

# Online Learning and Online Convex Optimization

**Chapter 2 in Shai Shalev Shwartz / Online Learning and Online convex Optimization**

# Outline

## Online Convex Optimization

### Follow The Leader

- Quadratic Optimization

- Failure of Follow the Leader

### Follow The Regularized Leader

- FTRL for linear functions

- Online Gradient Descent

- Doubling Trick

- Strong Convexity

- General Theorem regarding FTRL with Strong Convexity

- Applications to expert advice

### Mirror Descent

- Potential Based Gradient descent

- OMD for linear cost functions

- sub-gradients

### Duality

### Algorithms for specific potentials

# Online Convex Optimization (OCO)

## Algorithm

**Input:** A convex set  $S$

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

- ▶ Predict a vector  $w_t \in S$
- ▶ Receive a convex loss function  $f_t : S \rightarrow \mathbb{R}$
- ▶ Suffer loss  $f_t(w_t)$

## Regret Definition

**Regret of the Algorithm:**

$$\text{Regret}_T(u) = \sum_{t=1}^T f_t(w_t) - \sum_{t=1}^T f_t(u). \quad (1)$$

**Regret relative to a set of vectors  $U$ :**

$$\text{Regret}_T(U) = \max_{u \in U} \text{Regret}_T(u). \quad (2)$$

# Follow-the-Leader Algorithm

## FTL Strategy

At round  $t$ , select:

$$w_t = \operatorname{argmin}_{w \in S} \sum_{i=1}^{t-1} f_i(w)$$

- ▶ Natural approach: Choose best performer on past data
- ▶ Simple but can be unstable
- ▶ Requires solving optimization problem each round

# FTL Regret Analysis

## Theorem (Lemma 2.1)

For any  $u \in S$ :

$$\text{Regret}_T(u) = \sum_{t=1}^T (f_t(w_t) - f_t(u)) \leq \sum_{t=1}^T (f_t(w_t) - f_t(w_{t+1})).$$

**proof**

**Step 1:** Equivalent to

$$\sum_{t=1}^T f_t(w_{t+1}) \leq \sum_{t=1}^T f_t(u)$$

**Step 2:** By induction on  $T$ :

- ▶ **Base case:**  $T = 1$  trivial as  $f_1(w_1) - f_1(u) \leq 0$
- ▶ **Inductive step:** Assume holds for  $T - 1$ , then

$$\begin{aligned} & \sum_{t=1}^T [f_t(w_t) - f_t(u)] \\ &= \underbrace{\sum_{t=1}^{T-1} [f_t(w_t) - f_t(u)]}_{\leq \sum_{t=1}^{T-1} [f_t(w_t) - f_t(w_{t+1})]} + [f_T(w_T) - f_T(u)] \\ &\leq \sum_{t=1}^T [f_t(w_t) - f_t(w_{t+1})] \end{aligned}$$

using  $w_{T+1} = \operatorname{argmin}_w \sum_{t=1}^T f_t(w)$

## FTL for Quadratic Optimization

For  $f_t(w) = \frac{1}{2} \|w - z_t\|_2^2$ :

- ▶ FTL update:  $w_t = \frac{1}{t-1} \sum_{i=1}^{t-1} z_i$
- ▶ Regret bound:  $O(\log T)$

Regret Calculation for quadratic optimization.

$$\begin{aligned} \text{Regret}_T(u) &\leq \sum_{t=1}^T \frac{1}{t} \|w_t - z_t\|^2 \\ &\leq \sum_{t=1}^T \frac{(2L)^2}{t} = 4L^2(\log T + 1) \end{aligned}$$

where  $L = \max_t \|z_t\|$





## Failure of follow the leader

$$f_t(w) = w \cdot z:$$



$$z_t = \begin{cases} -0.5 & \text{if } t = 1 \\ 1 & \text{if } t \text{ is even} \\ -1 & \text{if } t > 1 \text{ and } t \text{ is odd} \end{cases}$$

- ▶  $w_t = -1, 1, -1, 1, \dots$
- ▶ Cumulative loss is  $T$ .
- ▶ Cumulative loss of 0 is 0
- ▶ Regret is  $T$ .
- ▶ **Reason:** prediction is unstable
- ▶ We need to regularize.
- ▶  $R(W)$  penalizes vectors which are large.

## Follow-the-Regularized-Leader (FTRL)

$$\forall t, \quad w_t = \arg \min_{w \in S} \sum_{i=1}^{t-1} f_i(w) + R(w)$$

- ▶ For bad case above:  $w_t = 0, 0, 0, 0, \dots$
- ▶ Each step requires solving a minimization problem.

