a_1) \leq D^2$ , choosing  $\eta$  appropriately gives regret no more than  $2GD\sqrt{n}$ .

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

Consider  $R(a) = \frac{1}{2} ||a||^2$ ,  $a_1 = 0$ , and 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.

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  $2GD\sqrt{n}$ .

Consider  $\mathcal{A} = \Delta^m$ ,  $R(a) = \sum_i a_i \ln a_i$ .



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).

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,

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

・ロト ・四ト ・ヨト ・ヨト

78/132

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,

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

・ロン ・雪 と ・ ヨ と ・ ヨ と

78 / 132

And *R* is strongly convex wrt  $\|\cdot\|_1$ .

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,

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

・ロト ・日下・ ・日下・ ・日下

78 / 132

And *R* is strongly convex wrt  $\|\cdot\|_1$ . Also,  $R(a) - R(a_1) \leq \ln m$ .

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,

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

78 / 132

And R is strongly convex wrt  $\|\cdot\|_1$ . Also,  $R(a) - R(a_1) \le \ln m$ . Suppose that  $\|g_t\|_{\infty} \le 1$ .

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,

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

イロト 不得下 イヨト イヨト

78 / 132

And *R* is strongly convex wrt  $\|\cdot\|_1$ . Also,  $R(a) - R(a_1) \le \ln m$ . Suppose that  $\|g_t\|_{\infty} \le 1$ . Then the regret is no more than  $2\sqrt{n \ln m}$ .

### Example

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

$$a_{t+1} = \Pi^R_{\mathcal{A}}(\tilde{a}_{t+1})$$

$$egin{aligned} eta_{t+1} &= \Pi^R_\mathcal{A}( ilde{a}_{t+1}) \ &= \Pi^R_\mathcal{A}(
abla R^*(
abla R( ilde{a}_t) - \eta g_t)) \end{aligned}$$

$$egin{aligned} & \boldsymbol{\mu}_{t+1} = \Pi^R_{\mathcal{A}}(\tilde{a}_{t+1}) \ & = \Pi^R_{\mathcal{A}}(
abla R^*(
abla R(\tilde{a}_t) - \eta g_t)) \ & = \Pi^R_{\mathcal{A}}(
abla R^*(\ln(\tilde{a}_t \exp(-\eta g_t)))) \end{aligned}$$

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

 $\mathcal{A} = \Delta^m$ ,  $R(a) = \sum_i a_i \ln a_i$ . What are the updates?

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

where the In and exp functions are applied component-wise.

 $\mathcal{A} = \Delta^m$ ,  $R(a) = \sum_i a_i \ln a_i$ . What are the updates?

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

where the ln and exp functions are applied component-wise. This is exponentiated gradient: mirror descent with  $\nabla R = \ln$ .

 $\mathcal{A} = \Delta^m$ ,  $R(a) = \sum_i a_i \ln a_i$ . What are the updates?

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

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,  $\Pi^{R}_{\mathcal{A}}(\tilde{a}) = \tilde{a}/\|\tilde{a}\|_{1}$ . Notice that when the losses are linear, exponentiated gradient is exactly the exponential weights strategy we discussed for a finite comparison class. 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  $O(\sqrt{mn})$ .

Instead of

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

Instead of

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

we can use

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

◆□ → < □ → < 亘 → < 亘 → < 亘 → < 亘 → ○ Q (~ 81/132) Instead of

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

we can use

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

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

# Online convex optimization

## Binary prediction

- ② General online convex
  - Empirical minimization fails
  - Gradient algorithm
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
- Minimax strategies

#### Key Point:

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

# Some intuition about time-varying $\eta$ :

+

Consider

$$\Phi_t(a) = \sum_{s=1}^{s} \eta_s \ell_s(a) + R(a),$$

$$a_{t+1} = \arg\min_{a\in\mathbb{R}^d}\Phi_t(a).$$

<ロト < 部 > < 言 > < 言 > 言 \* の < 0 84/132

# Some intuition about time-varying $\eta$ :

Consider

$$\Phi_t(a) = \sum_{s=1}^t \eta_s \ell_s(a) + R(a), \qquad a_{t+1} = \arg\min_{a \in \mathbb{R}^d} \Phi_t(a).$$

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

$$\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) = \sum_{t=1}^n \frac{1}{\eta_t} \left( D_{\Phi_t}(a_t, a_{t+1}) + D_{\Phi_{t-1}}(a, a_t) - D_{\Phi_t}(a, a_{t+1}) \right)$$

## Some intuition about time-varying $\eta$ :

Consider

$$\Phi_t(a) = \sum_{s=1}^t \eta_s \ell_s(a) + R(a), \qquad a_{t+1} = \arg \min_{a \in \mathbb{R}^d} \Phi_t(a).$$

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

$$\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) = \sum_{t=1}^n \frac{1}{\eta_t} \left( D_{\Phi_t}(a_t, a_{t+1}) + D_{\Phi_{t-1}}(a, a_t) - D_{\Phi_t}(a, a_{t+1}) \right)$$

(Easy to check. Use  $\nabla \Phi_t(a_{t+1}) = \nabla \Phi_{t-1}(a_t) = 0.$ )

## Some intuition about time-varying $\eta$ :

Consider

$$\Phi_t(a) = \sum_{s=1}^t \eta_s \ell_s(a) + R(a), \qquad a_{t+1} = \arg \min_{a \in \mathbb{R}^d} \Phi_t(a).$$

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

$$\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) = \sum_{t=1}^n \frac{1}{\eta_t} \left( D_{\Phi_t}(a_t, a_{t+1}) + D_{\Phi_{t-1}}(a, a_t) - D_{\Phi_t}(a, a_{t+1}) \right)$$

(Easy to check. Use  $\nabla \Phi_t(a_{t+1}) = \nabla \Phi_{t-1}(a_t) = 0.$ ) What keeps the last two terms small?

#### Some intuition about time-varying $\eta$ :

Consider

$$\Phi_t(a) = \sum_{s=1}^{l} \eta_s \ell_s(a) + R(a), \qquad a_{t+1} = \arg \min_{a \in \mathbb{R}^d} \Phi_t(a).$$

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

$$\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) = \sum_{t=1}^n \frac{1}{\eta_t} \left( D_{\Phi_t}(a_t, a_{t+1}) + D_{\Phi_{t-1}}(a, a_t) - D_{\Phi_t}(a, a_{t+1}) \right)$$

(Easy to check. Use  $\nabla \Phi_t(a_{t+1}) = \nabla \Phi_{t-1}(a_t) = 0.$ ) What keeps the last two terms small? If we linearize the  $\ell_t$ , we have

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a) \leq \sum_{t=1}^{n} \frac{1}{\eta_t} \left( D_R(a_t, a_{t+1}) + D_R(a, a_t) - D_R(a, a_{t+1}) \right),$$

#### Some intuition about time-varying $\eta$ :

Consider

$$\Phi_t(a) = \sum_{s=1}^{l} \eta_s \ell_s(a) + R(a), \qquad a_{t+1} = \arg \min_{a \in \mathbb{R}^d} \Phi_t(a).$$

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

$$\sum_{t=1}^n (\ell_t(a_t) - \ell_t(a)) = \sum_{t=1}^n \frac{1}{\eta_t} \left( D_{\Phi_t}(a_t, a_{t+1}) + D_{\Phi_{t-1}}(a, a_t) - D_{\Phi_t}(a, a_{t+1}) \right)$$

(Easy to check. Use  $\nabla \Phi_t(a_{t+1}) = \nabla \Phi_{t-1}(a_t) = 0.$ ) What keeps the last two terms small? If we linearize the  $\ell_t$ , we have

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a) \leq \sum_{t=1}^{n} \frac{1}{\eta_t} \left( D_R(a_t, a_{t+1}) + D_R(a, a_t) - D_R(a, a_{t+1}) \right),$$

which requires  $\eta_t \approx \text{constant}$ . But what if  $\ell_t$  are strongly convex?

84 / 132

#### Theorem

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

$$\ell_t(a) \geq \ell_t(b) + 
abla \ell_t(b) \cdot (a-b) + rac{\sigma}{2} D_{\mathcal{R}}(a,b),$$

#### Theorem

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

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

and *R* is strongly convex wrt  $\|\cdot\|$ ,

#### Theorem

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

$$\ell_t(a) \geq \ell_t(b) + 
abla \ell_t(b) \cdot (a-b) + rac{\sigma}{2} D_{\mathcal{R}}(a,b),$$

and *R* is strongly convex wrt  $\|\cdot\|$ , then for any  $a \in \mathcal{A}$ , mirror descent,

$$a_{t+1} = \Pi_{\mathcal{A}}^{R}\left( (\nabla R)^{-1} \left( \nabla R(a_t) - \eta_t \nabla \ell_t(a_t) \right) \right)$$

with  $\eta_t \geq \frac{2}{t\sigma}$  has regret

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a) \leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) \leq \sum_{t=1}^{n} \eta_t \|\nabla \ell_t(a_t)\|_*^2.$$

[B., Hazan, Rakhlin, 2007] < □ → < 酉 → < 喜 → < 喜 → 85 / 132

$$\sum_{t=1}^n \left(\ell_t(a_t) - \ell_t(a)\right) \leq \sum_{t=1}^n \left(\nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t)\right).$$

#### Proof

$$\sum_{t=1}^n \left(\ell_t(a_t) - \ell_t(a)\right) \leq \sum_{t=1}^n \left(\nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t)\right).$$

Define:  $\tilde{a}_{t+1}$  so that  $a_{t+1} = \prod_{\mathcal{A}}^{R} (\tilde{a}_{t+1})$ :

#### Proof

$$\sum_{t=1}^n \left(\ell_t(a_t) - \ell_t(a)\right) \leq \sum_{t=1}^n \left(\nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t)\right).$$

