# Dual Descent

February 20, 2025

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

# Follow-the-Regularized-Leader (FTRL)

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

▶ Regularizer controls length of weight vector -¿ changes from iteration to iteration.

# Follow-the-Regularized-Leader (FTRL)

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

▶ Regularizer controls length of weight vector -¿ changes from iteration to iteration.

▶ Each step requires solving a constrained minimization problem.

# Review: Property of FoRel Algorithm

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

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

# Review: One step of Gradient Descent using strongly convex regularizer

**Lemma 2.10:**

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

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

# Main Theorem regarding FoReL using stongly convex regularizer

Let $f_1, \ldots, f_T$ be a sequence of convex functions with the following conditions:

- $f_t$ is $L_t$-Lipschitz with respect to some norm $\|\cdot\|$.

Then, for all $u \in S$,

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

# Main Theorem regarding FoReL using stongly convex regularizer

Let $f_1, \ldots, f_T$ be a sequence of convex functions with the following conditions:

- $f_t$ is $L_t$-Lipschitz with respect to some norm $\| \cdot \|$.

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

Then, for all $u \in S$,

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

# Main Theorem regarding FoReL using stongly convex regularizer

Let $f_1, \ldots, f_T$ be a sequence of convex functions with the following conditions:

- $f_t$ is $L_t$-Lipschitz with respect to some norm $\| \cdot \|$.

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

- FoReL is run on the sequence with a regularization function that is $\sigma$-strongly-convex with respect to the same norm.

Then, for all $u \in S$,

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

# Potential based gradient Descent

▶ $\text{Regret}_t = $ Regret vector $\text{Regret}_t(w) = L_{A,t} - L_t(w)$

# Potential based gradient Descent

- $\mathrm{Regret}_t =$ Regret vector $\mathrm{Regret}_t(\mathrm{w}) = L_{A,t} - L_t(\mathrm{w})$
- $\mathrm{Regret}_t =$ State of prediction algorithm at time $t$

# Potential based gradient Descent

▶ $\text{Regret}_t$ = Regret vector $\text{Regret}_t(\text{w}) = L_{A,t} - L_t(\text{w})$

▶ $\text{Regret}_t$ = State of prediction algorithm at time $t$

▶ Potential/Regularizer: $R(\text{Regret})$ Quantifies badness of the state.

# Potential based gradient Descent

- $\text{Regret}_t =$ Regret vector $\text{Regret}_t(\text{w}) = L_{A,t} - L_t(\text{w})$
- $\text{Regret}_t =$ 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.

# 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 $\ell_t$

# 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 $\ell_t$
- $w_t = \nabla R(\text{Regret}_t)$ is a good choice for the interior of $S$.

# 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 $\ell_t$
- $w_t = \nabla R(\text{Regret}_t)$ is a good choice for the interior of $S$.
- For finite number of experts, $\text{Regret}_t$ is finite dimensional and we can compute $w_t$ explicitly.

# Potential based gradient Descent

- $\text{Regret}_t = $ Regret vector $\text{Regret}_t(\mathsf{w}) = L_{A,t} - L_t(\mathsf{w})$
- $\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 $R(\text{Regret}_{t+1}) - R(\text{Regret}_t) + \mathsf{w}_t \cdot \ell_t$ is small for all $\ell_t$
- $\mathsf{w}_t = \nabla R(\text{Regret}_t)$ is a good choice for the interior of $S$.
- For finite number of experts, $\text{Regret}_t$ is finite dimensional and we can compute $\mathsf{w}_t$ explicitly.
- Here, $\text{Regret} = \{R(\mathsf{w})\}_{\mathsf{w} \in \mathbb{R}^d}$ is uncountably infinite.

# 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 $\ell_t$

- $w_t = \nabla R(\text{Regret}_t)$ is a good choice for the interior of $S$.

- 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).