## Lemma 2.3: Follow-the-Regularized-Leader

**Lemma 2.3.** Let  $w_1, w_2, \dots$  be the sequence of vectors produced by FoReL. Then, for all  $u \in S$  we have:

$$\sum_{t=1}^T (f_t(w_t) - f_t(u)) \leq R(u) - R(w_1) + \sum_{t=1}^T (f_t(w_t) - f_t(w_{t+1})).$$

## Proof of Lemma 2.3

*Proof.* Observe that running FoReL on  $f_1, \dots, f_T$  is equivalent to running FTL on  $f_0, f_1, \dots, f_T$  where  $f_0 = R$ . Using Lemma 2.1, we obtain:

$$\sum_{t=0}^T (f_t(w_t) - f_t(u)) \leq \sum_{t=0}^T (f_t(w_t) - f_t(w_{t+1})).$$

Rearranging the above and using  $f_0 = R$ , we conclude our proof.  $\square$

## FTRL Regret Bound for linear functions

**For linear**  $f_t(w) = \langle w, z_t \rangle$  and  $R(w) = \frac{1}{2\eta} \|w\|_2^2$

**Update rule**  $w_{t+1} = w_t - \eta z_t$  Then, for all  $u$  we have

$$\text{Regret}_T(u) \leq \frac{1}{2\eta} \|u\|_2^2 + \eta \sum_{t=1}^T \|z_t\|_2^2.$$

## Choice of $\eta$ and Final Bound for linear functions

### Tunings:

► Define the set  $U = \{u : \|u\| \leq B\}$ .

► Assume that

$$\frac{1}{T} \sum_{t=1}^T \|z_t\|_2^2 \leq L^2.$$

► Set  $\eta = \frac{B}{L\sqrt{2T}}$ .

### Conclusion:

$$\text{Regret}_T(U) \leq BL\sqrt{2T}.$$

## From linear functions to Online Gradient Descent

### Example (OGD from FTRL)

Consider the OCO setup where the functions  $f_1, f_2, \dots$  are differentiable. Let  $\eta$  be the learning rate.

$$w_{t+1} = w_t - \eta z_t, \quad z_t = \nabla f_t(w_t)$$

Identical to FTRL with regularization:  $R(w) = \frac{1}{2\eta} \|w\|_2^2$

**Regret bound on OGD:** From FTRL theorem:

$$\begin{aligned} \text{Regret} &\leq \frac{\|u\|^2}{2\eta} + \eta \sum_{t=1}^T \|z_t\|^2 \\ &\leq \frac{B^2}{2\eta} + \eta T L^2 \end{aligned}$$

## Regret Bound for OGD

If we further assume that each  $f_t$  is  $L_t$ -Lipschitz with respect to  $\|\cdot\|_2$ , and let  $L$  be such that

$$\frac{1}{T} \sum_{t=1}^T L_t^2 \leq L^2.$$

Then, for all  $\mathbf{u}$ , the regret of OGD satisfies

$$\text{Regret}_T(\mathbf{u}) \leq \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta T L^2.$$



## Bounding the norm of $\mathbf{u}$

In particular, if  $U = \{\mathbf{u} : \|\mathbf{u}\|_2 \leq B\}$  and  $\eta = \frac{B}{L\sqrt{2T}}$  then

$$\text{Regret}_T(U) \leq BL\sqrt{2T}.$$

## Practical Considerations

### Doubling Trick

- ▶ Removes need to know time horizon  $T$
- ▶ Divide time into epochs  $2^m, 2^{m+1} - 1$
- ▶ Regret increases by constant factor:

$$\sum_{m=0}^{\log T} \sqrt{2^m} = O(\sqrt{T})$$

### Example (Optimal $\eta$ )

Setting  $\eta = \frac{B}{L} \sqrt{\frac{2}{T}}$  gives:

$$BL\sqrt{2T}$$

## Definition 2.4: Strong Convexity

### Strong Convexity

A function  $f : S \rightarrow \mathbb{R}$  is  $\sigma$ -strongly convex over  $S$  with respect to a norm  $\|\cdot\|$  if for any  $w \in S$  we have:

$$\forall z \in \partial f(w), \quad \forall u \in S, \quad f(u) \geq f(w) + \langle z, u - w \rangle + \frac{\sigma}{2} \|u - w\|^2.$$