Define:  $\tilde{a}_{t+1}$  so that  $a_{t+1} = \prod_{\mathcal{A}}^{R} (\tilde{a}_{t+1})$ :

$$\widetilde{a}_{t+1} := \nabla R^{-1} \left( \nabla R(a_t) - \eta_t \nabla \ell_t(a_t) \right),$$

#### Proof

$$\sum_{t=1}^n \left(\ell_t(a_t) - \ell_t(a)\right) \leq \sum_{t=1}^n \left(\nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t)\right).$$

Define:  $\tilde{a}_{t+1}$  so that  $a_{t+1} = \prod_{\mathcal{A}}^{R} (\tilde{a}_{t+1})$ :

$$\widetilde{a}_{t+1} := \nabla R^{-1} \left( \nabla R(a_t) - \eta_t \nabla \ell_t(a_t) \right),$$

and hence

$$\nabla R^{-1}(\tilde{a}_{t+1}) := \nabla R(a_t) - \eta_t \nabla \ell_t(a_t).$$

$$egin{aligned} \nabla \ell_t(a_t) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) 
ight) \cdot (a_t - a) \end{aligned}$$

#### Proof

$$egin{aligned} 
abla \ell_t(a_t) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) 
ight) \cdot (a_t - a) \end{aligned}$$

where the first equality follows from the definition of  $\tilde{a}_{t+1}$ ,

#### Proof

$$egin{aligned} 
abla \ell_t(a_t) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) 
ight) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, ilde{a}_{t+1}) + D_R(a_t, ilde{a}_{t+1}) 
ight) \end{aligned}$$

where the first equality follows from the definition of  $\tilde{a}_{t+1}$ ,

#### Proof

$$egin{aligned} 
abla \ell_t(a_t) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) 
ight) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, ilde{a}_{t+1}) + D_R(a_t, ilde{a}_{t+1}) 
ight) \end{aligned}$$

where the first equality follows from the definition of  $\tilde{a}_{t+1}$ , the second follows from the definition of Bregman divergence,

#### Proof

$$\begin{aligned} \nabla \ell_t(a_t) \cdot (a_t - a) \\ &= \frac{1}{\eta_t} \left( \nabla R(a_t) - \nabla R(\tilde{a}_{t+1}) \right) \cdot (a_t - a) \\ &= \frac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, \tilde{a}_{t+1}) + D_R(a_t, \tilde{a}_{t+1}) \right) \\ &\leq \frac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, a_{t+1}) + D_R(a_t, \tilde{a}_{t+1}) \right) \end{aligned}$$

where the first equality follows from the definition of  $\tilde{a}_{t+1}$ , the second follows from the definition of Bregman divergence,

#### Proof

$$egin{aligned} 
abla \ell_t(a_t) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) 
ight) \cdot (a_t - a) \ &= rac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, ilde{a}_{t+1}) + D_R(a_t, ilde{a}_{t+1}) 
ight) \ &\leq rac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, a_{t+1}) + D_R(a_t, ilde{a}_{t+1}) 
ight) \ \end{aligned}$$

where the first equality follows from the definition of  $\tilde{a}_{t+1}$ , the second follows from the definition of Bregman divergence, and the inequality follows from the Pythagorean Theorem for  $D_R$  (for  $a^* = \prod_{\mathcal{A}}^{\Phi}(b)$  and  $a \in \mathcal{A}$ ,  $D_{\Phi}(a, b) \ge D_{\Phi}(a, a^*) + D_{\Phi}(a^*, b)$ .)

$$\sum_{t=1}^{n} (\ell_t(a_t) - \ell_t(a))$$
  
$$\leq \sum_{t=1}^{n} \left( \nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t) \right)$$

$$\begin{split} &\sum_{t=1}^{n} \left( \ell_t(a_t) - \ell_t(a) \right) \\ &\leq \sum_{t=1}^{n} \left( \nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t) \right) \\ &\leq \sum_{t=1}^{n} \left( \frac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, a_{t+1}) + D_R(a_t, \tilde{a}_{t+1}) \right) - \frac{\sigma}{2} D_R(a, a_t) \right) \end{split}$$

$$\begin{split} &\sum_{t=1}^{n} \left( \ell_t(a_t) - \ell_t(a) \right) \\ &\leq \sum_{t=1}^{n} \left( \nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t) \right) \\ &\leq \sum_{t=1}^{n} \left( \frac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, a_{t+1}) + D_R(a_t, \tilde{a}_{t+1}) \right) - \frac{\sigma}{2} D_R(a, a_t) \right) \\ &= \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) + \sum_{t=2}^{n} \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - \frac{\sigma}{2} \right) D_R(a, a_t) \\ &+ \left( \frac{1}{\eta_1} - \frac{\sigma}{2} \right) D_R(a, a_1). \end{split}$$

#### Proof

$$\begin{split} &\sum_{t=1}^{n} \left( \ell_t(a_t) - \ell_t(a) \right) \\ &\leq \sum_{t=1}^{n} \left( \nabla \ell_t(a_t) \cdot (a_t - a) - \frac{\sigma}{2} D_R(a, a_t) \right) \\ &\leq \sum_{t=1}^{n} \left( \frac{1}{\eta_t} \left( D_R(a, a_t) - D_R(a, a_{t+1}) + D_R(a_t, \tilde{a}_{t+1}) \right) - \frac{\sigma}{2} D_R(a, a_t) \right) \\ &= \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) + \sum_{t=2}^{n} \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - \frac{\sigma}{2} \right) D_R(a, a_t) \\ &+ \left( \frac{1}{\eta_1} - \frac{\sigma}{2} \right) D_R(a, a_1). \end{split}$$

And choosing  $\eta_t = c/t$  for  $c \ge 2/\sigma$  eliminates the second and third terms.

# Proof Also, $D_R(a_t, \tilde{a}_{t+1}) \leq D_R(a_t, \tilde{a}_{t+1}) + D_R(\tilde{a}_{t+1}, a_t)$

#### Proof

Also,

$$egin{aligned} D_R(a_t, ilde{a}_{t+1}) &\leq D_R(a_t, ilde{a}_{t+1}) + D_R( ilde{a}_{t+1}, a_t) \ &= (
abla R(a_t) - 
abla R( ilde{a}_{t+1})) \cdot (a_t - ilde{a}_{t+1}) \end{aligned}$$

#### Proof

Also,

$$egin{aligned} \mathcal{D}_R(a_t, ilde{a}_{t+1}) &\leq \mathcal{D}_R(a_t, ilde{a}_{t+1}) + \mathcal{D}_R( ilde{a}_{t+1}, a_t) \ &= (
abla R(a_t) - 
abla R( ilde{a}_{t+1})) \cdot (a_t - ilde{a}_{t+1}) \ &= \eta_t 
abla \ell_t (a_t) \cdot (a_t - ilde{a}_{t+1}) \end{aligned}$$

#### Proof

Also,

$$egin{aligned} \mathcal{D}_R(a_t, ilde{a}_{t+1}) &\leq \mathcal{D}_R(a_t, ilde{a}_{t+1}) + \mathcal{D}_R( ilde{a}_{t+1}, a_t) \ &= (
abla R(a_t) - 
abla R( ilde{a}_{t+1})) \cdot (a_t - ilde{a}_{t+1}) \ &= \eta_t 
abla \ell_t (a_t) \cdot (a_t - ilde{a}_{t+1}) \end{aligned}$$

where the second equality is from the definition of  $\tilde{a}_{t+1}$ 

Also,

$$egin{aligned} D_R(a_t, ilde{a}_{t+1}) &\leq D_R(a_t, ilde{a}_{t+1}) + D_R( ilde{a}_{t+1}, a_t) \ &= (
abla R(a_t) - 
abla R( ilde{a}_{t+1})) \cdot (a_t - ilde{a}_{t+1}) \ &= \eta_t 
abla \ell_t(a_t) \cdot (a_t - ilde{a}_{t+1}) \ &\leq \eta_t \| 
abla \ell_t(a_t) \|_* \| a_t - ilde{a}_{t+1} \| \end{aligned}$$

where the second equality is from the definition of  $\tilde{a}_{t+1}$ 

Also,

$$egin{aligned} D_R(a_t, ilde{a}_{t+1}) &\leq D_R(a_t, ilde{a}_{t+1}) + D_R( ilde{a}_{t+1}, a_t) \ &= (
abla R(a_t) - 
abla R( ilde{a}_{t+1})) \cdot (a_t - ilde{a}_{t+1}) \ &= \eta_t 
abla \ell_t(a_t) \cdot (a_t - ilde{a}_{t+1}) \ &\leq \eta_t \| 
abla \ell_t(a_t) \|_* \| a_t - ilde{a}_{t+1} \| \ &\leq \eta_t \| 
abla \ell_t(a_t) \|_* \| 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) \| \end{aligned}$$

where the second equality is from the definition of  $\tilde{a}_{t+1}$ 

Also,

$$egin{aligned} D_R(a_t, ilde{a}_{t+1}) &\leq D_R(a_t, ilde{a}_{t+1}) + D_R( ilde{a}_{t+1}, a_t) \ &= (
abla R(a_t) - 
abla R( ilde{a}_{t+1})) \cdot (a_t - ilde{a}_{t+1}) \ &= \eta_t 
abla \ell_t(a_t) \cdot (a_t - ilde{a}_{t+1}) \ &\leq \eta_t \| 
abla \ell_t(a_t) \|_* \| a_t - ilde{a}_{t+1} \| \ &\leq \eta_t \| 
abla \ell_t(a_t) \|_* \| 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) \| \end{aligned}$$

where the second equality is from the definition of  $\tilde{a}_{t+1}$ and the second inequality follows from the strong convexity of R wrt  $\|\cdot\|$ .

