

CS 726 Cheat Sheet Midterm

I. Setting the Stage

1. $\langle x, y \rangle = \sum_{i=1}^d x_i y_i$; $\|x\|_2 = \sqrt{\langle x, x \rangle} = \sqrt{\sum_{i=1}^d x_i^2}$; $\|x\|_p = (\sum_{i=1}^d |x_i|^p)^{1/p}$
 $\|x\|_1 = \sum_{i=1}^d |x_i|$; $\|x\|_\infty = \max_i |x_i|$; $\|z\|_1 = \sup_{x \in \mathbb{R}^d, \|x\|_2=1} \langle z, x \rangle$
 Prop 1.1 $\langle z, x \rangle \leq \|z\|_1 \cdot \|x\|_\infty$

dual norm $\| \cdot \|_q$ $\frac{1}{p} + \frac{1}{q} = 1$, $\| \cdot \|_p$ dual $\| \cdot \|_q$ $p \geq 1$
 $\forall p \geq 1 \forall r > p: \|x\|_r \leq \|x\|_p \leq d^{1/p-1/r} \|x\|_r$

2. Def 2.1 \mathcal{X} is convex if $x, y \in \mathcal{X}$ and $\lambda \in (0, 1)$:

(Ob, Ep) $(1-\lambda)x + \lambda y = x + \lambda(y-x) \in \mathcal{X}$

Def 3.1 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is lower semicontinuous at $\bar{x} \in \mathbb{R}^d$ if $\liminf_{x \rightarrow \bar{x}} f(x) = f(\bar{x})$ and lower semicontinuous on \mathbb{R}^d if it is lower semicontinuous at every $\bar{x} \in \mathbb{R}^d$

Def 3.2 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is Lipschitz-continuous on a set $\mathcal{X} \subseteq \mathbb{R}^d$ if there is $M < \infty$ s.t. $\forall x, y \in \mathcal{X}, |f(x) - f(y)| \leq M \|x - y\|$

Def 3.3 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is L-smooth wrt $\| \cdot \|$ on \mathcal{X} if $\exists L < \infty$ s.t. $\forall x, y \in \mathcal{X}: \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$

Def 3.4 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is (K, L)-weakly smooth for $K \in [1, 2]$ if $\exists L < \infty$ $\forall x, y \in \mathcal{X}: \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|^{K-1}$

Def 3.6 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is convex if $\forall x, y \in \mathbb{R}^d \forall \lambda \in (0, 1)$, $f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y)$

Lemma 3.7 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is convex iff its epigraph, i.e., $\text{epi}(f) = \{(x, a) : x \in \mathbb{R}^d, a \in \mathbb{R}, f(x) \leq a\}$ is convex

Lemma 3.8 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ differentiable, then f convex iff $\forall x, y \in \mathbb{R}^d$, $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$

3.1. $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ continuously differentiable ($\nabla f > 0$) ($\exists S > 0$) ($\forall x, y \in \mathbb{R}^d$)

$\|x - y\| \leq \delta \Rightarrow \|\nabla f(x) - \nabla f(y)\| \leq \epsilon$ (twice)
 operator norm for symmetric $A: \|A\|_2 = \sup_{x \in \mathbb{R}^d, \|x\|_2=1} \|Ax\|_2$

Thm 3.9 (Taylor) Let $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ continuously differentiable. $\forall x, y \in \mathbb{R}^d$:

- (i) $f(y) = f(x) + \int_0^1 \langle \nabla f(x + t(y-x)), y-x \rangle dt$
- (ii) $\exists t \in (0, 1)$ s.t. $f(y) = f(x) + \langle \nabla f(x + t(y-x)), y-x \rangle$ (MVT)
- If f twice continuously differentiable,
- (iii) $\nabla f(y) = \nabla f(x) + \int_0^1 \nabla^2 f(x + t(y-x))(y-x) dt$
- (iv) $\exists t \in (0, 1)$ s.t. $f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \frac{1}{2} \langle \nabla^2 f(x + t(y-x))(y-x), y-x \rangle$