## Lemma 2.8: Strong Convexity implication

### Lemma 2.8

Let  $S$  be a nonempty convex set. Let  $f : S \rightarrow \mathbb{R}$  be a  $\sigma$ -strongly convex function over  $S$  with respect to a norm  $\|\cdot\|$ . Let:

$$w = \arg \min_{v \in S} f(v).$$

Then, for all  $u \in S$ , we have:

$$f(u) - f(w) \geq \frac{\sigma}{2} \|u - w\|^2.$$

## Strong Convexity Condition

If  $R$  is twice differentiable, then it is easy to verify that a sufficient condition for strong convexity of  $R$  is that for all  $\mathbf{w}, \mathbf{x}$ ,

$$\langle \nabla^2 R(\mathbf{w}) \mathbf{x}, \mathbf{x} \rangle \geq \sigma \|\mathbf{x}\|^2,$$

where  $\nabla^2 R(\mathbf{w})$  is the Hessian matrix of  $R$  at  $\mathbf{w}$ , namely, the matrix of second-order partial derivatives of  $R$  at  $\mathbf{w}$  [39, Lemma 14].

## Example 2.4: Euclidean Regularization

The function

$$R(w) = \frac{1}{2} \|w\|_2^2$$

is 1-strongly-convex with respect to the  $\ell_2$  norm over  $\mathbb{R}^d$ . To see this, simply note that the Hessian of  $R$  at any  $w$  is the identity matrix.

## Example 2.5: Entropic Regularization

The function

$$R(w) = \sum_{i=1}^d w[i] \log(w[i])$$

is  $\frac{1}{B}$ -strongly-convex with respect to the  $\ell_1$  norm over the set

$$S = \{w \in \mathbb{R}^d : w > 0 \wedge \|w\|_1 \leq B\}.$$

In particular,  $R$  is 1-strongly-convex over the probability simplex, which is the set of positive vectors whose elements sum to 1.

## Proof of strong convexity for Entropic Regularization

$$\frac{\partial^2}{\partial w[i]^2} w[i] \log w[i] = \frac{1}{w[i]}$$

$$\begin{aligned}\langle \nabla^2 R(w) x, x \rangle &= \sum_i \frac{x[i]^2}{w[i]} \\&= \frac{1}{\|w\|_1} \left( \sum_i w[i] \right) \left( \sum_i \frac{x[i]^2}{w[i]} \right) \\&\geq \frac{1}{\|w\|_1} \left( \sum_i \sqrt{w[i]} \frac{x[i]}{\sqrt{w[i]}} \right)^2 = \frac{\|x\|_1^2}{\|w\|_1},\end{aligned}$$

where the inequality follows from Cauchy–Schwarz inequality.



## Single Step of FTRL with Strong Convexity

Let

$$R : S \rightarrow \mathbb{R}$$

be a  $\sigma$ -strongly-convex function over  $S$  with respect to a norm  $\|\cdot\|$ . Let  $w_1, w_2, \dots$  be the predictions of the FoReL algorithm. Then, for all  $t$ , if  $f_t$  is  $L_t$ -Lipschitz with respect to  $\|\cdot\|$ , then:

$$f_t(w_t) - f_t(w_{t+1}) \leq L_t \|w_t - w_{t+1}\| \leq \frac{L_t^2}{\sigma}.$$

## Proof (Single Step of FTRL with Strong Convexity)

For all  $t$  let

$$F_t(w) = \sum_{i=1}^{t-1} f_i(w) + R(w)$$

and note that the FoReL rule is

$$w_t = \arg \min_{w \in S} F_t(w).$$

Note also that  $F_t$  is  $\sigma$ -strongly-convex since the addition of a convex function to a strongly convex function keeps the strong convexity property. Therefore, Lemma 2.8 implies that:

$$F_t(w_{t+1}) \geq F_t(w_t) + \frac{\sigma}{2} \|w_t - w_{t+1}\|^2.$$

## Continuing the Proof (Single Step of FTRL with Strong Convexity)

Repeating the same argument for  $F_{t+1}$  and its minimizer  $w_{t+1}$ , we get:

$$F_{t+1}(w_t) \geq F_{t+1}(w_{t+1}) + \frac{\sigma}{2} \|w_t - w_{t+1}\|^2.$$

Taking the difference between the last two inequalities and rearranging, we obtain:

$$\sigma \|w_t - w_{t+1}\|^2 \leq f_t(w_t) - f_t(w_{t+1}). \quad (2.7)$$

## Final Steps (Single Step of FTRL with Strong Convexity)

Next, using the Lipschitzness of  $f_t$ , we get that:

$$f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1}) \leq L_t \|\mathbf{w}_t - \mathbf{w}_{t+1}\|.$$