Also,

$$egin{aligned} D_R(a_t, ilde{a}_{t+1}) &\leq D_R(a_t, ilde{a}_{t+1}) + D_R( ilde{a}_{t+1}, a_t) \ &= (
abla R(a_t) - 
abla R( ilde{a}_{t+1})) \cdot (a_t - ilde{a}_{t+1}) \ &= \eta_t 
abla \ell_t(a_t) \cdot (a_t - ilde{a}_{t+1}) \ &\leq \eta_t \| 
abla \ell_t(a_t) \|_* \| a_t - ilde{a}_{t+1} \| \ &\leq \eta_t \| 
abla \ell_t(a_t) \|_* \| 
abla R(a_t) - 
abla R( ilde{a}_{t+1}) \|_* \ &= \eta_t^2 \| 
abla \ell_t(a_t) \|_*, \end{aligned}$$

where the second equality is from the definition of  $\tilde{a}_{t+1}$ and the second inequality follows from the strong convexity of R wrt  $\|\cdot\|$ .

#### Theorem

If  $\ell_t$  is  $\sigma$ -strongly convex wrt R and R is strongly convex wrt  $\|\cdot\|$ , then for any  $a \in \mathcal{A}$ , mirror descent,  $a_{t+1} = \prod_{\mathcal{A}}^{R} \left( (\nabla R)^{-1} \left( \nabla R(a_t) - \eta_t \nabla \ell_t(a_t) \right) \right)$ with  $\eta_t \geq \frac{2}{t\sigma}$  has regret

$$\sum_{t=1}^n \ell_t(\mathsf{a}_t) - \sum_{t=1}^n \ell_t(\mathsf{a}) \leq \sum_{t=1}^n \eta_t \| 
abla \ell_t(\mathsf{a}_t) \|_*^2.$$

#### Theorem

If  $\ell_t$  is  $\sigma$ -strongly convex wrt R and R is strongly convex wrt  $\|\cdot\|$ , then for any  $a \in A$ , mirror descent,  $a_{t+1} = \prod_{\mathcal{A}}^{R} \left( (\nabla R)^{-1} \left( \nabla R(a_t) - \eta_t \nabla \ell_t(a_t) \right) \right)$ with  $\eta_t \geq \frac{2}{t\sigma}$  has regret

$$\sum_{t=1}^n \ell_t(\mathsf{a}_t) - \sum_{t=1}^n \ell_t(\mathsf{a}) \leq \sum_{t=1}^n \eta_t \| 
abla \ell_t(\mathsf{a}_t) \|_*^2.$$

#### Example

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

$$\sum_{t=1}^n \ell_t(a_t) - \inf_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) \le \sum_{t=1}^n \eta_t \|\nabla \ell_t\|_*^2 = O\left(\frac{G^2}{\sigma} \log n\right).$$

$$\sum_{t=1}^{n} (\ell_t(a_t) - \ell_t(a))$$
  
$$\leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) + \sum_{t=2}^{n} \left(\frac{\sqrt{t}}{c} - \frac{\sqrt{t-1}}{c}\right) D_R(a, a_t) + \frac{1}{c} D_R(a, a_1)$$

$$\begin{split} &\sum_{t=1}^{n} \left( \ell_t(a_t) - \ell_t(a) \right) \\ &\leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) + \sum_{t=2}^{n} \left( \frac{\sqrt{t}}{c} - \frac{\sqrt{t-1}}{c} \right) D_R(a, a_t) + \frac{1}{c} D_R(a, a_1) \\ &\leq \sum_{t=1}^{n} \eta_t \| \nabla \ell_t(a_t) \|_*^2 + \frac{D^2}{c} \sum_{t=1}^{n} \left( \sqrt{t} - \sqrt{t-1} \right) \end{split}$$

$$\begin{split} &\sum_{t=1}^{n} \left( \ell_t(a_t) - \ell_t(a) \right) \\ &\leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) + \sum_{t=2}^{n} \left( \frac{\sqrt{t}}{c} - \frac{\sqrt{t-1}}{c} \right) D_R(a, a_t) + \frac{1}{c} D_R(a, a_1) \\ &\leq \sum_{t=1}^{n} \eta_t \| \nabla \ell_t(a_t) \|_*^2 + \frac{D^2}{c} \sum_{t=1}^{n} \left( \sqrt{t} - \sqrt{t-1} \right) \\ &\leq \left( cG^2 + \frac{D^2}{c} \right) \sqrt{n} \end{split}$$

$$\begin{split} &\sum_{t=1}^{n} \left( \ell_t(a_t) - \ell_t(a) \right) \\ &\leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) + \sum_{t=2}^{n} \left( \frac{\sqrt{t}}{c} - \frac{\sqrt{t-1}}{c} \right) D_R(a, a_t) + \frac{1}{c} D_R(a, a_1) \\ &\leq \sum_{t=1}^{n} \eta_t \| \nabla \ell_t(a_t) \|_*^2 + \frac{D^2}{c} \sum_{t=1}^{n} \left( \sqrt{t} - \sqrt{t-1} \right) \\ &\leq \left( cG^2 + \frac{D^2}{c} \right) \sqrt{n} \\ &= O(DG\sqrt{n}). \end{split}$$

## Regularization methods: Convexity and Strong Convexity

$\ell_t$	$\eta_t$	R <sub>n</sub>
convex	$\frac{1}{\sqrt{t}}$	$O(\sqrt{n})$
$\sigma$ -strongly convex	$\frac{1}{\sigma t}$	$O\left(\frac{1}{\sigma}\log n\right)$

# Regularization methods: Convexity and Strong Convexity

t	$\eta_t$	l <sub>t</sub>
	$\frac{1}{\sqrt{t}}$	convex
$\frac{1}{\sigma} \log \left(\frac{1}{\sigma} \log \right)$	$\frac{1}{\sigma t}$	$\sigma$ -strongly convex

All that changes is the step-size.

# Regularization methods: Convexity and Strong Convexity

	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	D
ℓt	$\eta_t$	<u> </u>
convex	$\frac{1}{\sqrt{t}}$	$O(\sqrt{n})$
$\sigma$ -strongly convex	$\frac{1}{\sigma t}$	$O\left(\frac{1}{\sigma}\log n\right)$

イロト イポト イヨト イヨト

92 / 132

All that changes is the step-size.

What if we don't know  $\sigma$ ?

Can we adapt our step-size to give the right rate?

# Outline

## Binary prediction

- eneral online convex
  - Empirical minimization fails
  - Gradient algorithm
  - A regularization viewpoint
  - Bregman divergence
  - Properties of regularization
  - Linearization
  - Mirror descent
  - Regret bounds
  - Strongly convex losses
  - Adaptive regularization
    - Strong convexity (Adaptive Gradient)
    - Diagonal regularizers (AdaGrad)
- Olinimax strategies

## Adaptive regularization

## Adaptive regularization

#### Adaptive regularization

$$R_n = \sum_{t=1}^n \left( \ell_t(a_t) - \ell_t(a) \right)$$

$$\begin{aligned} R_n &= \sum_{t=1}^n \left( \ell_t(a_t) - \ell_t(a) \right) \\ &= \sum_{t=1}^n \left( \tilde{\ell}_t(a_t) - \tilde{\ell}_t(a) + \lambda_t(g(a) - g(a_t)) \right) \end{aligned}$$

$$egin{aligned} &R_n = \sum_{t=1}^n \left(\ell_t(a_t) - \ell_t(a)
ight) \ &= \sum_{t=1}^n \left( ilde{\ell}_t(a_t) - ilde{\ell}_t(a) + \lambda_t(g(a) - g(a_t))
ight) \ &\leq D^2 \sum_{t=1}^n \lambda_t + \sum_{t=1}^n \left( ilde{\ell}_t(a_t) - ilde{\ell}_t(a)
ight), \end{aligned}$$

Replace  $\ell_t(\cdot)$  with  $\tilde{\ell}_t(\cdot) := \ell_t(\cdot) + \lambda_t g(\cdot)$ , where g is strongly convex wrt R.

$$\begin{split} R_n &= \sum_{t=1}^n \left( \ell_t(a_t) - \ell_t(a) \right) \\ &= \sum_{t=1}^n \left( \tilde{\ell}_t(a_t) - \tilde{\ell}_t(a) + \lambda_t(g(a) - g(a_t)) \right) \\ &\leq D^2 \sum_{t=1}^n \lambda_t + \sum_{t=1}^n \left( \tilde{\ell}_t(a_t) - \tilde{\ell}_t(a) \right), \end{split}$$

where we've defined  $D^2 := \sup_{a,a_t} (g(a) - g(a_t)).$ 

Replace  $\ell_t(\cdot)$  with  $\tilde{\ell}_t(\cdot) := \ell_t(\cdot) + \lambda_t g(\cdot)$ , where g is strongly convex wrt R.

$$egin{aligned} &R_n = \sum_{t=1}^n \left(\ell_t(a_t) - \ell_t(a)
ight) \ &= \sum_{t=1}^n \left( ilde{\ell}_t(a_t) - ilde{\ell}_t(a) + \lambda_t(g(a) - g(a_t))
ight) \ &\leq D^2 \sum_{t=1}^n \lambda_t + \sum_{t=1}^n \left( ilde{\ell}_t(a_t) - ilde{\ell}_t(a)
ight), \end{aligned}$$

where we've defined  $D^2 := \sup_{a,a_t} (g(a) - g(a_t))$ . This is an approximation error term, plus the regret for the regularized losses  $\tilde{\ell}_t$ .

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \ldots, \lambda_n).$$

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \dots, \lambda_n).$$

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \dots, \lambda_n).$$

This is similar to a model selection problem.

• How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s?

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \ldots, \lambda_n).$$

This is similar to a model selection problem.

**1** How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s? (We'll give a bound.)

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \dots, \lambda_n).$$

- How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s? (We'll give a bound.)
- Ooes the best trade-off between the two terms above ensure the optimal rates for convex and strongly convex l<sub>t</sub>?

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \ldots, \lambda_n).$$

- How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s? (We'll give a bound.)
- Ooes the best trade-off between the two terms above ensure the optimal rates for convex and strongly convex l<sub>t</sub>? (Yes.)

