

Chapter 7 Smoothing Techniques

Convex Conjugate

- Definition 7.1 (Convex conjugate) For any function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, its *convex conjugate* is $f^*(y) = \sup_{x \in \text{dom}(f)} \{x^\top y - f(x)\}$
- Fenchel's inequality $f^*(y) \geq x^\top y - f(x), \forall x, y \Rightarrow x^\top y \leq f(x) + f^*(y), \forall x, y$.
 - which is a generalization of Young's inequality $x^\top y \leq \frac{\|x\|^2}{2} + \frac{\|y\|_*^2}{2}, \forall x, y$
- Lemma 7.2 (Sec.12, [Roc97], no proof here) If f is convex, lower semi-continuous and proper, then $(f^*)^* = f$.
 - Lower semi-continuous (l.s.c) means $f(\mathbf{x}) \leq \liminf_{t \rightarrow \infty} f(\mathbf{x}_t)$ for $\mathbf{x}_t \rightarrow \mathbf{x}$. Proper means $f(\mathbf{x}) \geq -\infty$.
- Proposition 7.3 If f is μ -strongly convex then f^* is continuously differentiable and μ^{-1} -Lipschitz smooth.
 - Proof
 - Differentiable
 - Since f strongly-convex \Rightarrow " f + linear term" strictly convex $\Rightarrow \sup_{x \in \text{dom}(f)} \{x^\top y - f(x)\}$ have a unique solution for all $y \in \mathbb{R}^n$, $\text{dom}(f^*)$ is then convex.
 - f^* is convex since

$$\sup_{x \in \text{dom}(f)} \{x^\top (\lambda y_1 + (1 - \lambda)y_2) - f(x)\} \geq \lambda \sup_{x \in \text{dom}(f)} \{x^\top y_1 - f(x)\} + (1 - \lambda) \sup_{x \in \text{dom}(f)} \{x^\top y_2 - f(x)\}$$
 - then ∂f^* is non-empty.
 - Since f strongly convex, its level set $\{x : f(x) \leq \alpha\}$ is closed, and its proper $f(x) > -\infty$. By lemma 7.2, $(f^*)^* = f$.
 - Consider f^* 's subdifferential, $g \in \partial f^*(z) \Leftrightarrow \forall y, f^*(y) \geq f^*(z) + g^\top(y - z) \Leftrightarrow \forall y, g^\top y - f^*(y) \leq g^\top z - f^*(z)$
 - this means $g^\top z - f^*(z) = (f^*)^*(g) = f(g)$, equivalently,
 - $f^*(z) = g^\top z - f(g) \Leftrightarrow g \in \arg \max_{x \in \text{dom} f} \{y^\top x - f(x)\}$.
 - Since f strongly convex, $\arg \max$ is unique, this is a single-valued mapping.
 - By Lemma 26.1 in [Roc97], ∂f is single values mapping iff f essentially smooth. In that case ∂f reduce to gradient mapping ∇f .
 - This means f^* is differentiable.
 - μ^{-1} -smoothness
 - By similar argument, we know $y \in \partial f(g_y)$, by strong convexity $f(g_z) \geq f(g_y) + y^\top(g_z - g_y) + \mu\|g_y - g_z\|_a^2/2$
 - and also $f(g_y) \geq f(g_z) + z^\top(g_y - g_z) + \mu\|g_y - g_z\|_a^2/2$
 - combine two, we have $\mu\|g_y - g_z\|_a^2 \leq (g_z - g_y)^\top(z - y) \leq \|g_z - g_y\|_a\|z - y\|_a \Rightarrow \|g_y - g_z\|_a \leq \mu^{-1}\|z - y\|_a$
 - Lemma 7.4 (Exercise 49) Let f and g be two proper, convex and semi-continuous functions, then
 - (a) $(f + g)^*(x) = \inf_y \{f^*(y) + g^*(x - y)\}$
 - (b) $(\alpha f)^*(x) = \alpha f^*\left(\frac{x}{\alpha}\right)$ for $\alpha > 0$.
 - Proof
 - (a)
 - $(f + g)^*(x) = \sup_{y \in \text{dom}(f \cap g)} \{x^\top y - f(y) - g(y)\} = \sup_{y \in \text{dom}(f \cap g)} \inf_z \{(x - z)^\top y - g(y) + f^*(z)\}$
 - Since function $h(y, z) := (x - z)^\top y - g(y) + f^*(z)$ is convex w.r.t z and concave w.r.t y , then by von Neumann-Fan minimax theorem $\min \max = \max \min$
 - then $(f + g)^*(x) = \inf_z \sup_{y \in \text{dom}(f \cap g)} \{(x - z)^\top y - g(y) + f^*(z)\} = \inf_z \{g^*(x - z) + f^*(z)\}$
 - (b) $(\alpha f)^*(x) = \sup_{y \in \text{dom}(f \cap g)} \{x^\top y - \alpha f(y)\} = \alpha \sup_{y \in \text{dom}(f \cap g)} \{\alpha^{-1}x^\top y - f(y)\} = \alpha f^*(x/\alpha)$.
 - Examples
 - $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top Q\mathbf{x}, Q \succ 0, f^*(y) = \frac{1}{2}\mathbf{y}^\top Q^{-1}\mathbf{y}$.
 - $f(\mathbf{x}) = \sum_{i=1}^n x_i \log(x_i), f^*(\mathbf{y}) = \sum_{i=1}^n e^{y_i - 1}$.
 - $f(\mathbf{x}) = -\sum_{i=1}^n \log(x_i), f^*(\mathbf{y}) = -\sum_{i=1}^n \log(-y_i) - n$.
 - $f(\mathbf{x}) = \|\mathbf{x}\|, f^*(\mathbf{y}) = \begin{cases} 0, & \|\mathbf{y}\|_* \leq 1 \\ +\infty, & \|\mathbf{y}\|_* > 1. \end{cases}$