Combining with Equation (2.7) and rearranging, we get:

$$\|\mathbf{w}_t - \mathbf{w}_{t+1}\| \leq L/\sigma.$$

Together with the above, we conclude our proof. ■

## Main theorem regarding $\sigma$ -strongly convex regularization functions

Let  $f_1, \dots, f_T$  be a sequence of convex functions such that  $f_t$  is  $L_t$ -Lipschitz with respect to some norm  $\|\cdot\|$ . Let  $L$  be such that

$$\frac{1}{T} \sum_{t=1}^T L_t^2 \leq L^2.$$

Assume that FoReL is run on the sequence with a regularization function which is  $\sigma$ -strongly-convex with respect to the same norm. Then, for all  $u \in S$ ,

$$\text{Regret}_T(u) \leq R(u) - \min_{v \in S} R(v) + \frac{TL^2}{\sigma}.$$

## Corollary for $l_2$ regularization

Let  $f_1, \dots, f_T$  be a sequence of convex functions such that  $f_t$  is  $L_t$ -Lipschitz with respect to  $\|\cdot\|_2$ . Let  $L$  be such that

$$\frac{1}{T} \sum_{t=1}^T L_t^2 \leq L^2.$$

Assume that FoReL is run on the sequence with the regularization function

$$R(w) = \frac{1}{2\eta} \|w\|_2^2.$$

Then, for all  $u$ ,

$$\text{Regret}_T(u) \leq \frac{1}{2\eta} \|u\|_2^2 + \eta T L^2.$$

## Applications to expert advice

- ▶ Distribution  $w_t$
- ▶ Action Losses:  $x_t \in [0, 1]^d$
- ▶ Algorithm Loss:  $\langle x_t, w_t \rangle$
- ▶ We want to bound regret.
- ▶ we will compare  $l_2$  regularization with Entropic Regularization.

## Experts using $l_2$ regularization (1)

$S$  be a convex set and consider running FoReL with the regularization function:

$$R(w) = \begin{cases} \frac{1}{2\eta} \|w\|_2^2 & \text{if } w \in S \\ \infty & \text{if } w \notin S \end{cases}$$

Where  $S$  is the  $d$  dimensional simplex.

Then, for all  $u \in S$ ,

$$\text{Regret}_T(u) \leq \frac{1}{2\eta} \|u\|_2^2 + \eta TL^2.$$



## Experts using $l_2$ regularization (2)

If

$$B \geq \max_{u \in S} \|u\|_2$$

Setting

$$B = 1; L = \sqrt{d}; \eta = \frac{B}{L\sqrt{2T}} = \frac{1}{\sqrt{2dT}}$$

then,

$$\text{Regret}_T(S) \leq \sqrt{2dT}.$$

## Entropic Regularization

Let  $f_1, \dots, f_T$  be a sequence of convex functions such that  $f_t$  is  $L_t$ -Lipschitz with respect to  $\|\cdot\|_1$ . Let  $L$  be such that  $\frac{1}{T} \sum_{t=1}^T L_t^2 \leq L^2$ . Assume that FoReL is run on the sequence with the regularization function

$$R(w) = \frac{1}{\eta} \sum_i w[i] \log(w[i])$$

and with the set

$$S = \{w : \|w\|_1 = B \wedge w > 0\} \subset \mathbb{R}^d.$$

Then,

$$\text{Regret}_T(S) \leq \frac{B \log(d)}{\eta} + \eta B T L^2.$$

In particular, setting  $\eta = \frac{\sqrt{\log d}}{L\sqrt{2T}}$  yields

$$\text{Regret}_T(S) \leq B L \sqrt{2 \log(d) T}.$$

## Entropic regularization for Experts

The Entropic regularization is strongly convex with respect to the  $\ell_1$  norm, and therefore the Lipschitzness requirement of the loss functions is also with respect to the  $\ell_1$ -norm.

For linear functions,

$$f_t(w) = \langle w, x_t \rangle,$$

we have by Hölder's inequality that,

$$|f_t(w) - f_t(u)| = |\langle w - u, x_t \rangle| \leq \|w - u\|_1 \|x_t\|_\infty.$$

Therefore, the Lipschitz parameter grows with the  $\ell_\infty$  norm of  $x_t$  rather than the  $\ell_2$  norm of  $x_t$ .

expert advice:  $B = 1$  and  $L = 1$ ), we obtain the regret bound of

$$\sqrt{2 \log(d) T}.$$

## Comparison between regularizations

- ▶ entropic regularization vs.  $\ell_2$  regularization.
- ▶  $\log d$  vs  $\sqrt{d}$
- ▶  $L$ :  $\|x_t\|_\infty \geq \|x_t\|_2$  Lipschitz condition carries heavier penalty with entropic regularization.
- ▶  $B$ :  $\|u\|_1 \leq \|u\|_2$  Comparator length carries heavier penalty with  $l_2$  norm.