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \ldots, \lambda_n).$$

- **1** How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s? (We'll give a bound.)
- ② Does the best trade-off between the two terms above ensure the optimal rates for convex and strongly convex  $\ell_t$ ? (Yes.)
- Or an we choose \(\lambda\_t\) online to obtain the best trade-off between these two terms?

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \ldots, \lambda_n).$$

- How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s? (We'll give a bound.)
- ② Does the best trade-off between the two terms above ensure the optimal rates for convex and strongly convex  $\ell_t$ ? (Yes.)
- Can we choose λ<sub>t</sub> online to obtain the best trade-off between these two terms? (Yes.)

#### Theorem

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

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

and *R* is strongly convex wrt  $\|\cdot\|$ ,

see, e.g., [B., Hazan, Rakhlin, 2007]

96 / 132

イロト イポト イヨト イヨト

#### Theorem

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

$$\ell_t(a) \geq \ell_t(b) + 
abla \ell_t(b) \cdot (a-b) + rac{\sigma_t}{2} D_R(a,b),$$

and *R* is strongly convex wrt  $\|\cdot\|$ , then for any  $a \in \mathbb{R}^d$ , mirror descent with  $\eta_t = 2/\sum_{s=1}^t \sigma_s$  has regret

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a) \leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) \leq 2 \sum_{t=1}^{n} \frac{\|\nabla \ell_t(a_t)\|_*^2}{\sum_{s=1}^t \sigma_s}.$$

see, e.g., [B., Hazan, Rakhlin, 2007]

96 / 132

#### Theorem

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

$$\ell_t(a) \geq \ell_t(b) + 
abla \ell_t(b) \cdot (a-b) + rac{\sigma_t}{2} D_R(a,b),$$

and R is strongly convex wrt  $\|\cdot\|$ , then for any  $a \in \mathbb{R}^d$ , mirror descent with  $\eta_t = 2/\sum_{s=1}^t \sigma_s$  has regret

$$\sum_{t=1}^{n} \ell_t(a_t) - \sum_{t=1}^{n} \ell_t(a) \leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) \leq 2 \sum_{t=1}^{n} \frac{\|\nabla \ell_t(a_t)\|_*^2}{\sum_{s=1}^{t} \sigma_s}.$$

Notice:  $\eta_t$  is used to update  $a_t$  to  $a_{t+1}$ , so it uses only past information.

see, e.g., [B., Hazan, Rakhlin, 2007]

## Proof idea

As before (when  $\sigma_t$  was constant), we have

$$egin{aligned} &\sum_{t=1}^n \left(\ell_t(a_t) - \ell_t(a)
ight) \ &\leq \sum_{t=1}^n rac{1}{\eta_t} D_R(a_t, ilde{a}_{t+1}) + \sum_{t=2}^n \left(rac{1}{\eta_t} - rac{1}{\eta_{t-1}} - rac{\sigma_t}{2}
ight) D_R(a, a_t) \ &+ \left(rac{1}{\eta_1} - rac{\sigma_1}{2}
ight) D_R(a, a_1). \end{aligned}$$

・ロ ・ < 回 ト < 三 ト < 三 ト ミ の Q ()
97 / 132

#### Proof idea

As before (when  $\sigma_t$  was constant), we have

$$\begin{split} &\sum_{t=1}^{n} \left( \ell_t(a_t) - \ell_t(a) \right) \\ &\leq \sum_{t=1}^{n} \frac{1}{\eta_t} D_R(a_t, \tilde{a}_{t+1}) + \sum_{t=2}^{n} \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - \frac{\sigma_t}{2} \right) D_R(a, a_t) \\ &\quad + \left( \frac{1}{\eta_1} - \frac{\sigma_1}{2} \right) D_R(a, a_1). \end{split}$$

And the choice of  $\eta_t$  eliminates the second and third terms.

Work with  $\tilde{\ell}_t(\cdot) := \ell_t(\cdot) + \lambda_t g(\cdot)$  (where g is strongly convex wrt R). If the  $\ell_t$  are  $\sigma_t$ -strongly convex wrt R, then  $\tilde{\ell}_t$  are  $(\sigma_t + \lambda_t)$ -strongly convex.

Work with  $\tilde{\ell}_t(\cdot) := \ell_t(\cdot) + \lambda_t g(\cdot)$  (where g is strongly convex wrt R). If the  $\ell_t$  are  $\sigma_t$ -strongly convex wrt R, then  $\tilde{\ell}_t$  are  $(\sigma_t + \lambda_t)$ -strongly convex. Using mirror descent for the  $\tilde{\ell}_t$ s, we choose steps

$$\eta_t = \frac{2}{\sum_{s=1}^t (\sigma_s + \lambda_s)}.$$

(日) (同) (日) (日)

98 / 132

This strategy incurs regret

$$\sum_{t=1}^n \ell_t(a_t) - \inf_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) \le D^2 \sum_{t=1}^n \lambda_t + 2 \sum_{t=1}^n \left( \tilde{\ell}_t(a_t) - \tilde{\ell}_t(a) \right)$$

This strategy incurs regret

$$\sum_{t=1}^{n} \ell_t(a_t) - \inf_{a \in \mathbb{R}^d} \sum_{t=1}^{n} \ell_t(a) \le D^2 \sum_{t=1}^{n} \lambda_t + 2 \sum_{t=1}^{n} \left( \tilde{\ell}_t(a_t) - \tilde{\ell}_t(a) \right)$$
$$\le D^2 \sum_{t=1}^{n} \lambda_t + 2 \sum_{t=1}^{n} \frac{\|\nabla \tilde{\ell}_t(a_t)\|_*^2}{\sum_{s=1}^{t} (\sigma_s + \lambda_s)}$$

This strategy incurs regret

$$\begin{split} \sum_{t=1}^n \ell_t(a_t) &- \inf_{a \in \mathbb{R}^d} \sum_{t=1}^n \ell_t(a) \le D^2 \sum_{t=1}^n \lambda_t + 2 \sum_{t=1}^n \left( \tilde{\ell}_t(a_t) - \tilde{\ell}_t(a) \right) \\ &\le D^2 \sum_{t=1}^n \lambda_t + 2 \sum_{t=1}^n \frac{\|\nabla \tilde{\ell}_t(a_t)\|_*^2}{\sum_{s=1}^t (\sigma_s + \lambda_s)} \\ &\le D^2 \sum_{t=1}^n \lambda_t + 2 \sum_{t=1}^n \frac{(G_t + \lambda_t B)^2}{\sum_{s=1}^t (\sigma_s + \lambda_s)}, \end{split}$$

This strategy incurs regret

$$\sum_{t=1}^{n} \ell_t(a_t) - \inf_{a \in \mathbb{R}^d} \sum_{t=1}^{n} \ell_t(a) \le D^2 \sum_{t=1}^{n} \lambda_t + 2 \sum_{t=1}^{n} \left( \tilde{\ell}_t(a_t) - \tilde{\ell}_t(a) \right)$$
$$\le D^2 \sum_{t=1}^{n} \lambda_t + 2 \sum_{t=1}^{n} \frac{\|\nabla \tilde{\ell}_t(a_t)\|_*^2}{\sum_{s=1}^{t} (\sigma_s + \lambda_s)}$$
$$\le D^2 \sum_{t=1}^{n} \lambda_t + 2 \sum_{t=1}^{n} \frac{(G_t + \lambda_t B)^2}{\sum_{s=1}^{t} (\sigma_s + \lambda_s)},$$

where  $\|\nabla \ell_t(a_t)\|_* \leq G_t$  and  $\|\nabla g(a_t)\|_* \leq B$ .

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \ldots, \lambda_n).$$

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \ldots, \lambda_n).$$

• How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s?

Obes the best trade-off between the two terms above ensure the optimal rates for convex and strongly convex l<sub>t</sub>?

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \dots, \lambda_n).$$

- How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s?
- Does the best trade-off between the two terms above ensure the optimal rates for convex and strongly convex l<sub>t</sub>?
- Solution Can we choose λ<sub>t</sub> online to obtain the best trade-off between these two terms?

And the best choice of  $\lambda_1, \ldots, \lambda_n$  is good here in the convex case:

#### Example

Assume  $\sigma_t \geq 0$ . Choose

$$\lambda_1 = \sqrt{\frac{\sum_{t=1}^n G_t^2}{B^2 + D^2}}$$

and  $\lambda_2 = \cdots = \lambda_n = 0$ .

And the best choice of  $\lambda_1, \ldots, \lambda_n$  is good here in the convex case:

#### Example

Assume  $\sigma_t \geq 0$ . Choose

$$\lambda_1 = \sqrt{\frac{\sum_{t=1}^n G_t^2}{B^2 + D^2}}$$

and  $\lambda_2 = \cdots = \lambda_n = 0$ . Then the bound gives

$$egin{aligned} &R_n \leq D^2 \sum_{t=1}^n \lambda_t + 2 \sum_{t=1}^n rac{(G_t + \lambda_t B)^2}{\sum_{s=1}^t (\sigma_s + \lambda_s)^2} \ &= O\left(\sqrt{(B^2 + D^2) \sum_{t=1}^n G_t^2}
ight). \end{aligned}$$

And the best choice of  $\lambda_1, \ldots, \lambda_n$  is good here in the convex case:

#### Example

Assume  $\sigma_t \geq 0$ . Choose

$$\lambda_1 = \sqrt{\frac{\sum_{t=1}^n G_t^2}{B^2 + D^2}}$$

and  $\lambda_2 = \cdots = \lambda_n = 0$ . Then the bound gives

$$egin{aligned} &R_n \leq D^2 \sum_{t=1}^n \lambda_t + 2 \sum_{t=1}^n rac{(G_t + \lambda_t B)^2}{\sum_{s=1}^t (\sigma_s + \lambda_s)^2} \ &= O\left(\sqrt{(B^2 + D^2) \sum_{t=1}^n G_t^2}
ight). \end{aligned}$$

If  $G_t \leq G$ , this is  $R_n = O\left(\sqrt{B^2 + D^2}G\sqrt{n}\right)$ .