Smoothing Techniques

Common Smoothing techs

- Nesterov's smoothing based on conjugate function $f_\mu(x) = \max_{y \in \text{dom}(f^*)} \{x^\top y - f^*(y) - \mu \cdot d(y)\} = (f^* + \mu d)^*(x)$, where $d(y)$ is some proximity function, d is (i) 1-strongly convex (ii) non-negative everywhere and (possibly iii) $\min d(y) = 0$.
 - By Prop 7.3, f_μ is continuously differentiable and μ^{-1} -Lipschitz.
 - Examples of d :

- $d(\mathbf{y}) = \frac{1}{2} \|\mathbf{y} - \mathbf{y}_0\|_2^2$
- $d(\mathbf{y}) = \frac{1}{2} \sum w_i (y_i - y_{0,i})^2$ with $w_i \geq 1$
- $d(\mathbf{y}) = V_\omega(\mathbf{y}, \mathbf{y}_0)$.

- Moreau-Yosida smoothing/regularization $f_\mu(x) = \min_{y \in \text{dom}(f)} \left\{ f(y) + \frac{1}{2\mu} \|x - y\|_2^2 \right\}$.
 - when $d(\mathbf{y}) = \frac{1}{2} \|\mathbf{y}\|_2^2$, Nesterov's smoothing turns to Moreau-Yosida regularization (by Lemma 7.4).
- Lasry-Lions regularization double application of M-Y reg $f_{\mu,\delta}(x) = \max_y \min_z \left\{ f(z) + \frac{1}{2\mu} \|z - y\|_2^2 - \frac{1}{2\delta} \|y - x\|_2^2 \right\}$
 - If f 1-Lipschitz, then choosing $(\delta, \mu) = O(\epsilon)$ guarantees ϵ approximation error, and $f_{\mu,\delta}(\mathbf{x})$ is $O(1/\epsilon)$ -smooth.
 - Computation inefficiency due to solving nonconvex problems.
- Randomized smoothing $f_\mu(x) = \mathbb{E}_Z f(x + \mu Z)$ where Z is isotropic Gaussian or uniform.
 - Choosing $\mu = O(\epsilon)$ guarantees ϵ approximation error.
 - $f_\mu(\mathbf{x})$ is $O\left(\frac{\sqrt{d}}{\epsilon}\right)$ -smooth, dimension dependent.
 - Stochastic gradient is computationally efficient through sampling.
- Ben-Tal-Teboulle smoothing based on recession function. Only applicable to particular class $f(x) = F(f_1(x), f_2(x), \dots, f_m(x))$, where $F(y) = \max_{x \in \text{dom}(g)} \{g(x + y) - g(x)\}$ is the recession function of some func g .
 - The smoothed function is $f_\mu(x) = \mu g\left(\frac{f_1(x)}{\mu}, \dots, \frac{f_m(x)}{\mu}\right)$.