# 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 $\ell_t$
- $w_t = \nabla R(\text{Regret}_t)$ is a good choice for the interior of $S$.
- 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 $z_{1:t} = \sum_{i=1}^{t} z_i$, the FoReL update rule can be written as

$$w_{t+1} = \arg\min_{w \in S} R(w) + \sum_{i=1}^{t} \langle w, z_i \rangle$$

$$= \arg\min_{w \in S} R(w) + \langle w, z_{1:t} \rangle$$

$$= \arg\max_{w \in S} \langle w, -z_{1:t} \rangle - R(w).$$

# Mirror Descent Update for linear functions

Update rule
$$w_{t+1} = \arg\max_{w \in S} \langle w, -z_{1:t} \rangle - R(w).$$

Link Function:
$$g(\theta) = \arg\max_{w \in S} \langle w, \theta \rangle - R(w),$$

Update rule can be re-written as

1. $\theta_0 = 0$

# Mirror Descent Update for linear functions

Update rule

$$w_{t+1} = \arg \max_{w \in S} \langle w, -z_{1:t} \rangle - R(w).$$

Link Function:

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

Update rule can be re-written as

1. $\theta_0 = 0$
2. $\theta_{t+1} = \theta_t - z_t$

# Mirror Descent Update for linear functions

Update rule

$$w_{t+1} = \arg\max_{w \in S} \langle w, -z_{1:t} \rangle - R(w).$$

Link Function:

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

Update rule can be re-written as

1. $\theta_0 = 0$

2. $\theta_{t+1} = \theta_t - z_t$

3. $w_{t+1} = g(\theta_{t+1})$

   **Identical** update to FTRL for linear loss functions.
   What about general convex loss functions?

## Sub-Gradients

- ▶ we can reduce general convex to linear using the gradient.
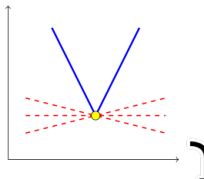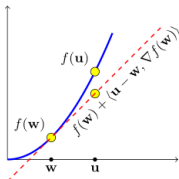
## Sub-Gradients

- ▶ we can reduce general convex to linear using the gradient.
- ▶ What can we do if $f(x)$ is convex but not differentiable at $x$?

## 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 $\forall y, l(y) \leq f(x)$ and $l(x) = f(x)$

## 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 $\forall y, l(y) \leq f(x)$ 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 $f_1, f_2, \ldots$ are convex (but not necessarily differentiable). Let $\eta$ be the learning rate.

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

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

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

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

$$\leq \frac{B^2}{2\eta} + \eta T L^2$$

# Gradient based Online Mirror Descent (OMD)

**parameter:** a link function $g : \mathbb{R}^d \to S$
**initialize:** $\theta_1 = 0$
**for** $t = 1, 2, \ldots$

▶ **project** $w_t = g(\theta_t)$

Dual Decent: Instead of minimizing $f$, minimize $\nabla f$.
Convexity implies equivalence of goals.

# Gradient based Online Mirror Descent (OMD)

**parameter:** a link function $g : \mathbb{R}^d \to S$
**initialize:** $\theta_1 = 0$
**for** $t = 1, 2, \ldots$

- ▶ **project** $w_t = g(\theta_t)$

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

Dual Decent: Instead of minimizing $f$, minimize $\nabla f$.
Convexity implies equivalence of goals.

# Duality

- OMD can be analyzed using elementary tools from duality.

# Duality

- OMD can be analyzed using elementary tools from duality.

- Using Duality Gives better intuition, more general analysis, tighter bounds.

# Dual Vector Spaces

- $V$ is a vector space, with a norm $\|v\|$

# Dual Vector Spaces

- $V$ is a vector space, with a norm $\|v\|$
- $U$ is the set of all linear mappings from $V$ to $V$

# Dual Vector Spaces

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

$$\|u\|_* \doteq \max_{v \in V} \frac{\|u(v)\|}{\|v\|}$$

# Dual Vector Spaces

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

$$\|u\|_* \doteq \max_{v \in V} \frac{\|u(v)\|}{\|v\|}$$

- $V$ is equivalent to the set of all linear mappings from $U$ to $U$.

