

1. Conjugate Duality.

V. 2.2
2022.04.06

A. Definition. Let $F: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$. Define $F^*: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{-\infty, +\infty\}$.

by

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

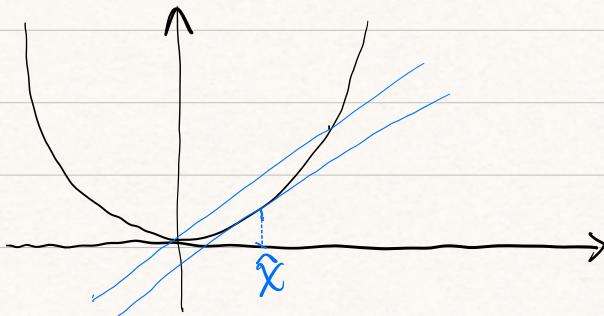
This is the convex conjugate of f .

B. Intuition:

① $\forall x \in \mathbb{R}^n$, 过点 $(x, f(x))$ 作斜率的直线, 那么

$$\hat{x} = \underset{x \in \mathbb{R}^n}{\arg \max} y^T x - f(x) \text{ 为让直线的纵截距最小}$$

的 x .



② \hat{x} 为让直线 $y = y^T x$ 与 $f(x)$ 交最多的那个 x .

C. Example

Ex. 1.2. Let $f(x) = cx$ for some $c \in \mathbb{R}^d$. Then $f^* = \delta_{\{c\}}$, i.e.

$$f^*(y) = \begin{cases} +\infty, & y \neq c \\ 0, & y = c. \end{cases}$$

Ex. 1.5 (Negative Entropy). Define $f: \mathbb{R}_+^n \rightarrow \mathbb{R}$ by

$f(x) = \sum_{i=1}^d x_i \ln x_i$. Then $f^*(y) = \sum_{i=1}^d e^{y_i - 1}$ by Proposition 1.9.

Ex. 1.6 Let $\|\cdot\|$ be a norm on \mathbb{R} , and let $f(x) = \frac{1}{2}\|x\|^2$. Then
 $f^*(x) = \frac{1}{2}\|x\|_*^2$.

D. Properties.

Proposition 1.7 [Young-Fenchel Inequality]. $\forall x, y \in \mathbb{R}^n$.

$$f^*(y) + f(x) \geq y^T x.$$

$\text{epi}(f)$ is closed



Proposition 1.8. Regardless of whether f is, f^* is closed and convex.

Proposition 1.9. (Conjugate of separable function).

Let $f: \mathbb{R}^a \times \mathbb{R}^b \rightarrow \mathbb{R}$ be defined by $f(x_1, x_2) = f_1(x_1) + f_2(x_2)$.
 Then $f^*(y_1, y_2) = f_1^*(y_1) + f_2^*(y_2)$.

Proposition 1.10. If f is closed and convex, then $f^{**} = f$.

Proposition 1.11. If f is closed and convex then the following are equivalent.

(a) $y \in \partial f(x)$.

* $f(x)$ 是关于 primal variable x 的函数.

(b) $f(x) + f^*(y) = y^T x$.

* $f^*(y)$ 是关于 dual variable y 的

(c) $x \in \partial f^*(y)$.

函数.

Proof. 1) $(a) \Rightarrow (b)$.

$$y^T x - f(x) \geq y^T u - f(u), \forall u \in \mathbb{R}^n \text{ (Definition of Subgradient)}$$
$$\Rightarrow y^T x - f(x) \geq \sup_{u \in \mathbb{R}^n} y^T u - f(u) = f^*(y). \quad (\text{P})$$

$$f^*(y) \geq y^T x - f(x) \quad (\text{Proposition 1.7}) \quad (2)$$

$$(1)(2) \Rightarrow y^T x = f^*(y) + f(x).$$

2) $(b) \Rightarrow (c)$.