Nesterov's Smoothing

- Proposition 7.10 $\forall \mu > 0$, let $D_Y^2 = \max_{y \in Y} d(y)$, $f(x) - \mu D_Y^2 \leq f_\mu(x) \leq f(x)$.
 - Proof
 - Left: $f_\mu(x) = \max_{y \in \text{dom}(f^*)} \{x^\top y - f^*(y) - \mu \cdot d(y)\} \geq \max \{x^\top y - f^*(y)\} - \mu \max \{d(y)\} = f(x) - \mu D^2$
 - Right: $f_\mu(x) = \max_{y \in \text{dom}(f^*)} \{x^\top y - f^*(y) - \mu \cdot d(y)\} \leq f_\mu(x) = \max_{y \in \text{dom}(f^*)} \{x^\top y - f^*(y)\} = f(x)$
- Remark Tradeoff b.t.w. approximation error and optimization efficiency $f(\mathbf{x}) - f^* \leq \underbrace{f(\mathbf{x}) - f_\mu(\mathbf{x})}_{\text{approximation error}} + \underbrace{f_\mu(\mathbf{x}) - \min_x f_\mu(\mathbf{x})}_{\text{optimization error}}$.
- Suppose we can access $\nabla f_\mu(\mathbf{x})$ and apply gradient methods to solve $\min_{x \in X} f_\mu(x)$
 - PGD gives $f(x_t) - f_* \leq O\left(\frac{R^2}{\mu t} + \mu D_Y^2\right)$
 - to achieve error of ϵ , we have to set $\mu = O\left(\frac{\epsilon}{D_Y^2}\right)$, and the number of iteration is $T = O\left(\frac{R^2 D_Y^2}{\epsilon^2}\right)$
 - Accelerate GD gives $f(x_t) - f_* \leq O\left(\frac{R^2}{\mu t^2} + \mu D_Y^2\right)$
 - similarly to achieve error of ϵ , we have to set $\mu = O\left(\frac{\epsilon}{D_Y^2}\right)$ and $T = O\left(\frac{RD_Y}{\epsilon}\right)$

Example

- The function form is generalized as $f(x) = \max_{y \in Y} \{\langle Ax + b, y \rangle - \phi(y)\}$, where ϕ convex and continuous, Y convex and compact.
 - example like $f(x) = \max_{1 \leq i \leq m} |a_i^\top x - b_i|$ can be re-written as $f(x) = \max_{y \in Y} \sum_{i=1}^m (a_i^\top x - b_i)y_i$, where $Y := \{y \in \mathbb{R}^m : \sum_{i=1}^m |y_i| \leq 1\}$.
 - The smoothed function is $f_\mu(x) = \max_{y \in Y} \{\langle Ax + b, y \rangle - \phi(y) - \mu \cdot d(y)\} = (\phi + \mu \cdot d)^*(Ax + b)$
- As an example, $f(x) = |x| = \sup_{|y| \leq 1} yx = \sup_{\substack{y_1, y_2 \geq 0 \\ y_1 + y_2 = 1}} (y_1 - y_2)x$,
 - $d(y) = y^2/2$, then $f_\mu(x) = \sup_{|y| \leq 1} \{yx - \frac{\mu}{2} y^2\} = \begin{cases} \frac{x^2}{2\mu}, & |x| \leq \mu \\ |x| - \frac{\mu}{2}, & |x| > \mu \end{cases}$
 - $d(y) = 1 - \sqrt{1 - y^2}$ (half circle), clearly 1-strongly convex. $f_\mu(x) = \sup_{|y| \leq 1} \{yx - \mu(1 - \sqrt{1 - y^2})\} = \sqrt{x^2 + \mu^2} - \mu$.
 - This is a special case of Ben-Tal & Teboulle's smoothing $|x| = \sup_y \{g(x + \mu) - g(y)\}$, $g(y) = \sqrt{1 + y^2}$ and $f_\mu(x) = \mu g(x/\mu) = \sqrt{x^2 + \mu^2}$
 - $d(y) = y_1 \log y_1 + y_2 \log y_2 + \log 2$, where $Y = \{(y_1, y_2) : y_1, y_2 \geq 0, y_1 + y_2 = 1\}$, $f_\mu(x) = \mu \log\left(\frac{e^{-\frac{x}{\mu}} + e^{\frac{x}{\mu}}}{2}\right)$
 - This can be seen as $|x| = \max\{x, -x\} = \sup_y \{g(x + \mu) - g(y)\}$, $g(y) = \log(e^{y_1} + e^{y_2})$.