Lemma 3.10 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ L-smooth. $\forall x, y \in \mathbb{R}^d$, both hold:
 $f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2$; $f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2$

Lemma 3.11 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ twice continuously differentiable function. If $\forall x \in \mathbb{R}^d$ $\|\nabla^2 f(x)\|_p \leq \sup_{y: \|y\|_q=1} \|\nabla^2 f(x)y\|_p \leq L < \infty$, $\frac{1}{p} + \frac{1}{q} = 1$, $p \geq 1$, then f is L-smooth wrt $\| \cdot \|_p$

Lemma 3.12 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ twice continuously differentiable, then f is L-smooth wrt $\| \cdot \|_2$ iff $-\infty < \lambda \leq \nabla^2 f(x) \leq \infty$

Def 3.13 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ strongly convex with modulus $\mu > 0$ if $\forall x, y \in \mathbb{R}^d$, $\lambda \in (0, 1)$: $f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y) - \lambda(1-\lambda) \frac{\mu}{2} \|y-x\|^2$

Lemma 3.14 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ continuously differentiable, then f is μ -strongly convex iff $\forall x, y \in \mathbb{R}^d: f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2$

Cor f twice continuously differentiable, then μ -strongly convex wrt $\| \cdot \|_2$ iff $\nabla^2 f(x) \succeq \mu I$

Def 3.15 $p \geq 2$, $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ p -uniformly convex with modulus $M > 0$ wrt $\| \cdot \|_1$ if $\forall x, y \in \mathbb{R}^d, \lambda \in (0, 1)$: (p, M)-uniformly convex

$f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y) - \lambda(1-\lambda) \frac{M}{p} \|y-x\|^p$

Cor f is (p, M) -uniformly convex, then $\forall x, y \in \mathbb{R}^d$ $f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{M}{p} \|y-x\|^p$

Def 3.17 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ satisfy (r, μ)-sharp or tojasiewicz inequality on \mathcal{X} if $\mathcal{X}^* = \text{argmin}_{x \in \mathbb{R}^d} f(x)$ is nonempty and $\exists \mu > 0, r > 0$, $\forall x \in \mathcal{X}: f(x) - \min_y f(y) \geq \frac{\mu}{r} \text{dist}(x, \mathcal{X}^*)^r$ where $\text{dist}(x, \mathcal{X}^*) = \min_{x^* \in \mathcal{X}^*} \|x - x^*\|$

Cor μ -strongly convex \Rightarrow (2, μ)-sharp

Def 3.18 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ differentiable satisfy Pt condition on \mathcal{X} if \mathcal{X}^* is nonempty and $\exists \mu, r > 0$ s.t. $\forall x \in \mathcal{X}: f(x) - \min_y f(y) \leq \frac{\mu}{r} \|\nabla f(x)\|^r$

Local and Global Solns

Def 4.1 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ differentiable, \mathcal{X} closed and convex. $x \in \mathcal{X}$ stationary point for (P) if $\forall y \in \mathcal{X}, \langle \nabla f(x), y-x \rangle \geq 0$ iff $\exists x \in \mathcal{X}$ s.t. $\nabla f(x) = 0$

Def 4.2 Given (P) and $x^* \in \mathcal{X}$, x^* local minimum if \exists nbhd $N_{x^*} \forall x \in N_{x^*} \cap \mathcal{X}, f(x) \geq f(x^*)$

x^* global soln if $\forall x \in \mathcal{X}, f(x) \geq f(x^*)$

Thm 4.3 (1st-order Necessity Optimality Conditions) f continuously differentiable and \mathcal{X} closed and convex; if $x^* \in \mathcal{X}$ global soln, x^* must be a stationary point;

If $x^* \in \mathcal{X}$ local soln, $\exists N$ of x^* s.t. $\forall y \in N \cap \mathcal{X}, \langle \nabla f(x^*), y-x^* \rangle \geq 0$

Thm 4.4 (2nd-order Necessary and Sufficient Conditions for Unconstrained) f twice continuously differentiable, $\mathcal{X} \equiv \mathbb{R}^d$:

• If $x^* \in \mathbb{R}^d$ local soln, then $\|\nabla f(x^*)\|_2 = 0$ and $\nabla^2 f(x^*) \succeq 0$

• If $\exists x^* \in \mathbb{R}^d$ s.t. $\|\nabla f(x^*)\|_2 = 0$ and $\nabla^2 f(x^*) \succ 0$, then x^* strict local minimizer

Thm 4.5 (Optimality Conditions for Convex) f convex, \mathcal{X} closed and convex:

- (i) Every local soln is also global
- (ii) Set of solns is convex
- (iii) If f differentiable, then x^* global soln iff x^* stationary point

Thm 4.6 (Optimality Conditions for Convex Probl Strongly Convex Obj)

f μ -strongly convex and continuous on its domain; \mathcal{X} closed, convex, and non-empty intersection w/ effective domain of f , then soln to (P) is attained and unique.

II Basic Descent Methods $x_{k+1} = x_k + \alpha_k z_k$

Def 4.1 $z \in \mathbb{R}^d$ is a descent direction of $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ at x if $\exists t > 0$ s.t. $f(x + t'z) \leq f(x) \forall 0 < t' \leq t$

Prop 4.2 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ is continuously differentiable (in nbhd of $x \in \mathbb{R}^d$), then $z \in \mathbb{R}^d$ s.t. $\langle \nabla f(x), z \rangle < 0$ is a descent direction for f at x

Gradient Descent: $\nabla f(x) \neq 0, x_{k+1} = x_k - \alpha_k \nabla f(x_k)$

Lemma 2.1 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ L-smooth wrt $\| \cdot \|_2$ and $x_0 \in \mathbb{R}^d$, consider GD with step size $\alpha_k = \alpha \in (0, \frac{1}{L}]$. For all $k \geq 0$, $f(x_{k+1}) - f(x_k) \leq -\frac{\alpha}{2} \|\nabla f(x_k)\|_2^2$
 further if $f(x) \geq f^* > -\infty, \forall x \in \mathbb{R}^d, k \geq 0, \min_{k \geq 0} \|\nabla f(x_k)\|_2 \leq \frac{2(f(x_0) - f^*)}{\alpha(k+1)}$
 so $\forall \epsilon > 0$, GD takes at most $k = \lceil \frac{2(f(x_0) - f^*)}{\alpha \epsilon} \rceil$ iterations to construct $x_i, i \in [k, k]$ s.t. $\|\nabla f(x_i)\|_2 \leq \epsilon$

Lemma 2.2 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ convex and L-smooth wrt $\| \cdot \|_2$ and $x_0 \in \mathbb{R}^d$, consider GD with step size $\alpha_k = \alpha \in (0, \frac{1}{L}]$. Then $\forall k \geq 0$, $f(x_{k+1}) - f(x^*) \leq \frac{\|x_0 - x^*\|_2^2}{2\alpha(k+1)}$. So $\forall \epsilon > 0$, GD after at most

$k = \lceil \frac{\|x_0 - x^*\|_2^2}{2\alpha \epsilon} \rceil$ iterations we have $f(x_k) - f(x^*) \leq \epsilon$

Lemma 2.3 $f: \mathbb{R}^d \rightarrow \bar{\mathbb{R}}$ μ -strongly convex and L-smooth wrt $\| \cdot \|_2$ and $x_0 \in \mathbb{R}^d$, GD w/ step size $\alpha_k = \alpha \in (0, \frac{1}{L}]$. $\forall k \geq 0$, $\|x_k - x^*\|_2 \leq (1 - 2\mu\alpha)^k \|x_0 - x^*\|_2$

So $\forall \epsilon > 0$, GD after at most $k = \max\{0, \lceil \frac{1}{2\mu\alpha} \log(\frac{\|x_0 - x^*\|_2}{\epsilon}) \rceil\}$ iterations we have $\|x_k - x^*\|_2 \leq \epsilon$