Let $g = f^*$, then $g^* = f$.

$$x \in \partial g(y) \sim \forall v \in \mathbb{R}^n, g(v) - x^T v \geq g(y) - x^T y.$$

$$\text{Since } x^T y - f^*(y)$$

$$= x^T y - g(y)$$

$$= f(x)$$

$$= g^*(x)$$

$$= \sup_{v \in \mathbb{R}^n} x^T v - g(v).$$

$$\Rightarrow x^T y - g(y) \geq x^T v - g(v), \quad \forall v \in \mathbb{R}^n.$$

3) $(c) \Rightarrow (a)$.

Let $g = f^*$.

$$\Rightarrow g^*(x) + g(y) = y^T x \quad (\text{by } (c) \Rightarrow (b))$$

$$\Rightarrow y \in \partial g^*(x) = \partial f(x). \quad (\text{by } (b) \Rightarrow (c))$$

III

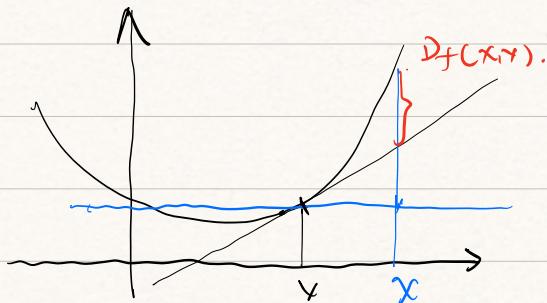
2. Bregman Divergence

A Definition 2.1. The Bregman divergence is defined to be

$$D_f(x, y) = f(x) - f(y) - \langle \nabla f(y), x - y \rangle.$$

B. Intuition.

The difference between the value of function and the first order estimation.



C. Example.

Ex. 2.2. Define $f: \mathbb{R}^n \rightarrow \mathbb{R}$ by $f(x) = \|x\|_2^2$. Then

$$D_f(x, y) = \|x - y\|_2^2.$$

Ex. 2.3 (Negative Entropy).

Define $f: \mathbb{R}_+^n \rightarrow \mathbb{R}$ by $f(x) = \sum_{i=1}^d x_i \ln x_i$. Then

$$D_f(x, y) = D_{KL}(x, y),$$

where $D_{KL}(x, y)$ is the generalized relative entropy.

Proof. $D_f(x, y) = f(x) - f(y) - \langle \nabla f(y), x - y \rangle$

$$= -H(x) + H(y) - \sum_{i=1}^d (1 + \ln y_i) \cdot (x_i - y_i)$$

$$\begin{aligned}
 &= -H(x) + H(y) - \sum_{i=1}^d x_i - \sum_{i=1}^d x_i \cdot \ln y_i + \sum_{i=1}^d y_i + \sum_{i=1}^d y_i \ln y_i \\
 &= H(x, y) - H(x) - \sum_{i=1}^d x_i + \sum_{i=1}^d y_i \\
 &= D_{KL}(x, y).
 \end{aligned}$$

□

Proposition 2.4.

Negative entropy is 1-strongly-convex with respect to L_1 -norm.
Proof.

Theorem [Pinsker's Inequality]

For any distribution p, q , $D_{KL}(p, q) \geq \frac{1}{2} \|p - q\|_1^2$.

$$\begin{aligned}
 D_f(p, q) &= f(p) - f(q) - \langle \nabla f(q), p - q \rangle \\
 &= D_{KL}(p, q) \\
 &\geq \frac{1}{2} \|p - q\|_1^2.
 \end{aligned}$$

□

D. Properties.

Proposition 2.8.

$D_f(x, y)$ is convex in x .

Proposition 2.9.

Let f be closed, convex and differential. Fix any $x, y \in \mathbb{X}$,

Define $\hat{x} = \nabla f(x)$, $\hat{y} = \nabla f(y)$. Then

$$\nabla f^*(\hat{x}) = x. \quad \textcircled{1}$$

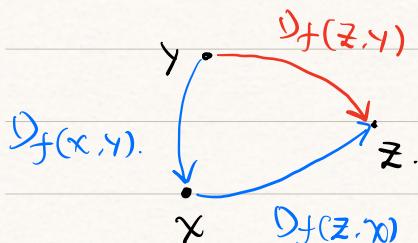
$$D_f(x, y) = D_{f^*}(\hat{y}, \hat{x}). \quad \textcircled{2}$$