Moreau-Yosida Regularization

- Definition 7.11 the *proximal operator* of convex f at a given point x is defined as $\text{prox}_f(x) = \operatorname{argmin}_y \{f(y) + \frac{1}{2} \|x - y\|^2\}$.
- Examples
 - Indicator function $f(\mathbf{x}) = \delta_X(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in X \\ +\infty, & \mathbf{x} \notin X \end{cases}$ $\text{prox}_f(\mathbf{x}) = \Pi_X(\mathbf{x})$ is the projection.

- $f(\mathbf{x}) = \mu \|\mathbf{x}\|_1$, $\text{prox}_{\mu|\cdot|}(x_i) = \begin{cases} x_i - \mu & \text{if } x_i > \mu \\ 0 & \text{if } |x_i| \leq \mu \\ x_i + \mu & \text{if } x_i < -\mu \end{cases} = \text{sign}(\mathbf{x}) \otimes [|\mathbf{x}| - \lambda] + \text{soft thresholding operator.}$
- 2-norm $f(\mathbf{x}) := \|\mathbf{x}\|_2$, $\text{prox}_{\lambda f}(\mathbf{x}) = [1 - \lambda/\|\mathbf{x}\|_2] + \mathbf{x}$
- Support function $f(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{C}} \mathbf{x}^T \mathbf{y}$, $\text{prox}_{\lambda f}(\mathbf{x}) = \mathbf{x} - \lambda \pi_{\mathcal{C}}(\mathbf{x})$
- Square norm $f(\mathbf{x}) := (1/2)\|\mathbf{x}\|_2^2$, $\text{prox}_{\lambda f}(\mathbf{x}) = (1/(1+\lambda))\mathbf{x}$
- log function $f(\mathbf{x}) := -\log(x)$, $\text{prox}_{\lambda f}(x) = \left((x^2 + 4\lambda)^{1/2} + x \right) / 2$
- $\log \det f(\mathbf{x}) := -\log \det(\mathbf{X})$, $\text{prox}_{\lambda f}(x)$ is the log function prox applied to the individual eigenvalues of \mathbf{X} .
- Proposition 7.13 f convex, then
 - (a) fixed point \mathbf{x}^* minimized $f(\mathbf{x})$ iff $\mathbf{x}^* = \text{prox}_f(\mathbf{x}^*)$
 - Proof (one direction is obvious)
 - $\mathbf{x}^* = \text{prox}_f(\mathbf{x}^*) \Leftrightarrow \forall \mathbf{y}, f(\mathbf{x}^*) \leq f(\mathbf{y}) + \frac{1}{2} \|\mathbf{x}^* - \mathbf{y}\|^2$
 - If $\exists \mathbf{y}_0$ s.t. $f(\mathbf{y}) < f(\mathbf{x}^*)$, by convexity, $f(\mathbf{y}_\lambda) := f(\mathbf{x}^* + \lambda(\mathbf{y} - \mathbf{x}^*)) \leq f(\mathbf{x}^*) - \lambda(f(\mathbf{x}^*) - f(\mathbf{y})) \Rightarrow f(\mathbf{x}^*) \geq f(\mathbf{x}^* + \lambda(\mathbf{y} - \mathbf{x}^*)) + \lambda(f(\mathbf{x}^*) - f(\mathbf{y}))$
 - Setting $\lambda \min\{1, \alpha 2(f(\mathbf{x}^*) - f(\mathbf{y}))/\|\mathbf{x}^* - \mathbf{y}\|^2\}$, we get $f(\mathbf{x}^*) \geq f(\mathbf{y}_\lambda) + \frac{\alpha}{2} \|\mathbf{x}^* - \mathbf{y}_\lambda\|^2$, setting $\alpha > 1$ and we get contradiction.
 - Another path is to prove $0 \in \partial f(\mathbf{x}_*) + (\mathbf{x}_* - \mathbf{x}_*) \Rightarrow 0 \in \partial f(\mathbf{x}_*)$.
 - (b) Non-expansive $\|\text{prox}_f(\mathbf{x}) - \text{prox}_f(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|$
 - Proof
 - by optimial condition (suppose $\text{dom}(f) = \mathbb{R}^n$), $\nabla f(\text{prox}(\mathbf{x})) + \text{prox}(\mathbf{x}) - \mathbf{x} = 0$ and $\nabla f(\text{prox}(\mathbf{y})) + \text{prox}(\mathbf{y}) - \mathbf{y} = 0$
 - By convexity $(\nabla f(\text{prox}(\mathbf{y})) - \nabla f(\text{prox}(\mathbf{x})))^\top (\text{prox}(\mathbf{y}) - \text{prox}(\mathbf{x})) \geq 0$
 - This means $(\mathbf{y} - \mathbf{x})^\top (\text{prox}(\mathbf{y}) - \text{prox}(\mathbf{x})) \geq \|\text{prox}(\mathbf{y}) - \text{prox}(\mathbf{x})\|_2^2 \Rightarrow \text{QED}$
 - (c) Moreau Decomposition $\mathbf{x} = \text{prox}_f(\mathbf{x}) + \text{prox}_{f^*}(\mathbf{x})$
 - Proof
 - By fact of $\partial f^*(g) = \arg \max_x \{g^T x - f(x)\} = \{x | \nabla f(x) = g\}$
 - we already have $\nabla f(\text{prox}(\mathbf{x})) + \text{prox}(\mathbf{x}) - \mathbf{x} = 0$
 - which can be seen as $g = \nabla f(\text{prox}(\mathbf{x}))$, so $\nabla f^*(g) = \text{prox}(\mathbf{x})$
 - then this means $g + \nabla f^*(g) - \mathbf{x} = 0$, which is the optimial condition for $\text{prox}_{f^*}(\mathbf{x})$, QED.