Cor 2.4 $f: \mathbb{R}^d \rightarrow \mathbb{R}$ μ -strongly convex and L -smooth w.r.t $\|\cdot\|_2$ and $x_0 \in \mathbb{R}^d$, consider GD w/ step size $\alpha_k = \alpha \in (0, 1/L]$. Then,
 • After at most $k = \max\{0, \lceil \frac{2}{\alpha\mu} \log(\frac{\|x_0 - x^*\|_2}{\epsilon}) \rceil\}$ iterations we have $f(x_k) - f(x^*) \leq \epsilon$
 • After at most $k = \max\{0, \lceil \frac{2}{\alpha\mu} \log(\frac{L\|x_0 - x^*\|_2}{\epsilon}) \rceil\}$ iterations we have $\|\nabla f(x_k)\|_2 \leq \epsilon$

Nesterov Acceleration

1. Smooth and Convex Setting

Algo 1 Nesterov-Smooth (x_0, f, L, K)
 Init: $a_0 = A_0 = \frac{1}{L}$, $v_0 = x_0 - a_0 \nabla f(x_0)$, $y_0 = x_0 - \frac{1}{L} \nabla f(x_0)$
 for $k=1$ to K do
 $a_k = \frac{k+1}{2L}$, $A_k = A_{k-1} + a_k$
 $x_k = \frac{A_{k-1}}{A_k} y_{k-1} + \frac{a_k}{A_k} v_{k-1}$
 $v_k = v_{k-1} - a_k \nabla f(x_k)$
 $y_k = x_k - \frac{1}{L} \nabla f(x_k)$
 end for
 return y_k

$f(y_k) - f(x^*) \leq G_k \leq \frac{2L\|x^* - x_0\|_2^2}{(k+1)(k+1)}$ choosing $a_k = \frac{k+1}{2L}$

2. Smooth and Strongly Convex Optimization

2.1 Restart

Algo 2 Nesterov-Restart (x_0, f, L, μ, R)

Init: $\bar{x}_0 = x_0$
 for $r=1$ to R do
 $\bar{x}_r = \text{Nesterov-Smooth}(\bar{x}_{r-1}, f, L, \lceil \frac{R}{\mu} \rceil)$
 end for
 return \bar{x}_R

Motivation: Lemma 3.14 w/ x^*

Lemma 2.1 Let f smooth, convex, $(2, \mu)$ -sharp, $x^* = \text{argmin}_{x \in \mathbb{R}^d} f(x)$
 Give $x_0 \in \mathbb{R}^d$, consider running Algo 2 for $R \geq 1$ iterations. Then
 give $\epsilon > 0$, if $R \geq 2 \log_2(\frac{\text{dist}(x_0, x^*)}{\epsilon})$, we have $\text{dist}(\bar{x}_R, x^*) \leq \epsilon$.

The total number of operations is $O(\frac{L}{\mu} \log_2(\frac{\text{dist}(x_0, x^*)}{\epsilon}))$

2.2 Direct Algorithm

Algo 3 Nesterov-Smooth-Strongly-Convex (x_0, f, L, μ, K)

Init: $a_0 = A_0 = \frac{1}{L-\mu}$, $v_0 = \frac{x_0 + \mu a_0 x_0 - a_0 \nabla f(x_0)}{1 + \mu A_0}$, $y_0 = x_0 - \frac{1}{L} \nabla f(x_0)$
 for $k=1$ to K do
 $a_k = \text{positive soln to } \frac{a_k}{(A_{k-1} + a_k)(1 + \mu(A_{k-1} + a_k))} = \frac{1}{L}$
 $A_k = A_{k-1} + a_k$
 $x_k = \frac{A_{k-1}}{A_k(1-\mu/L)} y_{k-1} + \frac{(1-A_{k-1})}{A_k(1-\mu/L)} v_{k-1}$
 $v_k = \frac{(1+\mu A_{k-1})v_{k-1} + a_k(\mu x_k - \nabla f(x_k))}{1 + \mu A_k}$
 $y_k = x_k - \frac{1}{L} \nabla f(x_k)$
 end for
 return y_k

$f(y_k) - f(x^*) \leq G_k \leq \frac{\|x^* - x_0\|_2^2}{k^2}$; $f(y_k) - f(x^*) < (1 - \frac{\mu}{L})^k \frac{(L-\mu)\|x^* - x_0\|_2^2}{k^2}$
 $k \geq \frac{L}{\mu} \log_2(\frac{(L-\mu)\|x^* - x_0\|_2^2}{\epsilon})$ then $f(y_k) - f(x^*) \leq \epsilon$