Proof of $\textcircled{2}$.

$$\begin{aligned} D_{f^*}(\hat{y}, \hat{x}) &= f^*(\hat{y}) - f^*(\hat{x}) - \langle \nabla f^*(\hat{x}), \hat{y} - \hat{x} \rangle \\ &= \hat{y}^\top y - f(y) - (\hat{x}^\top x - f(x)) - \langle x, \hat{y} - \hat{x} \rangle \\ &= \hat{y}^\top y - f(y) + f(x) - x^\top \hat{y} \\ &= f(x) - f(y) - \langle \nabla f(y), x - y \rangle. \\ &= D_f(x, y). \end{aligned}$$

\textcircled{2}

Lemma 2.10. [Generalized Pythagoras Identity].



$$D_f(z, x) + D_f(x, y) - D_f(z, y) = (\nabla f(y) - \nabla f(x))^\top (z - x)$$

Proof.

$$\begin{aligned} \text{LHS} &= f(z) - f(x) - \langle \nabla f(x), z - x \rangle + f(x) - f(y) - \langle \nabla f(y), x - y \rangle \\ &\quad - (f(z) - f(y)) - \langle \nabla f(y), z - y \rangle \\ &= \langle \nabla f(y), z - x \rangle - \langle \nabla f(x), z - x \rangle \\ &= (\nabla f(y) - \nabla f(x))^\top (z - x). \end{aligned}$$

\textcircled{3}

Ex. 2.11

Consider the case of $f(x) = \|x\|_2^2$. Let $\vec{a} = z-x$, $\vec{b} = y-x$, $\vec{c} = z-y$.

Then

$$\vec{a}^2 + \vec{b}^2 - \vec{c}^2 = \nabla f(z, x)^T (\vec{a} + \vec{b}) - \nabla f(z, y) = (\nabla f(y) - \nabla f(x))^T (\vec{a} + \vec{b}),$$

which recovers the law of cosine.

Proposition. 2.12. $\nabla \nabla f(x, y) = \nabla f(x) - \nabla f(y)$

3. Projections.

Let $\mathbb{X} \subseteq \mathbb{R}^n$ be a closed, convex set. Assume that $f: \mathbb{X} \rightarrow \mathbb{R}$ is a strictly convex function, which implies that $D_f(x, y)$ is strictly convex in x .

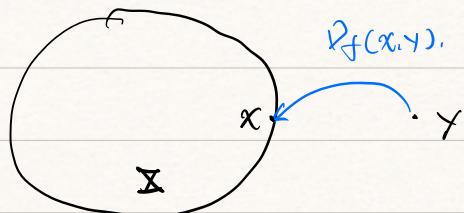
A. Definition 2.13.

The projection of y onto \mathbb{X} under the Bregman divergence is

$$\Pi_{\mathbb{X}}^f(y) = \underset{x \in \mathbb{X}}{\operatorname{argmin}} D_f(x, y),$$

where the minimizer is uniquely determined since $D_f(x, y)$ is strictly convex in x .

B.



C. Properties.

Proposition 2.14.

Suppose that f is differentiable. Fix any $x \in \mathbb{R}^n$ and let $\pi = \Pi_{\mathbb{X}}^f(y)$.
A $x \in \mathbb{X}$,

$$(\nabla f(\pi) - \nabla f(y))^T \cdot (x - \pi) \geq 0.$$

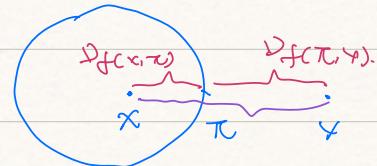
Proposition 2.15. [Anti-triangle Inequality].

Fix any $y \in \mathbb{R}^n$, and let $\pi = T_{\mathcal{X}}^+(y)$. Then $\forall x \in \mathcal{X}$

$$\underline{D_f(x, \pi)} + \underline{D_f(\pi, y)} \leq \underline{D_f(x, y)}.$$

Proof.

$$\begin{aligned} & D_f(x, \pi) + D_f(\pi, y) - D_f(x, y) \\ &= (\nabla f(y) - \nabla f(\pi))^T \cdot (x - \pi) \quad (\text{Generalized Pythagoras Identity}) \\ &= -\nabla D_f(\pi, y)^T \cdot (x - \pi) \\ &\leq 0. \quad (1^{\text{st}}\text{-order optimality condition}) \end{aligned}$$



□

Remark. [2.15.1]

In Proposition 2.15, the equality holds when \mathcal{X} is a hyperplane.