# Dual Vector Spaces

- $V$ is a vector space, with a norm $\|v\|$

- $U$ is the set of all linear mappings from $V$ to $V$

- The norm of $u \in U$ is defined as

$$\|u\|_* \doteq \max_{v \in V} \frac{\|u(v)\|}{\|v\|}$$

- $V$ is equivalent to the set of all linear mappings from $U$ to $U$.

- $U$ and $V$ are dual vector spaces, with dual norms.

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$
- The linear operation is the dot product $\mathbf{u} \cdot \mathbf{v}$

# Dual Norms

- ▶ The space is always $U, V = \mathbb{R}^n$
- ▶ The linear operation is the dot product $u \cdot v$
- ▶ $\|u\|_*$ is the dual norm to $\|v\|$,

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$

- The linear operation is the dot product $u \cdot v$

- $\|u\|_*$ is the dual norm to $\|v\|$,

- $L_2$ norm: $\sqrt{\sum_{i=1}^n x_i^2}$

# Dual Norms

▶ The space is always $U, V = \mathbb{R}^n$

▶ The linear operation is the dot product $\mathbf{u} \cdot \mathbf{v}$

▶ $\|u\|_*$ is the dual norm to $\|v\|$,

▶ $L_2$ norm: $\sqrt{\sum_{i=1}^n x_i^2}$

▶ $L_1$ norm: $\sum_{i=1}^n |x_i|$

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$

- The linear operation is the dot product $u \cdot v$

- $\|u\|_*$ is the dual norm to $\|v\|$,

- $L_2$ norm: $\sqrt{\sum_{i=1}^n x_i^2}$

- $L_1$ norm: $\sum_{i=1}^n |x_i|$

- $L_\infty$ norm: $\max_i |x_i|$

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$

- The linear operation is the dot product $\mathbf{u} \cdot \mathbf{v}$

- $\|u\|_*$ is the dual norm to $\|v\|$,

- $L_2$ 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: $\left(\sum_{i=1}^{n} x_i^p\right)^{\frac{1}{p}}$

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$

- The linear operation is the dot product $u \cdot v$

- $\|u\|_*$ is the dual norm to $\|v\|$,

- $L_2$ 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: $\left( \sum_{i=1}^{n} x_i^p \right)^{\frac{1}{p}}$

- $L_p, L_q$ are dual norms if $p, q \geq 1$, and $\frac{1}{p} + \frac{1}{q} = 1$

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$

- The linear operation is the dot product $u \cdot v$

- $\|u\|_*$ is the dual norm to $\|v\|$,

- $L_2$ 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: $\left(\sum_{i=1}^{n} x_i^p\right)^{\frac{1}{p}}$

- $L_p, L_q$ are dual norms if $p, q \geq 1$, and $\frac{1}{p} + \frac{1}{q} = 1$

- $L_1, L_\infty$ are dual.

# Dual Norms

- The space is always $U, V = \mathbb{R}^n$
- The linear operation is the dot product $\mathbf{u} \cdot \mathbf{v}$
- $\|u\|_*$ is the dual norm to $\|v\|$,
- $L_2$ 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: $\left(\sum_{i=1}^n x_i^p\right)^{\frac{1}{p}}$
- $L_p, L_q$ 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.

# Lipschitz condition and the dual norm

**Lemma 2.6:**
Let $f : S \to \mathbb{R}$ be a convex function. Then, $f$ is $L$-Lipschitz over $S$ with respect to a norm $\| \cdot \|$ if and only if for all $w \in S$ and $z \in \partial f(w)$ we have:

$$\|z\|_* \leq L$$

where $\| \cdot \|_*$ denotes the dual norm.

## Proof of Lemma 2.6

**Proof:**

Assume that $f$ is $L$-Lipschitz. For any $w \in S$ and $z \in \partial f(w)$, choose $u$ such that $u - w = \arg\max_{\|v\|=1} \langle v, z \rangle$. Then,

$$\langle z, u - w \rangle = \|z\|_*$$

By the sub-gradient definition,

$$f(u) - f(w) \geq \langle z, u - w \rangle = \|z\|_*$$

Since $f$ is $L$-Lipschitz,

$$f(u) - f(w) \leq L\|u - w\| = L$$

Combining the inequalities:

$$\|z\|_* \leq L$$

For the converse, assume $\|z\|_* \leq L$ for all $z \in \partial f(w)$. Then,

$$f(u) - f(w) \leq \langle z, u - w \rangle \leq \|z\|_* \|u - w\| \leq L\|u - w\|$$

Hence, $f$ is $L$-Lipschitz.