Proximal Point Algorithm

- Proximal Point Algorithm (PPA) $\mathbf{x}_{t+1} = \text{prox}_{\gamma_t f}(\mathbf{x}_t)$
 - This comes from $\mathbf{x}_{t+1} = \mathbf{x}_t - L^{-1} \nabla f_\mu(\mathbf{x}_t)$, (Danskin's thm used here)
- Example $\min_{\mathbf{x}} \|A\mathbf{x} + \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_1 = \min_{\mathbf{x}} f_1 + f_2$, where f_1 differentiable while f_2 is not. Usually we approximate $\mathbf{x}_{t+1} = \text{prox}_{f_2}(\mathbf{x}_t - \frac{1}{L_{f_1}} \nabla f_1(\mathbf{x}_t))$.
- Theorem 7.14 If f convex, then PPA satisfies $f(\mathbf{x}_t) - f^* \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2 \sum_{\tau=0}^{t-1} \gamma_\tau}$
 - Proof
 - by the optimality condition of \mathbf{x}_{t+1} , $f(\mathbf{x}_{t+1}) + \frac{1}{2\gamma_t} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|_2^2 \leq f(\mathbf{x}_t)$, this is something like sufficient decrease.
 - by first order optimality condition, $0 \in \partial f(\mathbf{x}_{t+1}) + \frac{1}{\gamma_t} (\mathbf{x}_{t+1} - \mathbf{x}_t) \Rightarrow \frac{\mathbf{x}_t - \mathbf{x}_{t+1}}{\gamma_t} \in \partial f(\mathbf{x}_{t+1})$,
 - so by convexity/subgradient
 - $f(\mathbf{x}_{t+1}) - f^* \leq \frac{1}{\gamma_t} (\mathbf{x}_t - \mathbf{x}_{t+1})^T (\mathbf{x}_{t+1} - \mathbf{x}_*) = \frac{1}{\gamma_t} \left[(\mathbf{x}_t - \mathbf{x}_*)^T (\mathbf{x}_{t+1} - \mathbf{x}_*) - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2 \right]$
 - By Young's inequality, $(\mathbf{x}_t - \mathbf{x}_*)^T (\mathbf{x}_{t+1} - \mathbf{x}_*) \leq \frac{1}{2} [\|\mathbf{x}_t - \mathbf{x}_*\|^2 + \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2]$, plug this in and we get
 - $\gamma_t (f(\mathbf{x}_{t+1}) - f^*) \leq \frac{1}{2} [\|\mathbf{x}_t - \mathbf{x}_*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}_*\|^2]$
 - Then we can happily do our summation over $t \in [0 : T - 1]$
 - Since we have monotonicity in f , this algorithm have last iteration convergence.
 - Comment prox step is more like a second order step.... more computation, but also more effective.
 - Algorithm converge a.l.a $\sum_t \gamma_t \rightarrow \infty$.
 - Larger γ_t makes algo converge faster, but harder to solve each step.