Proof. Since \mathcal{X} is a hyperplane, then there must exist $v \in \mathbb{R}^n$ s.t. $v^T(x - \pi) = 0$ (v is actually the normal vector of \mathcal{X}), together with the 1^{st} -order optimality condition in the proof of Prop. 2.15 concludes the proof.

□

4. The Mirror Descent Algorithm.

A. Motivations.

① Optimize the function f , which is Lipschitz with respect to a more general norm instead of L_2 -norm. Suppose f is L -Lipschitz w.r.t. $\|\cdot\|_\infty$, then f is $Td \cdot L$ -Lipschitz w.r.t. L_2 norm., but the factor Td might be undesirably large.

② In gradient descent, $x_{t+1} = x_t - \gamma \cdot \nabla f(x_t)$. However, x_t and $\nabla f(x_t)$ lie in different vector spaces respectively.

B. Main Ideas.

① **Mapping:** Find a mirror map Φ . Use its gradient $\nabla \Phi$ which is bijection and $\nabla \Phi$'s inverse $\nabla \Phi^*$ to map back and forth between primal and dual points

② **Ensuring Feasibility:** If the goal is to optimize f over a constraint set X which is convex, project that point onto X under the Bregman divergence D_Φ .

C. The mirror map Φ .

Consider the mirror map $\Phi: D \rightarrow \mathbb{R}$, where $D \subseteq \mathbb{R}^n$ is an open set.

A1: Φ is strictly convex and differential on all of D .

A2: The dual space of Φ is all of \mathbb{R}^n . That is, $\{\nabla \Phi(x) : x \in D\} = \mathbb{R}^n$.

A3: The gradient of Φ diverges on the boundary of D . That is,

$$\lim_{x \rightarrow \partial D} \|\nabla \Phi(x)\| = +\infty.$$

A function satisfying the above 3 assumptions is called Legendre function.

Proposition 3.2.

Suppose the assumptions A1 and A2 hold. Then $\nabla \Phi: D \rightarrow \mathbb{R}^n$ ^{is a} bijection.

Pruf.

1) Surjection:

By Assumption A2, $\forall y \in \mathbb{R}^n, \exists x \in D$ st. $\nabla \Phi(x) = y$.

2) injection:

Suppose $\exists x_1, x_2 \in D$ st. $\nabla \Phi(x_1) = \nabla \Phi(x_2) = y$. Then by

Proposition 1.11, $y \in \partial \Phi(x_1) \cap \partial \Phi(x_2)$, which contradicts with Proposition A.1.

■

D. The constraint set \mathcal{X} .

Typically we will optimize a constraint set \mathcal{X} , which is supposed to have the following assumptions:

A.4. \mathcal{X} is a closed convex set.

A.5. $\mathcal{X} \subseteq \bar{D}$, where \bar{D} is the closure of D .

A.6. $\mathcal{X} \cap D \neq \emptyset$.

E. Projection in MD

A primal point $y \in D$ could be projected back onto the constraint set S as

$$\Pi_{S \cap D}^{\Phi}(y) = \underset{x \in S \cap D}{\operatorname{argmin}} D_{\Phi}(x, y).$$

Proposition 3.3. Let $\mathbb{D}(x) = \sum_i x_i \ln x_i$. Suppose that $y \in D \setminus S$, i.e.,

$y \in \mathbb{R}_{>0}^n$ but $\|y\|_1 \neq 1$. Then $\Pi_{S \cap D}^{\Phi}(y) = \frac{y}{\|y\|_1}$.