## Potential based gradient Descent

- ▶  $\text{Regret}_t$  = Regret vector  $\text{Regret}_t(w) = L_{A,t} - L_t(w)$
- ▶  $\text{Regret}_t$  = State of prediction algorithm at time  $t$
- ▶ Potential/Regularizer:  $R(\text{Regret})$  Quantifies **badness** of the state.
- ▶ A state is bad if adversary can force high regret in the future.
- ▶ Choose prediction so that  $R(\text{Regret}_{t+1}) - R(\text{Regret}_t) + w_t \cdot \ell_t$  is small for all possible  $\ell_t$
- ▶  $w_t = \nabla R(\text{Regret}_t)$  is a good choice.
- ▶ For finite number of experts,  $\text{Regret}_t$  is finite dimensional and we can compute  $w_t$  explicitly.
- ▶ Here,  $\text{Regret} = \{R(w)\}_{w \in \mathbb{R}^d}$  is uncountably infinite.
- ▶ If Experts correspond to exponential distributions and loss is log loss- we can use conjugate priors. (recall: biased coins).
- ▶ We need a new trick to compute  $w_t = \nabla R(\text{Regret}_t)$  efficiently.

## FoReL Update Rule for linear cost function

Define  $\mathbf{z}_{1:t} = \sum_{i=1}^t \mathbf{z}_i$ , the FoReL update rule can be written as

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} R(\mathbf{w}) + \sum_{i=1}^t \langle \mathbf{w}, \mathbf{z}_i \rangle$$

$$= \arg \min_{\mathbf{w}} R(\mathbf{w}) + \langle \mathbf{w}, \mathbf{z}_{1:t} \rangle$$

$$= \arg \max_{\mathbf{w}} \langle \mathbf{w}, -\mathbf{z}_{1:t} \rangle - R(\mathbf{w}).$$

## Mirror Descent Update for linear functions

Update rule

$$\mathbf{w}_{t+1} = \arg \max_{\mathbf{w}} \langle \mathbf{w}, -\mathbf{z}_{1:t} \rangle - R(\mathbf{w}).$$

Link Function:

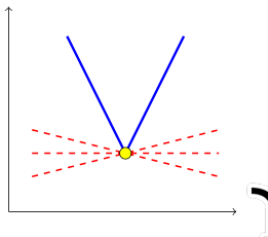
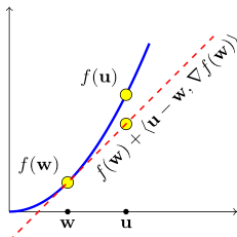
$$g(\theta) = \arg \max_{\mathbf{w}} \langle \mathbf{w}, \theta \rangle - R(\mathbf{w}),$$

Update rule can be re-written as

1.  $\theta_0 = 0$
2.  $\theta_{t+1} = \theta_t - \mathbf{z}_t$
3.  $\mathbf{w}_{t+1} = g(\theta_{t+1})$

## Sub-Gradients

- ▶ we can reduce general convex to linear using the gradient.
- ▶ What can we do if  $f(x)$  is convex but not differentiable at  $x$ ?
- ▶ Use the sub-gradients at  $x \doteq \partial f(x)$ : the set of linear functions such that  $l(x) = \langle w, x \rangle + o$  such that  $\forall y, l(y) \leq f(y)$  and  $l(x) = f(x)$
- ▶ if gradient  $\nabla f(x)$  exists, then  $\partial f(x) = \{\nabla f(x)\}$





## Example Generalized Online Gradient Descent

Consider the  $\ell_2$  setup where the functions  $f_1, f_2, \dots$  are convex (but not necessarily differentiable). Let  $\eta$  be the learning rate.

$$w_{t+1} = w_t - \eta z_t, \quad z_t \in \partial f_t(w_t)$$

Identical to FTRL with regularization:  $R(w) = \frac{1}{2\eta} \|w\|_2^2$

**Regret bound on OGD:** From FTRL theorem:

$$\begin{aligned} \text{Regret} &\leq \frac{\|u\|^2}{2\eta} + \eta \sum_{t=1}^T \|z_t\|^2 \\ &\leq \frac{B^2}{2\eta} + \eta T L^2 \end{aligned}$$

## Online Mirror Descent (OMD)

**parameter:** a link function  $g : \mathbb{R}^d \rightarrow S$

**initialize:**  $\theta_1 = 0$

**for**  $t = 1, 2, \dots$

▶ **predict**  $w_t = g(\theta_t)$

▶ **update**  $\theta_{t+1} = \theta_t - z_t$  where  $z_t \in \partial f_t(w_t)$

## Duality

- ▶ OMD can be analyzed using elementary tools.
- ▶ Using Duality Gives better intuition, more general analysis, tighter bounds.