101/132

And the best choice of  $\lambda_1, \ldots, \lambda_n$  is good here in the strongly convex case:

#### Example

Assume  $\sigma_t \geq \sigma$  and  $G_t \leq G$ . Choose  $\lambda_1 = \cdots = \lambda_n = 0$ .

And the best choice of  $\lambda_1, \ldots, \lambda_n$  is good here in the strongly convex case:

#### Example

Assume  $\sigma_t \geq \sigma$  and  $G_t \leq G$ . Choose  $\lambda_1 = \cdots = \lambda_n = 0$ . Then the bound gives

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + 2 \sum_{t=1}^n \frac{(G_t + \lambda_t B)^2}{\sum_{s=1}^t (\sigma_s + \lambda_s)}$$
$$= O\left(\frac{G^2}{\sigma} \log n\right).$$

We can also obtain a spectrum of rates with the best choice of  $\lambda_1, \ldots, \lambda_n$ :

# Example Suppose $\sigma_t = t^{-\alpha}$ and $G_t \leq G$ . Then the bound gives $R_n = \begin{cases} O(\log n) & \text{if } \alpha = 0, \\\\ O(\sqrt{n}) & \text{if } \alpha > 1/2. \end{cases}$

We can also obtain a spectrum of rates with the best choice of  $\lambda_1, \ldots, \lambda_n$ :

# Example Suppose $\sigma_t = t^{-\alpha}$ and $G_t \leq G$ . Then the bound gives $R_n = \begin{cases} O(\log n) & \text{if } \alpha = 0, \\ O(n^{\alpha}) & \text{if } 0 < \alpha \leq 1/2, \\ O(\sqrt{n}) & \text{if } \alpha > 1/2. \end{cases}$

We can also obtain a spectrum of rates with the best choice of  $\lambda_1, \ldots, \lambda_n$ :

#### Example

Suppose  $\sigma_t = t^{-\alpha}$  and  $G_t \leq G$ . Then the bound gives

$$R_n = \begin{cases} O(\log n) & \text{if } \alpha = 0, \\ O(n^{\alpha}) & \text{if } 0 < \alpha \le 1/2, \\ O(\sqrt{n}) & \text{if } \alpha > 1/2. \end{cases}$$

(Choose  $\lambda_1 = n^{\alpha}$  and  $\lambda_2 = \cdots \lambda_n = 0$ .)

$$R_n \leq D^2 \sum_{t=1}^n \lambda_t + \tilde{R}_n(\lambda_1, \dots, \lambda_n).$$

• How does  $\tilde{R}_n$  depend on the  $\lambda_t$ s?

- 2 Does the best trade-off between the two terms above ensure the optimal rates for convex and strongly convex l<sub>t</sub>?
- Solution Can we choose λ<sub>t</sub> online to obtain the best trade-off between these two terms?

# Regularization methods: adapting to strong convexity

#### Theorem

Choosing

$$\lambda_t = \frac{1}{2} \left( \sqrt{\left( \sum_{s=1}^{t-1} (\sigma_s + \lambda_s) + \sigma_t \right)^2 + \frac{16G_t^2}{D^2 + B^2}} - \left( \sum_{s=1}^{t-1} (\sigma_s + \lambda_s) + \sigma_t \right) \right)$$

with this regularized mirror descent strategy

[B., Hazan, Rakhlin, 2007]

#### Theorem

Choosing

$$\lambda_t = \frac{1}{2} \left( \sqrt{\left( \sum_{s=1}^{t-1} (\sigma_s + \lambda_s) + \sigma_t \right)^2 + \frac{16G_t^2}{D^2 + B^2}} - \left( \sum_{s=1}^{t-1} (\sigma_s + \lambda_s) + \sigma_t \right) \right)$$

with this regularized mirror descent strategy gives regret

$$R_n = O\left(\inf_{\lambda_1,\dots,\lambda_n}\left((D^2 + B^2)\sum_{t=1}^n \lambda_t + \sum_{t=1}^n \frac{(G_t + \lambda_t B)^2}{\sum_{s=1}^t (\sigma_s + \lambda_s)}\right)\right).$$

[B., Hazan, Rakhlin, 2007]

105 / 132

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

# Regularization methods: adapting to strong convexity

• Notice that we're using information about each  $\ell_t$  only after we see it.

## Regularization methods: adapting to strong convexity

- Notice that we're using information about each  $\ell_t$  only after we see it.
- Compare this to the simple gradient method that we saw earlier, which chooses  $\eta = D/(G\sqrt{n})$ . Here, we don't need to know the upper bound G (or n): we choose  $\lambda_t$  as a function of information about past losses, and we can compete with the optimal bounds.

- Notice that we're using information about each  $\ell_t$  only after we see it.
- Compare this to the simple gradient method that we saw earlier, which chooses  $\eta = D/(G\sqrt{n})$ . Here, we don't need to know the upper bound G (or n): we choose  $\lambda_t$  as a function of information about past losses, and we can compete with the optimal bounds.
- For instance, for the case of convex functions that satisfy a gradient dual norm bound *G*,

$$R_n = O\left(\sqrt{B^2 + D^2}G\sqrt{n}\right).$$

(And similarly for the stronger version that replaces G by the rms dual norm of the gradients.)

#### Proof Idea

We prove that balancing the two terms is near-optimal: Consider

$$H_n(\{\lambda_t\}) := \sum_{t=1}^n \lambda_t + \sum_{t=1}^n \frac{C_t}{\sum_{s=1}^t (\sigma_s + \lambda_s)}$$

Then choosing  $\lambda_t$  to solve the quadratic equation

$$\lambda_t = \frac{C_t}{\sum_{s=1}^t (\sigma_s + \lambda_s)}$$

ensures that

$$H_n(\{\lambda_t\}) \leq 2 \inf_{\{\lambda_t^*\}} H_n(\{\lambda_t^*\}).$$

#### Proof Idea

There is an inductive proof of this balancing result, which considers separately the cases

$$\sum_{s=1}^t \lambda_s < \sum_{s=1}^t \lambda_s^*$$

and

and exploits the fact that the two terms of  $H_t$  are monotonic in  $\sum_{s=1}^t \lambda_s$ . And the choice of  $\lambda_t$  in the theorem is the positive solution to the appropriate quadratic equation.

 $\sum_{s=1}^{t} \lambda_s > \sum_{s=1}^{t} \lambda_s^*,$ 

#### Theorem

#### Choosing

$$\lambda_t = \frac{1}{2} \left( \sqrt{\left(\sum_{s=1}^{t-1} (\sigma_s + \lambda_s) + \sigma_t\right)^2 + \frac{16G_t^2}{D^2 + B^2}} - \left(\sum_{s=1}^{t-1} (\sigma_s + \lambda_s) + \sigma_t\right) \right)$$

with this regularized mirror descent strategy gives regret

$$R_n = O\left(\inf_{\lambda_1,\dots,\lambda_n}\left((D^2 + B^2)\sum_{t=1}^n \lambda_t + \sum_{t=1}^n \frac{(G_t + \lambda_t B)^2}{\sum_{s=1}^t (\sigma_s + \lambda_s)}\right)\right).$$

◆□ → ◆□ → ◆ 三 → ◆ 三 → ○ へ ○
109/132

# Outline

#### Binary prediction

#### eneral online convex

- Empirical minimization fails
- Gradient algorithm
- A regularization viewpoint
- Bregman divergence
- Properties of regularization
- Linearization
- Mirror descent
- Regret bounds
- Strongly convex losses
- Adaptive regularization
  - Strong convexity (Adaptive Gradient)
  - Diagonal regularizers (AdaGrad)
- Minimax strategies

We considered mirror descent where we added an adaptively chosen component of a regularizer g that is strongly convex wrt R. To simplify, assume g = R.

We considered mirror descent where we added an adaptively chosen component of a regularizer g that is strongly convex wrt R. To simplify, assume g = R. We can view this in two ways:

$$a_{t+1} = \arg\min_{a \in \mathcal{A}} \left( \sum_{s=1}^{t} \eta_s \nabla(\ell_s + \lambda_s R)(a_s) \cdot (a - a_t) + R(a) \right)$$

We considered mirror descent where we added an adaptively chosen component of a regularizer g that is strongly convex wrt R. To simplify, assume g = R. We can view this in two ways:

$$a_{t+1} = \arg\min_{a \in \mathcal{A}} \left( \sum_{s=1}^{t} \eta_s \nabla(\ell_s + \lambda_s R)(a_s) \cdot (a - a_t) + R(a) \right)$$
  
=  $\arg\min_{a \in \mathcal{A}} \left( \eta_t \nabla(\ell_t + \lambda_t R)(a_t) \cdot (a - a_t) + D_R(a, \tilde{a}_t) \right).$ 

We considered mirror descent where we added an adaptively chosen component of a regularizer g that is strongly convex wrt R. To simplify, assume g = R. We can view this in two ways:

$$a_{t+1} = \arg\min_{a \in \mathcal{A}} \left( \sum_{s=1}^{t} \eta_s \nabla(\ell_s + \lambda_s R)(a_s) \cdot (a - a_t) + R(a) \right)$$
  
=  $\arg\min_{a \in \mathcal{A}} \left( \eta_t \nabla(\ell_t + \lambda_t R)(a_t) \cdot (a - a_t) + D_R(a, \tilde{a}_t) \right).$ 

Rather than minimizing the sum of the linearization of  $\ell_t + \lambda_t R$  plus the regularizer R, we could instead minimize the linearization of  $\ell_t$  plus the regularizer  $(1 + \lambda_t)R$ :

$$a_{t+1} = \arg\min_{a \in \mathcal{A}} \left( \eta_t \nabla \ell_t(a_t) \cdot (a - a_t) + D_{(1+\lambda_t)R}(a, \tilde{a}_t) \right).$$

• Adaptive regularization:  $R_t(a) = (1 + \lambda_t)R(a)$ .