Proof.

$$\underset{x \in S \cap D}{\operatorname{argmin}} D_{\Phi}(x, y) = \underset{x \in S \cap D}{\operatorname{argmin}} \sum_{i=1}^n (x_i \cdot \ln \frac{x_i}{y_i} - x_i + y_i) = \underset{x \in S \cap D}{\operatorname{argmin}} \sum_{i=1}^n x_i \cdot \ln \frac{x_i}{y_i}$$

$D_{\Phi}(x, y)$ is the generalized KL divergence $D_{KL}^G(x \| y)$ when Φ is the negative entropy.

Let $f(x) = x \cdot \ln x$. It's clear that $f(x)$ is convex function since $f''(x) = \frac{1}{x} > 0$. Then

→ Jensen's Ineq.

$$\begin{aligned} \sum_i x_i \cdot \ln x_i &= \sum_i y_i f\left(\frac{x_i}{y_i}\right) = \sum_i \|y\|_1 \cdot \frac{y_i}{\|y\|_1} \cdot f\left(\frac{x_i}{y_i}\right) \geq \|y\|_1 \cdot f\left(\frac{\sum_i y_i}{\|y\|_1} \cdot \frac{x_i}{y_i}\right) \\ &= \|y\|_1 \cdot f\left(\frac{1}{\|y\|_1}\right) = \ln \frac{1}{\|y\|_1}, \end{aligned}$$

where the equality holds iff $\forall i, j. \frac{x_i}{y_i} = \frac{x_j}{y_j}$. Therefore,

$$\underset{x \in S \cap D}{\operatorname{argmin}} D_{\Phi}(x, y) \text{ satisfies } x = \frac{y}{\|y\|_1}.$$

□

F. The MD Algorithm

- For $i = 1, 2, \dots$, do

- Incur cost $f_i(x_i)$, receive subgradient $g_i \in \partial f_i(x_i)$

- $\hat{x}_i \in \nabla \Phi(x_i)$ (map primal point to dual)

- $\hat{y}_{i+1} = \hat{x}_i - \gamma \cdot g_i$. (take a gradient descent step in the dual)

- $y_{i+1} = \nabla \Phi^*(\hat{y}_{i+1})$ (map new dual point back to a primal point in \mathcal{D})

- $x_{i+1} \in \Pi_{\mathcal{X} \cap \mathcal{D}}^{\Phi}(y_{i+1})$ (project new primal point onto \mathcal{X})

 : $y_{i+1} = \operatorname{argmin}_{y \in \mathcal{D}} \gamma \cdot f_i(x_i) + D_{\Phi}(y, x_i)$

 : $x_{i+1} = \operatorname{argmin}_{x \in \mathcal{X} \cap \mathcal{D}} D_{\Phi}(x, y_{i+1})$

Theorem 3.5

Given any $\gamma > 0$ and

- A mirror map $\Phi: \mathcal{D} \rightarrow \mathbb{R}$ satisfies A_1, A_2, A_3 where \mathcal{D} is an open subset of \mathbb{R}^n .

- A feasible set $\mathcal{X} \subseteq \mathbb{R}^n$ satisfying A_4, A_5, A_6 .

- Convex functions $f_1, f_2, \dots: \mathcal{X} \rightarrow \mathbb{R}$.

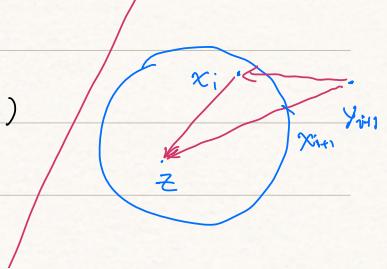
Furthermore, we assume Φ is φ -strongly w.r.t. norm $\|\cdot\|$. Then the MD algorithm satisfies

$$\begin{aligned} \sum_{i=1}^t (f_i(x_i) - f_i(x^*)) &\leq D_{\Phi}(x^*, x_1) + \sum_{i=1}^t (\langle g_i, x_i - y_{i+1} \rangle - \frac{\varphi}{2\gamma} \|x_i - y_{i+1}\|^2) \\ &\leq \frac{D_{\Phi}(x^*, x_1)}{\gamma} + \frac{\gamma}{2\varphi} \sum_{i=1}^t \|g_i\|_*^2. \end{aligned}$$

Remark. In linear opt. setting.