# Fenchel Duality

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

# Fenchel Duality

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

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

# Fenchel Duality

- Suppose $F : A \to \mathbb{R}$ 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 Equivalent to Legendre duality for differentiable functions.

# Visualization of the Legendre Dual

# Visualization of the Legendre Dual

# Visualization of the Legendre Dual
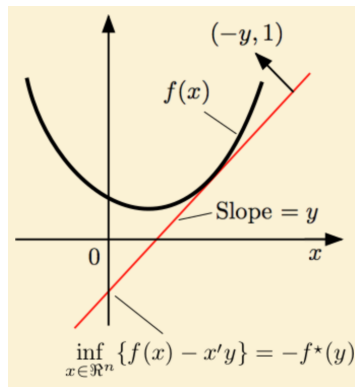
- $x, y \mathbb{R}$

# Visualization of the Legendre Dual

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

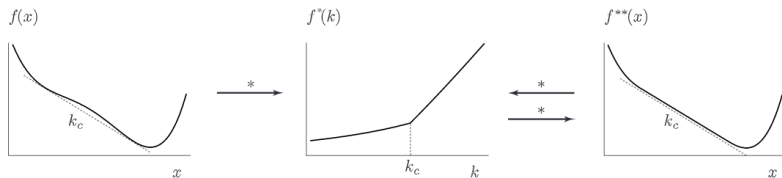# Visualization of the Legendre Dual

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

# Visualization of the Legendre Dual

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



$(-y, 1)$
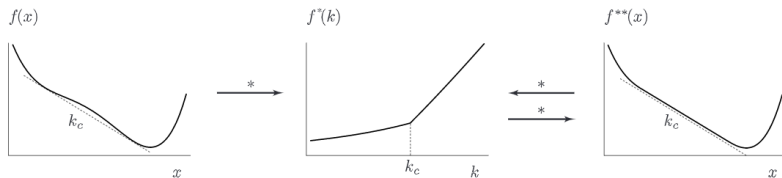
$f(x)$

Slope $= y$

$0$

$x$

$\inf_{x \in \mathbb{R}^n} \{f(x) - x'y\} = -f^\star(y)$

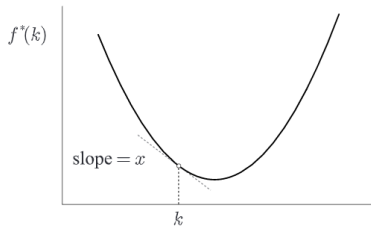## Dual of Dual

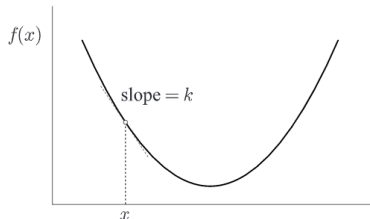- The dual of any function is convex.

# Dual of Dual

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

# Gradient Duality (legendre only)

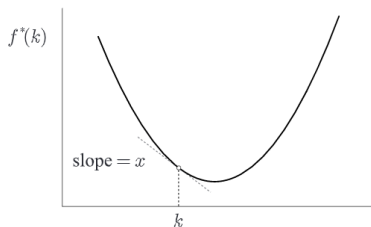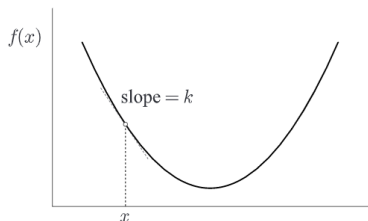▶ If the gradient of $f$ at $x$ is $k$ then
the gradient of $f^*$ at $k$ is $x$

# Gradient Duality (legendre only)

► 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: $F(u) = \sum_{i=1}^{d} e^{u_i}$

# Example: Exponential Potential

- Potential: $F(u) = \sum_{i=1}^{d} e^{u_i}$
- Gradient: $\nabla F(u)_i = e^{u_i}$ or $\nabla F(u) = F(u)$.