- Adaptive regularization:  $R_t(a) = (1 + \lambda_t)R(a)$ .
- We could be more ambitious, and consider more than a single parameter (λ<sub>t</sub>).

- Adaptive regularization:  $R_t(a) = (1 + \lambda_t)R(a)$ .
- We could be more ambitious, and consider more than a single parameter  $(\lambda_t)$ .
- For example, generalizing the gradient case (where  $R(a) = ||a||_2^2$ ), we could consider

$$R_t(a) = a^\top M_t a,$$

- Adaptive regularization:  $R_t(a) = (1 + \lambda_t)R(a)$ .
- We could be more ambitious, and consider more than a single parameter  $(\lambda_t)$ .
- For example, generalizing the gradient case (where  $R(a) = ||a||_2^2$ ), we could consider

$$R_t(a) = a^{\top} M_t a,$$

• with 
$$M_t = (1+\lambda_t) I$$
 (as before),

- Adaptive regularization:  $R_t(a) = (1 + \lambda_t)R(a)$ .
- We could be more ambitious, and consider more than a single parameter  $(\lambda_t)$ .
- For example, generalizing the gradient case (where  $R(a) = ||a||_2^2$ ), we could consider

$$R_t(a) = a^{ op} M_t a,$$

- with  $M_t = (1 + \lambda_t) I$  (as before),
- with  $M_t$  a positive diagonal matrix, or

- Adaptive regularization:  $R_t(a) = (1 + \lambda_t)R(a)$ .
- We could be more ambitious, and consider more than a single parameter  $(\lambda_t)$ .
- For example, generalizing the gradient case (where  $R(a) = ||a||_2^2$ ), we could consider

$$R_t(a) = a^{\top} M_t a,$$

- with  $M_t = (1 + \lambda_t) I$  (as before),
- with  $M_t$  a positive diagonal matrix, or
- with  $M_t \succ 0$  (an arbitrary positive definite matrix).

- Adaptive regularization:  $R_t(a) = (1 + \lambda_t)R(a)$ .
- We could be more ambitious, and consider more than a single parameter  $(\lambda_t)$ .
- For example, generalizing the gradient case (where  $R(a) = ||a||_2^2$ ), we could consider

$$R_t(a) = a^{\top} M_t a,$$

- with  $M_t = (1 + \lambda_t) I$  (as before),
- with  $M_t$  a positive diagonal matrix, or
- with  $M_t \succ 0$  (an arbitrary positive definite matrix).
- We can view this as adapting the step-size in different directions.

Consider the following version of mirror descent (also called 'proximal gradient': stay close to  $a_t$  instead of  $\tilde{a}_t$ ):

```
a_{t+1} = \arg\min_{a \in \mathcal{A}} \left( \eta \nabla \ell_t(a_t) \cdot a + D_{R_t}(a, a_t) \right).
```

Consider the following version of mirror descent (also called 'proximal gradient': stay close to  $a_t$  instead of  $\tilde{a}_t$ ):

```
a_{t+1} = \arg\min_{a \in \mathcal{A}} \left( \eta \nabla \ell_t(a_t) \cdot a + D_{R_t}(a, a_t) \right).
```

Similar arguments give the following theorem.

#### Theorem

For  $R_t$  strongly-convex wrt some norm  $\|\cdot\|_{R_t}$ ,

$$\begin{aligned} R_n &\leq \frac{1}{\eta} D_{R_1}(a^*, a_1) + \frac{1}{\eta} \sum_{t=1}^{n-1} \left( D_{R_{t+1}}(a, a_{t+1}) - D_{R_t}(a, a_{t+1}) \right) \\ &+ \frac{\eta}{2} \sum_{t=1}^n \| \nabla \ell_t(a_t) \|_{R_{t,*}}^2 \,. \end{aligned}$$

For  $R_t(a) = a^{\top} M_t a$  with  $M_t$  a positive diagonal matrix, say,  $M_t = \text{diag}(s_t)$ , we have

$$D_{R_t}(a,b) = (a-b)^{\top} M_t(a-b) = \sum_i (a_i - b_i)^2 s_{t,i}.$$

For  $R_t(a) = a^{\top} M_t a$  with  $M_t$  a positive diagonal matrix, say,  $M_t = \text{diag}(s_t)$ , we have

$$D_{R_t}(a,b)=(a-b)^ op M_t(a-b)=\sum_i(a_i-b_i)^2s_{t,i}.$$

And  $D_{R_t}$  is strongly convex wrt the norm  $||a||_{R_t}^2 = 2a^\top M_t a$ .

For  $R_t(a) = a^{\top} M_t a$  with  $M_t$  a positive diagonal matrix, say,  $M_t = \text{diag}(s_t)$ , we have

$$D_{R_t}(a,b) = (a-b)^{ op} M_t(a-b) = \sum_i (a_i - b_i)^2 s_{t,i}.$$

And  $D_{R_t}$  is strongly convex wrt the norm  $||a||_{R_t}^2 = 2a^\top M_t a$ . Also

$$\|g\|_{R_{t},*}^{2} = \frac{1}{2}g^{\top}M_{t}^{-1}g. = \frac{1}{2}\sum_{i}\frac{g_{i}^{2}}{s_{t,i}}.$$

Applying the theorem, the regret satisfies

$$\begin{split} R_n &\leq \frac{1}{\eta} D_{R_1}(a^*, a_1) + \frac{1}{\eta} \sum_{t=1}^{n-1} \left( D_{R_{t+1}}(a, a_{t+1}) - D_{R_t}(a, a_{t+1}) \right) \\ &+ \frac{\eta}{2} \sum_{t=1}^n \| \nabla \ell_t(a_t) \|_{R_{t,*}}^2 \end{split}$$

Applying the theorem, the regret satisfies

$$\begin{split} R_n &\leq \frac{1}{\eta} D_{R_1}(a^*, a_1) + \frac{1}{\eta} \sum_{t=1}^{n-1} \left( D_{R_{t+1}}(a, a_{t+1}) - D_{R_t}(a, a_{t+1}) \right) \\ &+ \frac{\eta}{2} \sum_{t=1}^n \| \nabla \ell_t(a_t) \|_{R_{t,*}}^2 \\ &\leq \frac{1}{\eta} D_{R_1}(a^*, a_1) + \frac{1}{\eta} \sum_{t=1}^{n-1} \max_i (a_i^* - a_{t+1,i})^2 \| s_{t+1} - s_t \|_1 \\ &+ \frac{\eta}{4} \sum_{t=1}^n \nabla \ell_t(a_t)^\top \operatorname{diag}(s_t)^{-1} \nabla \ell_t(a_t). \end{split}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

#### Adagrad

[Duchi, Hazan, Singer, 2011]

If we insist that the regularization increases (that is, the components of  $s_t$  are monotonically non-decreasing with t), we can choose

$$s_{t,i} = \sqrt{\sum_{s=1}^t 
abla \ell_t(a_t)_i^2},$$

#### Adagrad

[Duchi, Hazan, Singer, 2011]

If we insist that the regularization increases (that is, the components of  $s_t$  are monotonically non-decreasing with t), we can choose

$$S_{t,i} = \sqrt{\sum_{s=1}^{t} \nabla \ell_t(a_t)_i^2},$$
  
$$\eta = D_{\infty} := \sup_{a^*, a_t} \|a^* - a_t\|_{\infty},$$

#### Adagrad

[Duchi, Hazan, Singer, 2011]

If we insist that the regularization increases (that is, the components of  $s_t$  are monotonically non-decreasing with t), we can choose

$$S_{t,i} = \sqrt{\sum_{s=1}^{t} \nabla \ell_t(a_t)_i^2},$$
  
$$\eta = D_{\infty} := \sup_{a^*, a_t} \|a^* - a_t\|_{\infty},$$

to give an adaptivity result (versus constant s):

$$R_n \leq c \min_{\eta,s} \left( \frac{D_{\infty}^2}{\eta} s^{\top} \mathbf{1} + \eta \sum_{t=1}^n \nabla \ell_t(a_t)^{\top} \operatorname{diag}(s)^{-1} \nabla \ell_t(a_t) \right)$$

#### Adagrad

[Duchi, Hazan, Singer, 2011]

If we insist that the regularization increases (that is, the components of  $s_t$  are monotonically non-decreasing with t), we can choose

$$s_{t,i} = \sqrt{\sum_{s=1}^{t} \nabla \ell_t(a_t)_i^2},$$
  
$$\eta = D_{\infty} := \sup_{a^*, a_t} \|a^* - a_t\|_{\infty},$$

to give an adaptivity result (versus *constant s*):

$$R_n \leq c \min_{\eta, s} \left( \frac{D_{\infty}^2}{\eta} s^{\top} \mathbf{1} + \eta \sum_{t=1}^n \nabla \ell_t(a_t)^{\top} \operatorname{diag}(s)^{-1} \nabla \ell_t(a_t) \right)$$
$$= O\left( D_{\infty} \sum_{i=1}^d \sqrt{\sum_{t=1}^n \nabla \ell_t(a_t)_i^2} \right).$$



	•	0	•	æ	×	•	2	×.	æ.	500
									1	17 / 132

• The gradient term might be much smaller than  $\sqrt{nd}$ .

- The gradient term might be much smaller than  $\sqrt{nd}$ .
- For instance, if the gradients are sparse and bounded (for instance, for logistic regression with sparse {0,1}-valued features), then we expect the gradient terms to be much smaller.
   For features that appear more frequently, the s<sub>t,i</sub> will be larger (learning rate slower in those directions).

- The gradient term might be much smaller than  $\sqrt{nd}$ .
- For instance, if the gradients are sparse and bounded (for instance, for logistic regression with sparse {0,1}-valued features), then we expect the gradient terms to be much smaller.
   For features that appear more frequently, the s<sub>t,i</sub> will be larger (learning rate slower in those directions).
- More generally, for coordinate directions with large gradients, we can make the corresponding component of *s* large (to keep things more stable in those directions), and for coordinate directions with small gradients, we can use less regularization.