Proof. $\forall z \in \mathcal{X}$,

$$\begin{aligned} f_i(x_i) - f_i(z) &\leq g_i^\top (x_i - z) \quad \xrightarrow{\text{def.}} \hat{x}_i \quad \xrightarrow{\text{def.}} \hat{y}_{i+1} \\ &= \frac{1}{\gamma} \cdot (\nabla \Phi(x_i) - \nabla \Phi(y_{i+1}))^\top (x_i - z) \\ &= \frac{1}{\gamma} \cdot (\nabla \Phi(y_{i+1}) - \nabla \Phi(x_i))^\top (z - x_i) \end{aligned}$$



$$= \frac{1}{\eta} \cdot (\mathcal{D}_{\bar{\Phi}}(z, x_i) + \mathcal{D}_{\bar{\Phi}}(x_i, y_{i+1}) - \mathcal{D}_{\bar{\Phi}}(z, y_{i+1})). \quad (\text{GPJ}).$$

$$\leq \frac{1}{\eta} (\mathcal{D}_{\bar{\Phi}}(z, x_i) + \mathcal{D}_{\bar{\Phi}}(x_i, y_{i+1}) - \mathcal{D}_{\bar{\Phi}}(z, x_{i+1})). \quad (\text{Prop. 2.15})$$

$$\sum_{i=1}^t (f_i(x_i) - f(z)) \leq \frac{1}{\eta} \left[\mathcal{D}_{\bar{\Phi}}(z, x_1) + \sum_{i=1}^t \mathcal{D}_{\bar{\Phi}}(x_i, y_{i+1}) \right].$$

$$\mathcal{D}_{\bar{\Phi}}(x_i, y_{i+1}) = \bar{\Phi}(x_i) - \bar{\Phi}(y_{i+1}) - \nabla \bar{\Phi}(y_{i+1})^\top (x_i - y_{i+1})$$

$$= \bar{\Phi}(x_i) - \bar{\Phi}(y_{i+1}) + \nabla \bar{\Phi}(x_i)^\top (y_{i+1} - x_i) + (\nabla \bar{\Phi}(y_{i+1}) - \nabla \bar{\Phi}(x_i))^\top (y_{i+1} - x_i).$$

$$\leq -\frac{\rho}{2} \|y_{i+1} - x_i\|_2^2 - \eta \cdot g_v^\top (y_{i+1} - x_i)$$

$$\leq -\frac{\rho}{2} \|y_{i+1} - x_i\|_2^2 - \eta \cdot \|g_v\|_* \|y_{i+1} - x_i\|$$

$$\leq \frac{\eta^2 \cdot \|g_v\|_*^2}{2\rho}$$

III.

Corollary 3.6

Consider the offline setting where $f_i = f$, and f is L -Lipschitz w.r.t. $\|\cdot\|$.

Let $R^2 = \sup_{x \in \mathbb{R}^n} \mathcal{D}_{\bar{\Phi}}(x, x_i)$. Then

$$f\left(\frac{1}{t} \sum_{i=1}^t x_i\right) - f(x^*) \leq \frac{\mathcal{D}_{\bar{\Phi}}(x^*, x_1)}{\eta t} + \frac{\eta \cdot L^2}{2\rho} \stackrel{(2)}{\leq} RL \cdot \sqrt{\frac{2}{\rho t}}.$$

Proof. (1) Theorem A.2. + Theorem 3.5

(2) Minimize over η (choose $\eta = \frac{R}{L} \cdot \sqrt{\frac{2\rho}{t}}$)

IV

Corollary 3.7 [Quadratic Form Regret Bound]

The regret of OMD satisfies

$$R_n \leq \frac{\text{diam}(\mathcal{D}, \mathcal{A})}{\gamma} + \frac{\gamma}{2} \cdot \sum_{t=1}^n \|g_t\|_{\mathcal{F}}^{-2} \cdot \mathbb{E}(z_t),$$

Where $z_t = \alpha_t x_t + (1-\alpha_t) \tilde{x}_{t+1}$, $\alpha_t \in [0, 1]$ ($\tilde{x}_{t+1} := y_{t+1}$ in the above).