# Example: Exponential Potential

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

# Example: Exponential Potential

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

# Example: Exponential Potential

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

# Bregman Divergence

# Bregman Divergence

# Bregman Divergence

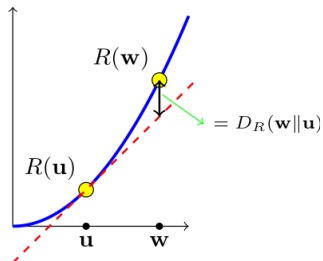- $R(x)$ is convex and differentiable.

# Bregman Divergence

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

# Bregman Divergence

- $R(x)$ is convex and differentiable.

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

- The error term of the first order Taylor expansion around $u$

# Bregman Divergence

- $R(x)$ is convex and differentiable.

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

- The error term of the first order Taylor expansion around $u$

# Fenchel and Bregman

- $F$: strictly convex with continuous first derivative.

## Fenchel and Bregman

- $F$: strictly convex with continuous first derivative.
- $F^*$ is the Fenchel Dual of $F$

# Fenchel and Bregman

- $F$: strictly convex with continuous first derivative.
- $F^*$ is the Fenchel Dual of $F$
- $D_F, D_{F^*}$ Bregman divergences wrt $F, F^*$

# 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 $v' = \nabla F(v)$

# 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 - Step 1

**Gradient Step in Dual Space:**

$$z_{t+1} = \nabla R(w_t) - \eta \nabla f_t(w_t)$$

Here, $\nabla R(w_t)$ maps the point into the dual space.

# Mirror Descent - Step 2

**Projection Back to Primal Space:**

$$w_{t+1} = \arg \min_{w \in S} D_R(w, z_{t+1})$$

Where $D_R(w, z)$ is the Bregman divergence:

$$D_R(w, z) = R(w) - R(z) - \langle \nabla R(z), w - z \rangle$$

This projection ensures $w_{t+1}$ stays within the feasible set $S$.

# Mirror Descent (alternative Notation)

▶ Gradient descent in dual space $\theta_t = \theta_{t-1} - \lambda \nabla \ell_t(\theta_{t-1})$

# Mirror Descent (alternative Notation)

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

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

# Mirror Descent (alternative Notation)

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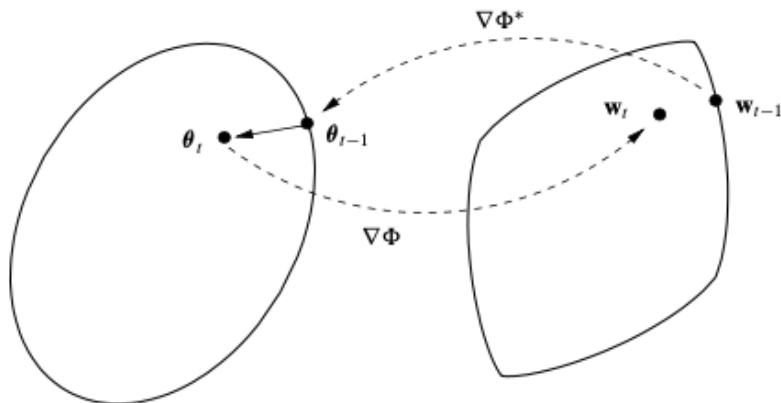
▶ Using duality can be rewritten as

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

▶ Projection: 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

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

# 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 $R(u) + \sum_t \langle u, z_t \rangle$.

# Proof: Step 1 - Fenchel–Young Inequality

Using the **Fenchel–Young inequality**, we have:

$$R(\mathsf{u}) + \sum_{t=1}^{T} \langle \mathsf{u}, \mathsf{z}_t \rangle = R(\mathsf{u}) - \langle \mathsf{u}, -\mathsf{z}_{1:T} \rangle \geq -R^*(-\mathsf{z}_{1:T}).$$

Equality holds for u that maximizes $\langle \mathsf{u}, -\mathsf{z}_{1:T} \rangle - R(\mathsf{u})$, hence minimizing $R(\mathsf{u}) + \langle \mathsf{u}, \mathsf{z}_{1:T} \rangle$.