## Dual Vector Spaces

- ▶  $V$  is a vector space, with a norm  $\|v\|$
- ▶  $U$  is the set of all linear mappings from  $V$  to  $\mathbb{R}$
- ▶ The norm of  $u \in U$  is defined as

$$\|u\|^* = \max_{v \in V} \frac{\|u(v)\|}{\|v\|}$$

- ▶  $U^*$  is equivalent to the set of all linear mappings from  $U$  to  $\mathbb{R}$ .
- ▶  $U$  and  $U^*$  are dual vector spaces, with dual norms.

## Dual Norms

- ▶ The space is always  $U, V = \mathbb{R}^n$
- ▶ The linear operation is the dot product  $\mathbf{u} \cdot \mathbf{v}$
- ▶  $L_2$  norm:  $\sqrt{\sum_{i=1}^n x_i^2}$
- ▶  $L_1$  norm:  $\sum_{i=1}^n |x_i|$
- ▶  $L_\infty$  norm:  $\max_i |x_i|$
- ▶  $L_p$  norm:  $(\sum_{i=1}^n x_i^p)^{\frac{1}{p}}$
- ▶  $L_p, L_q$  are dual norms if  $p, q \geq 1$ , and  $\frac{1}{p} + \frac{1}{q} = 1$
- ▶  $L_1, L_\infty$  are dual.
- ▶  $L_2$  is self-dual.

## Fenchel Duality

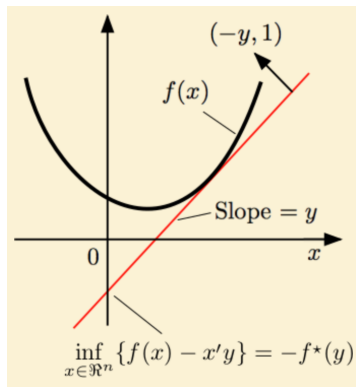
- ▶ Suppose  $F : A \rightarrow \mathbb{R}$  is a convex function over a convex set  $A \subseteq \mathbb{R}^n$ .
- ▶ The dual function to  $F$  is

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

- ▶ Fenchel duality Reduces to Legendre duality for differentiable functions

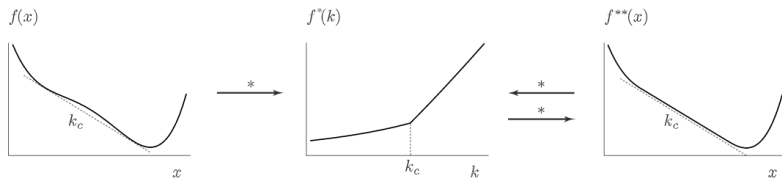
## Visualization of the Fenchel Dual

- ▶  $x, y \in \mathbb{R}$
- ▶  $f^*(y) = \sup_{x \in \mathbb{R}} (xy - f(x))$
- ▶  $-f^*(y) = \inf_{x \in \mathbb{R}} (f(x) - xy)$



# Dual of Dual

- ▶ The dual of any function is convex.
- ▶ if  $F$  is convex then  $F^{**} = F$

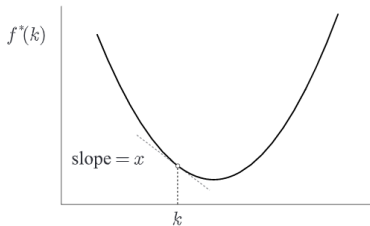
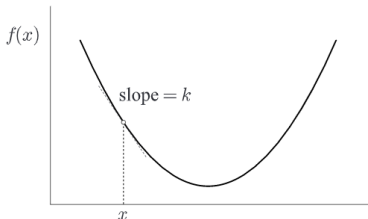




## Gradient Duality

- ▶ If the gradient of  $f$  at  $x$  is  $k$  then the gradient of  $f^*$  at  $k$  is  $x$
- ▶ In general:

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

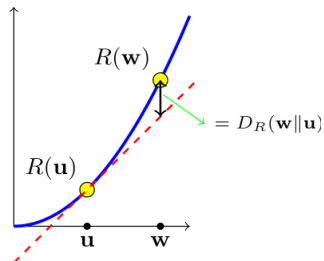


## Example: Exponential Potential

- ▶ Potential:  $F(u) = \sum_{i=1}^d e^{u_i}$
- ▶ Gradient:  $\nabla F(u)_i = e^{u_i}$  or  $\nabla F(u) = F(u)$ .
- ▶ Dual:  $F^*(v) = \sum_{i=1}^d v_i (\ln v_i - 1)$
- ▶ Gradient of dual:  $\nabla F^*(v)_i = \ln v_i$
- ▶ Note  $(\nabla F)^{-1} = \nabla F^*$