proof. The analysis of Theorem 3.5 shows that

$$R_n \leq \frac{\text{diam}(F, \mathcal{A})}{\gamma} + \frac{1}{\gamma} \cdot \sum_{t=1}^n D_F(x_t, \tilde{x}_{t+1}).$$

For $D_F(x_t, \tilde{x}_{t+1})$, we have

$$\begin{aligned} D_F(x_t, \tilde{x}_{t+1}) &= -D_F(\tilde{x}_{t+1}, x_t) + \langle \nabla F(x_t) - \nabla F(\tilde{x}_{t+1}), x_t - \tilde{x}_{t+1} \rangle \\ &= -D_F(\tilde{x}_{t+1}, x_t) + \langle -\gamma g_t, x_t - \tilde{x}_{t+1} \rangle \\ &\leq -D_F(\tilde{x}_{t+1}, x_t) + \frac{1}{2} \|\gamma g_t\|_{\mathcal{F}}^{-2} \cdot \mathbb{E}(z_t) + \frac{1}{2} \|x_t - \tilde{x}_{t+1}\|_{\mathcal{F}}^{-2} \cdot \mathbb{E}(z_t) \\ &= \frac{\gamma}{2} \|g_t\|_{\mathcal{F}}^{-2} \cdot \mathbb{E}(z_t), \end{aligned}$$

where the inequality follows from Hölder ineq. + AM-GM ineq. / Young-Fenchel ineq. + Lemma A.5, and the last equality from choosing z_t satisfies Lemma A.4.

□

Remark 3.8 [Advantage of choosing $z_t = \alpha_t x_t + (1-\alpha_t) \tilde{x}_{t+1}$ instead of $z_t = \alpha_t x_t + (1-\alpha_t) x_{t+1}$]

Example 4.1 [Online convex optimization over the simplex]

Suppose that $f_1, f_2, \dots, f_n : \Delta^n \rightarrow \mathbb{R}$ are convex functions, each of which is L -Lipschitz w.r.t. $\|\cdot\|_1$. Choose Φ as the negative entropy function.

$x_i = (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n})$ and $\gamma = \sqrt{\frac{2 \ln n}{t}}$. Then

$$\sum_{i=1}^t (f_i(x_i) - f_i(x^*)) \leq \sqrt{2t \ln n}.$$

Proof. Thm 3.5 (2) shows that

$$\begin{aligned} \sum_{i=1}^t (f_i(x_i) - f_i(x^*)) &\leq \frac{D_{KL}(x^*, x_i)}{\gamma} + \frac{\gamma}{2} \cdot \sum_{i=1}^t \|g_i\|_\infty^2 \\ &\leq \frac{\ln n}{\gamma} + \frac{\gamma}{2} t \\ &= \sqrt{2t \ln n}, \end{aligned}$$

where the second inequality is due to Lemma A.3, and Theorem A.2.

□

Proposition 4.1.

Suppose $\mathcal{A} = \Delta^k$. Denote $p^t = \operatorname{argmin}_{p \in \Delta^k} D_\psi(p, p^{t-1}) + \gamma_t \langle \hat{l}^{t-1}, p \rangle$.

$$\nabla \psi(p^{t+1}) = \nabla \psi(p^t) - \gamma_{t+1} \cdot \hat{l}^t - \lambda \cdot \mathbf{1},$$

where $\mathbf{1}$ is the one-vector. (The same also holds for FTRL)

Proof. $p^{t+1} = \operatorname{argmin}_{p \in \Delta^k} D_\psi(p, p^t) + \gamma_{t+1} \langle \hat{l}^t, p \rangle$ satisfying
 \rightarrow Lagrangian Multiplier.

$$\nabla_p (D_\psi(p, p^t) + \gamma_{t+1} \langle \hat{l}^t, p \rangle + \lambda \cdot (\langle p, \mathbf{1} \rangle - 1)) \Big|_{p=p^{t+1}} = 0.$$

□

Appendix A. Technical Stuff

Proposition A.1 Let $f: \mathbb{X} \rightarrow \mathbb{R}$ be strictly convex function, and $x, x' \in \mathbb{X}$ be distinct. Then $\partial f(x) \cap \partial f'(x) = \emptyset$.