# Proof: Step 2 - Bregman Divergence

Since $w_t = \nabla R^*(-z_{1:t-1})$ and using the definition of the Bregman divergence, we rewrite:

$$-R^*(-z_{1:T}) = -R^*(0) - \sum_{t=1}^{T} \left( R^*(-z_{1:t}) - R^*(-z_{1:t-1}) \right).$$

Rearranging, we get:

$$= -R^*(0) + \sum_{t=1}^{T} (\langle w_t, z_t \rangle - D_{R^*}(-z_{1:t} \| -z_{1:t-1})).$$

# Final Step

**Note:** Since

$$R^*(0) = \max_{\mathsf{w}} \langle 0, \mathsf{w} \rangle - R(\mathsf{w}) = -\min_{\mathsf{w}} R(\mathsf{w}) = -R(\mathsf{w}_1),$$

combining all the above, we conclude the proof. $\square$

# OMD for $\ell_2$

- ▶ Problem: OCO where $\forall z \in \partial f(x)\ \|z\|_2 \leq X_2$

# OMD for $\ell_2$

- Problem: OCO where $\forall z \in \partial f(x)$ $\|z\|_2 \leq X_2$
- The dual norm for the weights is $\|w\|_2$

# OMD for $\ell_2$

- Problem: OCO where $\forall z \in \partial f(x)\ \|z\|_2 \leq X_2$
- The dual norm for the weights is $\|w\|_2$
- We use the dual descend algorithm for the half quadratic regularizer

# OMD for $\ell_2$

- Problem: OCO where $\forall z \in \partial f(x) \ \|z\|_2 \leq X_2$
- The dual norm for the weights is $\|w\|_2$
- We use the dual descend algorithm for the half quadratic regularizer
- Regularizer: $R(u) = \frac{1}{2}\|u\|_2^2$

# OMD for $\ell_2$

- Problem: OCO where $\forall z \in \partial f(x) \; \|z\|_2 \leq X_2$

- The dual norm for the weights is $\|w\|_2$

- We use the dual descend algorithm for the half quadratic regularizer

- Regularizer: $R(u) = \frac{1}{2}\|u\|_2^2$

- Legendre Dual Regularizer $R^*(u) = \frac{1}{2}\|u\|_2^2$

# OMD for $\ell_2$

- Problem: OCO where $\forall z \in \partial f(x) \; \|z\|_2 \le X_2$
- The dual norm for the weights is $\|w\|_2$
- We use the dual descend algorithm for the half quadratic regularizer
- Regularizer: $R(u) = \frac{1}{2}\|u\|_2^2$
- Legendre Dual Regularizer $R^*(u) = \frac{1}{2}\|u\|_2^2$
- **Gradient Mapping:** $\nabla R(w) = w$

# OMD for $\ell_2$

- Problem: OCO where $\forall z \in \partial f(x) \, \|z\|_2 \leq X_2$
- The dual norm for the weights is $\|w\|_2$
- We use the dual descend algorithm for the half quadratic regularizer
- Regularizer: $R(\mathbf{u}) = \frac{1}{2}\|\mathbf{u}\|_2^2$
- Legendre Dual Regularizer $R^*(\mathbf{u}) = \frac{1}{2}\|\mathbf{u}\|_2^2$
- **Gradient Mapping:** $\nabla R(w) = w$
- Gradient step: $z_{t+1} = w_t - \eta \nabla f_t(w_t)$

# Projection Step for $\ell_2$ Norm

**Bregman Divergence:**

$$D_R(w, z) = \frac{1}{2}\|w - z\|_2^2$$

**Projection Back to Primal Space:**

$$w_{t+1} = \Pi_S(z_{t+1}) = \arg\min_{w \in S} \frac{1}{2}\|w - z_{t+1}\|_2^2$$

Where $\Pi_S$ denotes the Euclidean projection onto the feasible set $S$.