# Bregman Divergence

- ▶  $R(x)$  is convex and differentiable.
- ▶  $D_R(w||u) = R(w) - (R(u) + \langle \nabla R(u), (w - u) \rangle)$



## Fenchel and Bregman

- ▶  $F$ : strictly convex with continuous first derivative.
- ▶  $F^*$  is the Fenchel Dual of  $F$
- ▶  $D_F, D_{F^*}$  Bregman divergences wrt  $F, F^*$
- ▶  $u' = \nabla F(u)$  and  $v' = \nabla F(v)$
- ▶  $D_F(u, v) = D_{F^*}(u', v')$

## Mirror Descent

- ▶ Gradient descent in dual space  $\theta_t = \theta_{t-1} - \lambda \nabla \ell_t(\theta_{t-1})$
- ▶ Using duality can be rewritten as

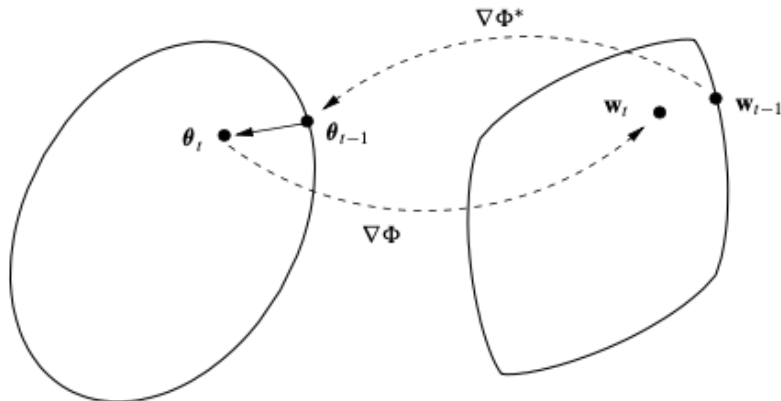
$$\nabla R^*(w_t) = \nabla R^*(w_{t-1}) - \lambda \nabla \ell_t(w_{t-1})$$

- ▶ As  $\nabla R$  is the inverse of  $\nabla R^*$  we get

$$w_t = \nabla R(\nabla R^*(w_{t-1}) - \lambda \nabla \ell_t(w_{t-1}))$$

## A picture of mirror descent

$$w_t = \nabla R(\nabla R^*(w_{t-1}) - \lambda \nabla \ell_t(w_{t-1}))$$



## Mirror Descent Regret Bound

**Goal:** Minimize regret in online convex optimization:

$$\text{Regret}(T) = \sum_{t=1}^T \langle w_t - w^*, z_t \rangle.$$

**Mirror Descent Update Rule:**

$$w_{t+1} = \arg \min_{w \in \mathcal{W}} [\eta \langle w, z_t \rangle + D_R(w \| w_t)].$$

where:

- ▶  $R(w)$  is a **strongly convex** potential function.
- ▶  $D_R(w \| u)$  is the **Bregman divergence**:

$$D_R(w \| u) = R(w) - R(u) - \langle \nabla R(u), w - u \rangle.$$

## Step 1: Understanding the First-Order Optimality Condition

### Why does the first-order optimality condition hold?

Since  $w_{t+1}$  minimizes the mirror descent objective:

$$w_{t+1} = \arg \min_w [\eta \langle w, z_t \rangle + D_R(w \| w_t)],$$

the function being minimized is **convex**. Hence, its **first-order optimality condition** states:

$$\langle w - w_{t+1}, \eta z_t + \nabla R(w_{t+1}) - \nabla R(w_t) \rangle \geq 0, \quad \forall w.$$

**Interpretation:** This inequality means that at  $w_{t+1}$ , the directional derivative of the objective function is **non-negative** for all feasible points, ensuring that  $w_{t+1}$  is a minimizer.



## Step 2: Applying the Optimality Condition to Regret

Setting  $w = w^*$  in the optimality condition:

$$\langle w^* - w_{t+1}, \eta z_t + \nabla R(w_{t+1}) - \nabla R(w_t) \rangle \geq 0.$$

Rearrange:

$$\langle w^* - w_{t+1}, \eta z_t \rangle \geq \langle w^* - w_{t+1}, \nabla R(w_{t+1}) - \nabla R(w_t) \rangle.$$

**Key Idea:** The \*\*gradient of  $R$  acts as a mirror map\*\*, ensuring that updates remain in the feasible region.

## Step 3: Using the Three-Point Bregman Identity

**Bregman divergence identity:**

$$D_R(w^* \| w_t) - D_R(w^* \| w_{t+1}) - D_R(w_{t+1} \| w_t) = \langle w^* - w_{t+1}, \nabla R(w_{t+1}) - \nabla R(w_t) \rangle.$$