Theorem A.2.

Let \mathbb{X} be convex and open, and $f: \mathbb{X} \rightarrow \mathbb{R}$ be convex. If f is L -Lipschitz w.r.t. $\|\cdot\|$, then

$$\|g\|_* \leq L, \forall x \in \mathbb{X}, \forall g \in \partial f(x).$$

Lemma A.3

Suppose $p \in \mathbb{R}_{>0}^n$ is a distribution. Let $q = (\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n}) \in \mathbb{R}_{>0}^n$ be the uniform distribution. Then $D_{KL}(p \parallel q) \leq \ln n$.

Lemma A.4 [Mean Value of Bregman Divergence]

If Ψ is twice differentiable, then

$$D_\Psi(x, y) = \frac{1}{2} \cdot \|y - x\|^2 \nabla^2 \Psi(z)$$

for some $z = \alpha x + (1-\alpha)y$, $\alpha \in [0, 1]$.

Proof. Taylor Theorem. (TBC).

Lemma A.5 [Convex conjugate of squared norms]

Let $f(x) = \frac{1}{2} \|x\|^2$ w.r.t. some norm $\|\cdot\|$. Then the convex conjugate of $f(x)$ is $f^*(y) = \frac{1}{2} \|y\|_*\|^2$.

Proof. (1) We show that $f^*(y) \leq \frac{1}{2} \|y\|_*^2$:

$$\text{LHS} = \sup_x x^T \cdot y - \frac{1}{2} \|x\|^2 \leq \sup_x \|x\| \cdot \|y\|_* - \frac{1}{2} \|x\|^2 \leq \frac{1}{2} \|y\|_*^2.$$

x s.t. $\|x\| = \|y\|_*$

(2) We show that $f^*(y) \geq \frac{1}{2} \|y\|_*^2$:

LHS = $\sup_x x^T \cdot y - \frac{1}{2} \|x\|^2$. Choosing x s.t. (a) $x^T \cdot y = \|x\| \cdot \|y\|_*$; and (b) $\|x\| = \|y\|_*$ proves the desired result.

■

Appendix. B. Tiny Stuff

Proposition B.1 [Negative Entropy and Unnormalized Negative Entropy].

If $F_1 = NE$. $F_2 = UNE$, then $D_{F_1}(.,.) = D_{F_2}(.,.)$.

Proof.

$$D_{F_1}(x,y) = F_1(x) - \underline{F_1(y)} - \langle \nabla F_1(y), x-y \rangle$$

$$D_{F_2}(x,y) = \underline{F_2(x)} - F_2(y) - \langle \nabla F_2(y), x-y \rangle$$

① If $x, y \in \Delta_{k-1}$.

$$\nabla F_1(y)_i = \ln y_i + 1$$

$$\nabla F_2(y)_i = \ln y_i$$

$$\begin{aligned} \text{but } \langle \nabla F_1(y), x-y \rangle &= \sum_i (\ln y_i + 1) \cdot (x_i - y_i) \\ &= \sum_i (\ln y_i) \cdot (x_i - y_i) + \sum_i (x_i - y_i) \\ &= \langle \nabla F_2(y), x-y \rangle \end{aligned}$$

② $x, y \in \Delta_{k-1}$ 不成立.

$$\begin{aligned} D_{F_2}(x,y) &= \sum_i (x_i \cdot \ln x_i - \underline{x_i}) - \sum_i (y_i \cdot \ln y_i - \underline{y_i}) - \sum_i (x_i - y_i) \\ &\quad - \sum_i (\ln y_i) \cdot (x_i - y_i) \end{aligned}$$

$$\begin{aligned} D_{F_1}(x,y) &= \sum_i (x_i \cdot \ln x_i) - \sum_i y_i \cdot \ln y_i \\ &\quad - \sum_i (\ln x_i + 1) \cdot (x_i - y_i) - \sum_i (x_i - y_i) \end{aligned}$$

Some again!