# Final Update Rule for $\ell_2$ Norm

Combining both steps, the final update rule becomes:

$$w_{t+1} = \Pi_S \left( w_t - \eta \nabla f_t(w_t) \right)$$

This is equivalent to the standard **Projected Gradient Descent** for the $\ell_2$ norm.

# Optimal Tuning for $\eta$ and Regret Bound

**Regret Bound:**

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

Assuming $\|u\|_2 \leq B$ and $\|\nabla f_t(w_t)\|_2 \leq L$, this simplifies to:

$$\text{Regret}_T(u) \leq \frac{B^2}{2\eta} + \frac{\eta L^2 T}{2}$$

**Optimal $\eta$:**

$$\eta^* = \frac{B}{L\sqrt{T}}$$

**Resulting Regret Bound:**

$$\text{Regret}_T(u) \leq BL\sqrt{T}$$

# $\ell_\infty$ OMD

- Problem: OCO where $\forall z \in \partial f(x)$ $\|z\|_\infty \leq X_\infty$

# $\ell_\infty$ OMD

- Problem: OCO where $\forall z \in \partial f(x) \; \|z\|_\infty \le X_\infty$
- The dual norm for the weights is $\|w\|_1$

# $\ell_\infty$ OMD

- ▶ Problem: OCO where $\forall z \in \partial f(x) \ \|z\|_\infty \leq X_\infty$

- ▶ The dual norm for the weights is $\|w\|_1$

- ▶ Suppose we use the dual descend algorithm for the exponential regularizer

# $\ell_\infty$ OMD

- Problem: OCO where $\forall z \in \partial f(x) \; \|z\|_\infty \leq X_\infty$

- The dual norm for the weights is $\|w\|_1$

- Suppose we use the dual descend algorithm for the exponential regularizer

- Regularizer: $R(u) = \frac{1}{\eta} \ln \sum_{i=1}^{d} e^{\eta u_i}$

# $\ell_\infty$ OMD

- Problem: OCO where $\forall z \in \partial f(x)\ \|z\|_\infty \leq X_\infty$

- The dual norm for the weights is $\|w\|_1$

- Suppose we use the dual descend algorithm for the exponential regularizer

- Regularizer: $R(u) = \frac{1}{\eta} \ln \sum_{i=1}^{d} e^{\eta u_i}$

- Legendre Dual Regularizer $R^*(u) = \sum_{i=1}^{d} u_i (\ln u_i - 1)$

▶ Weight update:

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

▶ Weight update:

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

▶ Unnormalized exponentiated gradient:

$$w_{i,t} = w_{i,t-1} e^{-\lambda \nabla \ell_t(w_i - 1)}$$

- Weight update:

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

- Unnormalized exponentiated gradient:

$$w_{i,t} = w_{i,t-1} e^{-\lambda \nabla \ell_t(w_i - 1)}$$

- Normalized exponentiated gradient:

$$w_{i,t} = \frac{w_{i,t-1} e^{-\lambda \nabla \ell_t(w_i - 1)}}{\sum_{j=1}^{d} w_{j,t-1} e^{-\lambda \nabla \ell_t(w_j - 1)}}$$

▶ Normalization corresponds to projection on the simplex using the Bregman divergence according to $R^*$.

▶ Normalization corresponds to projection on the simplex using the Bregman divergence according to $R^*$.

▶ The dual descend algorithm for the exponential regularizer function $R$ and the learning rate $\lambda = \frac{2\epsilon}{X_\infty^2}$ for some $0 < \epsilon < 1$

- Normalization corresponds to projection on the simplex using the Bregman divergence according to $R^*$.

- The dual descend algorithm for the exponential regularizer function $R$ and the learning rate $\lambda = \frac{2\epsilon}{X_\infty^2}$ for some $0 < \epsilon < 1$

- yields Loss Bound:

$$L_{A,T} \leq \frac{L_T(u)}{1 - \epsilon} + \frac{X_\infty^2 \ln d}{2\epsilon(1 - \epsilon)}$$