Substituting into our previous inequality:

$$\langle w^* - w_{t+1}, \eta z_t \rangle \geq D_R(w^* \| w_t) - D_R(w^* \| w_{t+1}) - D_R(w_{t+1} \| w_t).$$

Rearrange:

$$\langle w_t - w^*, z_t \rangle \leq \frac{1}{\eta} (D_R(w^* \| w_t) - D_R(w^* \| w_{t+1})).$$

## Step 4: Summing Over All Iterations

Summing from  $t = 1$  to  $T$ :

$$\sum_{t=1}^T \langle w_t - w^*, z_t \rangle \leq \frac{1}{\eta} \sum_{t=1}^T (D_R(w^* \| w_t) - D_R(w^* \| w_{t+1})).$$

**Why does this work?**

- ▶ The sum on the right forms a **telescoping series**.
- ▶ All intermediate terms cancel, leaving:

$$\sum_{t=1}^T \langle w_t - w^*, z_t \rangle \leq \frac{1}{\eta} (D_R(w^* \| w_1) - D_R(w^* \| w_{T+1})).$$

**Since Bregman divergence is non-negative:**

$$\sum_{t=1}^T \langle w_t - w^*, z_t \rangle \leq \frac{1}{\eta} D_R(w^* \| w_1).$$

## Final Regret Bound

### Conclusion: Mirror Descent Regret Bound

$$\sum_{t=1}^T \langle w_t - w^*, z_t \rangle \leq \frac{1}{\eta} D_R(w^* \| w_1).$$

### Key Takeaways:

- ▶ **\*\*Mirror Descent minimizes cumulative loss\*\*** via **regularized updates**.
- ▶ The **Bregman divergence** controls how much each step deviates.
- ▶ The regret bound depends on the choice of **\*\*regularizer  $R(w)$ \*\***.

**Implication:** Choosing an appropriate  $R(w)$  determines the efficiency of the algorithm.

## Polynomial Potential

- ▶ Potential:  $R_p(u) = \frac{1}{2} \|u\|_p^2 = \frac{1}{2} \left( \sum_{i=1}^d u_i^p \right)^{2/p}$
- ▶ Dual Potential  $R_p^* = R_q$  Where  $\frac{1}{p} + \frac{1}{q} = 1$
- ▶ Euclidean norm:  $q = p = 2$
- ▶ Suppose the sequence of examples  $(x_1, y_1), \dots, (x_T, y_T)$  satisfies  $\|x_t\|_p \leq X_p$  for all  $1 \leq t \leq T$
- ▶ Suppose we use the dual descend algorithm for the potential function  $R_p$  and the learning rate  $\lambda = \frac{2\epsilon}{(p-1)X_p^2}$  for some  $0 < \epsilon < 1$
- ▶ Loss Bound:
 
$$L_{A,T} \leq \frac{L_T(u)}{1-\epsilon} + \frac{\|u\|_q^2}{\epsilon(1-\epsilon)} \times \frac{(p-1)X_p^2}{4}$$

## Exponential Potential

- ▶ Potential:  $R(u) = \sum_{i=1}^d e^{u_i}$
- ▶ Dual Potential  $R^*(u) = \sum_{i=1}^d u_i (\ln u_i - 1)$
- ▶ Euclidean norm:  $q = p = 2$
- ▶ Suppose the sequence of examples  $(x_1, y_1), \dots, (x_T, y_T)$  satisfies  $\|x_t\|_\infty \leq X_p$  for all  $1 \leq t \leq T$
- ▶ Suppose we use the dual descend algorithm for the exponential potential function  $R$  and the learning rate  $\lambda = \frac{2\epsilon}{X_\infty^2}$  for some  $0 < \epsilon < 1$
- ▶ Loss Bound:  
$$L_{A,T} \leq \frac{L_T(u)}{1-\epsilon} + \frac{X_\infty^2 \ln d}{2\epsilon(1-\epsilon)}$$

## Lemma 2.20: Regret Bound for OMD

**Lemma 2.20.** Suppose that OMD is run with a link function  $g = \nabla R^*$ . Then, its regret is upper bounded by:

$$\sum_{t=1}^T \langle w_t - u, z_t \rangle \leq R(u) - R(w_1) + \sum_{t=1}^T D_{R^*}(-z_{1:t} \| -z_{1:t-1}).$$

Furthermore, equality holds for the vector  $u$  that minimizes  $R(u) + \sum_t \langle u, z_t \rangle$ .