A similar approach can be applied to matrices, with

$$M_t = \frac{\left(\sum_{s=1}^t \nabla \ell_t(a_t) \nabla \ell_t(a_t)^\top\right)^{1/2}}{\operatorname{tr}\left(\sum_{s=1}^t \nabla \ell_t(a_t) \nabla \ell_t(a_t)^\top\right)^{1/2}}$$

playing the role of  $s_t$ .

# Outline

## Binary prediction

### ② General online convex

- Empirical minimization fails
- Gradient algorithm
- A regularization viewpoint
- Bregman divergence
- Properties of regularization
- Linearization
- Mirror descent
- Regret bounds
- Strongly convex losses
- Adaptive regularization

### Minimax strategies

- Binary prediction
- ② General online convex
- Minimax strategies
  - Convex and strongly convex losses
  - The linear game

### The convex and linear games

For a convex set  $\mathcal{A} \subset \mathbb{R}^d$  and a sequence  $G_1, \ldots, G_n \ge 0$ , define  $\mathcal{G}_{conv}(\mathcal{A}, \{G_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

 $\ell_t \in \{\ell : \|\nabla \ell(a_t)\| \leq G_t, \ \ell \text{ convex}\}.$ 

### The convex and linear games

For a convex set  $\mathcal{A} \subset \mathbb{R}^d$  and a sequence  $G_1, \ldots, G_n \ge 0$ , define  $\mathcal{G}_{conv}(\mathcal{A}, \{G_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

 $\ell_t \in \{\ell : \|\nabla \ell(a_t)\| \leq G_t, \ \ell \text{ convex}\}.$ 

Define  $\mathcal{G}_{lin}(\mathcal{A}, \{G_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

$$\ell_t \in \left\{\ell: \ell(a) = v^ op(a-a_t) + c, \ v \in \mathbb{R}^d, \ c \in \mathbb{R}, \ \|v\| \leq G_t
ight\}.$$

イロト 不得下 イヨト イヨト 二日

### The convex and linear games

For a convex set  $\mathcal{A} \subset \mathbb{R}^d$  and a sequence  $G_1, \ldots, G_n \ge 0$ , define  $\mathcal{G}_{conv}(\mathcal{A}, \{G_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

 $\ell_t \in \{\ell : \|\nabla \ell(a_t)\| \leq G_t, \ \ell \text{ convex}\}.$ 

Define  $\mathcal{G}_{lin}(\mathcal{A}, \{G_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

$$\ell_t \in \left\{\ell: \ell(a) = v^ op(a-a_t) + c, \ v \in \mathbb{R}^d, \ c \in \mathbb{R}, \ \|v\| \leq G_t
ight\}.$$

• The adversary's constraints depend on the player's choices.

#### The strongly convex and quadratic games

For a convex set  $\mathcal{A} \subset \mathbb{R}^d$  and sequences  $G_1, \ldots, G_n \geq 0$  and  $\sigma_1, \ldots, \sigma_n \geq 0$ , define  $\mathcal{G}_{st-conv}(\mathcal{A}, \{G_t\}, \{\sigma_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

 $\ell_t \in \left\{\ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I\right\}.$ 

### The strongly convex and quadratic games

For a convex set  $\mathcal{A} \subset \mathbb{R}^d$  and sequences  $G_1, \ldots, G_n \geq 0$  and  $\sigma_1, \ldots, \sigma_n \geq 0$ , define  $\mathcal{G}_{st-conv}(\mathcal{A}, \{G_t\}, \{\sigma_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

 $\ell_t \in \left\{\ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I\right\}.$ 

Define  $\mathcal{G}_{quad}(\mathcal{A}, \{G_t\}, \{\sigma_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

$$\ell_t \in \left\{\ell : \ell(a) = v^\top (a - a_t) + \frac{\sigma_t}{2} \|a - a_t\|^2 + c, v \in \mathbb{R}^d, c \in \mathbb{R}, \|v\| \leq G_t\right\}$$

### The strongly convex and quadratic games

For a convex set  $\mathcal{A} \subset \mathbb{R}^d$  and sequences  $G_1, \ldots, G_n \geq 0$  and  $\sigma_1, \ldots, \sigma_n \geq 0$ , define  $\mathcal{G}_{st-conv}(\mathcal{A}, \{G_t\}, \{\sigma_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

 $\ell_t \in \left\{\ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I\right\}.$ 

Define  $\mathcal{G}_{quad}(\mathcal{A}, \{G_t\}, \{\sigma_t\})$  as the online convex optimization game with constraints  $a_t \in \mathcal{A}$  and

$$\ell_t \in \left\{\ell : \ell(a) = v^\top(a - a_t) + \frac{\sigma_t}{2} \|a - a_t\|^2 + c, v \in \mathbb{R}^d, c \in \mathbb{R}, \|v\| \leq G_t\right\}$$

• Again, the adversary's constraints depend on the player's choices.

#### Theorem

For fixed A,  $\{G_t\}$  and  $\{\sigma_t\}$ , we have

 $V_n(\mathcal{G}_{st-conv}(\mathcal{A}, \{G_t\}, \{\sigma_t\})) = V_n(\mathcal{G}_{quad}(\mathcal{A}, \{G_t\}, \{\sigma_t\})),$ 

[Abernethy, B., Rakhlin, Tewari, 2008]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

### Theorem

For fixed A,  $\{G_t\}$  and  $\{\sigma_t\}$ , we have

 $V_n(\mathcal{G}_{st-conv}(\mathcal{A}, \{G_t\}, \{\sigma_t\})) = V_n(\mathcal{G}_{quad}(\mathcal{A}, \{G_t\}, \{\sigma_t\})),$ 

$$V_n\left(\mathcal{G}_{conv}\left(\mathcal{A}, \{G_t\}\right)\right) = V_n\left(\mathcal{G}_{lin}\left(\mathcal{A}, \{G_t\}\right)\right).$$

[Abernethy, B., Rakhlin, Tewari, 2008]

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ ののの

123 / 132

Fix sets  $N_1, \ldots, N_n$  and  $M \subseteq N_t$ .

Suppose that for all  $\ell_t \in N_t$  and  $a_t \in A$  there is an  $\ell_t^* \in M$  such that for all  $a_1, \ell_1, \ldots, a_{t-1}, \ell_{t-1}$ , and  $a_{t+1}, \ell_{t+1}, \ldots, a_n, \ell_n$ ,

 $R_n(a_1,\ell_1,\ldots,a_t,\ell_t,\ldots,a_n,\ell_n) \leq R_n(a_1,\ell_1,\ldots,a_t,\ell_t^*,\ldots,a_n,\ell_n).$ 

Fix sets  $N_1, \ldots, N_n$  and  $M \subseteq N_t$ . Suppose that for all  $\ell_t \in N_t$  and  $a_t \in A$  there is an  $\ell_t^* \in M$  such that for all  $a_1, \ell_1, \ldots, a_{t-1}, \ell_{t-1}$ , and  $a_{t+1}, \ell_{t+1}, \ldots, a_n, \ell_n$ ,

$$R_n(a_1,\ell_1,\ldots,a_t,\ell_t,\ldots,a_n,\ell_n) \leq R_n(a_1,\ell_1,\ldots,a_t,\ell_t^*,\ldots,a_n,\ell_n).$$

Then

$$\inf_{a_1 \in \mathcal{A}} \sup_{\ell_1 \in N_1} \cdots \inf_{a_t \in \mathcal{A}} \sup_{\ell_t \in N_t} \cdots \inf_{a_n \in \mathcal{A}} \sup_{\ell_n \in N_n} R_n(a_1, \ell_1, \dots, a_n, \ell_n)$$

Fix sets  $N_1, \ldots, N_n$  and  $M \subseteq N_t$ . Suppose that for all  $\ell_t \in N_t$  and  $a_t \in A$  there is an  $\ell_t^* \in M$  such that for all  $a_1, \ell_1, \ldots, a_{t-1}, \ell_{t-1}$ , and  $a_{t+1}, \ell_{t+1}, \ldots, a_n, \ell_n$ ,

$$R_n(a_1,\ell_1,\ldots,a_t,\ell_t,\ldots,a_n,\ell_n) \leq R_n(a_1,\ell_1,\ldots,a_t,\ell_t^*,\ldots,a_n,\ell_n).$$

Then

$$\inf_{a_1 \in \mathcal{A}} \sup_{\ell_1 \in N_1} \cdots \inf_{a_t \in \mathcal{A}} \sup_{\ell_t \in N_t} \cdots \inf_{a_n \in \mathcal{A}} \sup_{\ell_n \in N_n} R_n(a_1, \ell_1, \dots, a_n, \ell_n)$$
  
= 
$$\inf_{a_1 \in \mathcal{A}} \sup_{\ell_1 \in N_1} \cdots \inf_{a_t \in \mathcal{A}} \sup_{\ell_t \in M} \cdots \inf_{a_n \in \mathcal{A}} \sup_{\ell_n \in N_n} R_n(a_1, \ell_1, \dots, a_n, \ell_n).$$

Fix sets  $N_1, \ldots, N_n$  and  $M \subseteq N_t$ . Suppose that for all  $\ell_t \in N_t$  and  $a_t \in A$  there is an  $\ell_t^* \in M$  such that for all  $a_1, \ell_1, \ldots, a_{t-1}, \ell_{t-1}$ , and  $a_{t+1}, \ell_{t+1}, \ldots, a_n, \ell_n$ ,

$$R_n(a_1,\ell_1,\ldots,a_t,\ell_t,\ldots,a_n,\ell_n) \leq R_n(a_1,\ell_1,\ldots,a_t,\ell_t^*,\ldots,a_n,\ell_n).$$

Then

$$\inf_{a_1 \in \mathcal{A}} \sup_{\ell_1 \in N_1} \cdots \inf_{a_t \in \mathcal{A}} \sup_{\ell_t \in N_t} \cdots \inf_{a_n \in \mathcal{A}} \sup_{\ell_n \in N_n} R_n(a_1, \ell_1, \dots, a_n, \ell_n)$$
  
= 
$$\inf_{a_1 \in \mathcal{A}} \sup_{\ell_1 \in N_1} \cdots \inf_{a_t \in \mathcal{A}} \sup_{\ell_t \in M} \cdots \inf_{a_n \in \mathcal{A}} \sup_{\ell_n \in N_n} R_n(a_1, \ell_1, \dots, a_n, \ell_n).$$

(Because  $M \subset N_t$ , and it contains  $\ell_t^*$  that's always at least as good as  $\ell_t$ .)

## Proof idea

For the strongly convex case, define

$$M := \left\{ \ell : \ell(\mathbf{a}) = \mathbf{v}^\top (\mathbf{a} - \mathbf{a}_t) + \frac{\sigma_t}{2} \|\mathbf{a} - \mathbf{a}_t\|^2 + \mathbf{c}, \|\mathbf{v}\| \leq G_t \right\},$$

### Proof idea

For the strongly convex case, define

$$\mathsf{M} := \left\{ \ell : \ell(\mathsf{a}) = \mathsf{v}^\top(\mathsf{a} - \mathsf{a}_t) + \frac{\sigma_t}{2} \|\mathsf{a} - \mathsf{a}_t\|^2 + \mathsf{c}, \, \|\mathsf{v}\| \leq \mathsf{G}_t \right\},$$

and notice that

 $M \subseteq N_t := \left\{ \ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I \right\}.$ 

### Proof idea

For the strongly convex case, define

$$\mathsf{M} := \left\{ \ell : \ell(\mathsf{a}) = \mathsf{v}^\top(\mathsf{a} - \mathsf{a}_t) + \frac{\sigma_t}{2} \|\mathsf{a} - \mathsf{a}_t\|^2 + \mathsf{c}, \, \|\mathsf{v}\| \leq \mathsf{G}_t \right\},$$

and notice that

$$M \subseteq N_t := \left\{ \ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I \right\}.$$

For  $\ell_t \in N_t$ , define  $\ell_t^*$  as

$$\ell_t^*(\mathsf{a}) = \ell_t(\mathsf{a}_t) + 
abla \ell_t(\mathsf{a}_t)^ op (\mathsf{a} - \mathsf{a}_t) + rac{\sigma_t}{2} \|\mathsf{a} - \mathsf{a}_t\|^2.$$

### Proof idea

For the strongly convex case, define

$$M := \left\{ \ell : \ell(\mathbf{a}) = \mathbf{v}^\top (\mathbf{a} - \mathbf{a}_t) + \frac{\sigma_t}{2} \|\mathbf{a} - \mathbf{a}_t\|^2 + \mathbf{c}, \|\mathbf{v}\| \leq G_t \right\},$$

and notice that

$$M \subseteq N_t := \left\{ \ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I \right\}.$$

For  $\ell_t \in N_t$ , define  $\ell_t^*$  as

$$\ell_t^*(\mathbf{a}) = \ell_t(\mathbf{a}_t) + \nabla \ell_t(\mathbf{a}_t)^\top (\mathbf{a} - \mathbf{a}_t) + \frac{\sigma_t}{2} \|\mathbf{a} - \mathbf{a}_t\|^2.$$

Notice that  $\ell_t^* \in M$ , since  $\ell_t^*(a_t) = \ell_t(a_t)$  and  $\nabla \ell_t(a_t) = \nabla \ell_t^*(a_t)$ .

### Proof idea

For the strongly convex case, define

$$\mathsf{M} := \left\{ \ell : \ell(\mathsf{a}) = \mathsf{v}^\top(\mathsf{a} - \mathsf{a}_t) + \frac{\sigma_t}{2} \|\mathsf{a} - \mathsf{a}_t\|^2 + \mathsf{c}, \, \|\mathsf{v}\| \leq \mathsf{G}_t \right\},$$

and notice that

$$M \subseteq N_t := \left\{ \ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I \right\}.$$

For  $\ell_t \in N_t$ , define  $\ell_t^*$  as

$$\ell_t^*(\mathbf{a}) = \ell_t(\mathbf{a}_t) + \nabla \ell_t(\mathbf{a}_t)^\top (\mathbf{a} - \mathbf{a}_t) + \frac{\sigma_t}{2} \|\mathbf{a} - \mathbf{a}_t\|^2.$$

Notice that  $\ell_t^* \in M$ , since  $\ell_t^*(a_t) = \ell_t(a_t)$  and  $\nabla \ell_t(a_t) = \nabla \ell_t^*(a_t)$ . Also,  $\ell_t(a) \ge \ell_t^*(a)$  for all a, so M and  $N_t$  satisfy the conditions of the lemma.

### Proof idea

For the strongly convex case, define

$$\mathsf{M} := \left\{ \ell : \ell(\mathsf{a}) = \mathsf{v}^\top(\mathsf{a} - \mathsf{a}_t) + \frac{\sigma_t}{2} \|\mathsf{a} - \mathsf{a}_t\|^2 + \mathsf{c}, \, \|\mathsf{v}\| \leq \mathsf{G}_t \right\},$$

and notice that

$$M \subseteq N_t := \left\{ \ell : \|\nabla \ell(a_t)\| \leq G_t, \, \nabla^2 \ell \succeq \sigma_t I \right\}.$$

For  $\ell_t \in N_t$ , define  $\ell_t^*$  as

$$\ell_t^*(\mathbf{a}) = \ell_t(\mathbf{a}_t) + \nabla \ell_t(\mathbf{a}_t)^\top (\mathbf{a} - \mathbf{a}_t) + \frac{\sigma_t}{2} \|\mathbf{a} - \mathbf{a}_t\|^2.$$

Notice that  $\ell_t^* \in M$ , since  $\ell_t^*(a_t) = \ell_t(a_t)$  and  $\nabla \ell_t(a_t) = \nabla \ell_t^*(a_t)$ . Also,  $\ell_t(a) \ge \ell_t^*(a)$  for all a, so M and  $N_t$  satisfy the conditions of the lemma. The convex/linear case uses a similar argument.

- Binary prediction
- ② General online convex
- Minimax strategies
  - Convex and strongly convex losses
  - The linear game

#### Theorem

For  $\mathcal{A} = \{a \in \mathbb{R}^d : ||a|| \le r\}$  with  $d \ge 3$ , and a fixed sequence  $\{G_t\}$ ,

$$V_n\left(\mathcal{G}_{conv}\left(\mathcal{A}, \{G_t\}
ight)
ight) = V_n\left(\mathcal{G}_{lin}\left(\mathcal{A}, \{G_t\}
ight)
ight)$$

$$= r \sqrt{\sum_{t=1}^{n} G_t^2}.$$

[Abernethy, B., Rakhlin, Tewari, 2008]

**9** Wlog, we can assume r = 1 and  $\ell_t(a) = w^\top a$  with  $||w|| \leq G_t$ .

• Wlog, we can assume r = 1 and  $\ell_t(a) = w^\top a$  with  $||w|| \leq G_t$ .

$$e Writing W_t := \sum_{s=1}^t w_s,$$

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

# The linear game

## Proof

The adversary can ensure

$$R_n \geq \sqrt{\sum_{t=1}^n G_t^2},$$

by playing  $w_t$  satisfying

$$w_t^{\top} a_t = 0, \quad w_t^{\top} W_{t-1} = 0, \quad ||w_t|| = G_t.$$

# The linear game

### Proof

The adversary can ensure

$$R_n \geq \sqrt{\sum_{t=1}^n G_t^2},$$

by playing  $w_t$  satisfying

$$w_t^{\top} a_t = 0, \quad w_t^{\top} W_{t-1} = 0, \quad ||w_t|| = G_t.$$

To see this, notice that this choice ensures  $\sum_{t=1}^{n} \ell_t(a_t) = 0$  and so  $R_n = ||W_n||$ .

# The linear game

### Proof

The adversary can ensure

$$R_n \geq \sqrt{\sum_{t=1}^n G_t^2},$$

by playing  $w_t$  satisfying

$$w_t^{\top} a_t = 0, \quad w_t^{\top} W_{t-1} = 0, \quad ||w_t|| = G_t.$$

To see this, notice that this choice ensures  $\sum_{t=1}^{n} \ell_t(a_t) = 0$  and so  $R_n = ||W_n||$ . But

$$\|W_t\| = \|W_{t-1} + w_t\| = \sqrt{\|W_{t-1}\|^2 + \|w_t\|^2} = \sqrt{\sum_{s=1}^t G_s^2}.$$

**(**) If the player defines  $W_0 = 0$  and chooses

$$a_t = rac{-W_{t-1}}{\sqrt{\|W_{t-1}\|^2 + \sum_{s=t}^n G_s^2}},$$

then

$$R_n \leq \sqrt{\sum_{t=1}^n G_t^2}.$$

・ロ ・ ・ 日 ・ ・ 目 ・ 日 ・ 日 ・ 130/132

This is equivalent to showing that, for this  $a_t$ , no matter what choices of  $w_t$  the adversary makes,

$$\sum_{t=1}^n w_t^\top a_t + \|W_n\| \leq \sqrt{\sum_{t=1}^n G_t^2}.$$

This is equivalent to showing that, for this  $a_t$ , no matter what choices of  $w_t$  the adversary makes,

$$\sum_{t=1}^n w_t^\top a_t + \|W_n\| \leq \sqrt{\sum_{t=1}^n G_t^2}.$$

The proof is by a backward induction, and involves a 2-dimensional geometric argument (since  $a_t$  is aligned with  $W_{t-1}$ , we need only consider the role of  $w_t$ ).

## Binary prediction

- ② General online convex
- Minimax strategies
  - Convex and strongly convex losses
  - The linear game