IV. Appendix

Inequalities

- triangle: $\|x+y\| \leq \|x\| + \|y\|$ (\mathbb{R}^2, L^p by Minkowski)
- reverse triangle: $\|x-y\| \geq |\|x\| - \|y\||$
- Jensen: $f(\sum_{i=1}^k a_i x_i) \leq \sum_{i=1}^k a_i f(x_i)$ given f convex, $a_i \geq 0$ s.t. $\sum_{i=1}^k a_i = 1$ (Cauchy-Schwarz) $p=q=2$
- Hölder: $\|fg\|_1 \leq \|f\|_p \|g\|_q$ given $p, q \in [1, \infty]$, $\frac{1}{p} + \frac{1}{q} = 1$
- Young: $ab \leq \frac{a^p}{p} + \frac{b^q}{q}$ (equal iff $a^p = b^q$) given $a, b \geq 0, p > 1, q > 1$, $\frac{1}{p} + \frac{1}{q} = 1$

Directional Derivative: $\nabla_v f(x) = \nabla f(x) \cdot v = \lim_{t \rightarrow 0} \frac{f(x+tv) - f(x)}{t}$

Complete Square: $\|u\|_2^2 + \|v\|_2^2 - 2\langle u, v \rangle = \|u-v\|_2^2 \leq n\|u-v\|_2^2$
 $\|u\|_2^2 + \|v\|_2^2 - 2\langle u, v \rangle \leq n\|u-v\|_2^2$

HW Pt holds, $\forall x \in \mathbb{R}^d: \frac{1}{2L} \|\nabla f(x)\|_2^2 \leq f(x) - f(x^*) \leq \frac{L}{2} \|x - x^*\|_2^2$

- 2.6 f L -smooth convex fn, $L \in (0, \infty)$
 - (a) $\frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|_2^2 \leq f(y) - f(x) - \langle \nabla f(x), y-x \rangle$
 - (b) $f(x_{k+1}) \leq f(x_k) - \frac{1}{2L} (\|\nabla f(x_{k+1})\|_2^2 + \|\nabla f(x_k)\|_2^2)$
 - (c) $\|\nabla f(x_{k+1})\|_2 \leq \|\nabla f(x_k)\|_2$
 - (d) $\|\nabla f(x_k)\|_2^2 \leq L(f(x_0) - f(x_k))$
- 2.3 X, Y convex, $X \cap Y$ must convex, ~~$X \cup Y$~~ may not
- 2.5 f Convex $\rightarrow f(a^T x)$ convex; $f(x) = \|x\|$ convex; $f(x) = \frac{1}{2} \|x\|_2^2$ convex
- 2.7 $x^T A x \geq \lambda_{\min} \|x\|_2^2$ $x^T A x \leq \lambda_{\max} \|x\|_2^2$
- 3.1 f differentiable, $x^* \in \text{argmin}_{x \in \mathbb{R}^d} f(x)$
 - (i) if f L -smooth, $\frac{1}{2L} \|\nabla f(x)\|_2^2 \leq f(x) - f(x^*) \leq \frac{L}{2} \|x - x^*\|_2^2$
 - (ii) if f L -smooth and convex, $\frac{1}{2L} \|\nabla f(x) - \nabla f(y)\|_2^2 \leq f(y) - f(x) - \langle \nabla f(x), y-x \rangle$
 - (iii) if f L -smooth and μ -strongly convex, $\langle \nabla f(x) - \nabla f(y), x-y \rangle \geq \frac{\mu}{\mu+L} \|x-y\|_2^2 + \frac{1}{\mu+L} \|\nabla f(x) - \nabla f(y)\|_2^2$
- 3.4 Equivalent Nesterov with $x_k = y_{k-1} + \frac{a_k}{A_k} (\frac{A_{k-1}}{A_{k-1}-1} (y_{k-1} - y_{k-2}))$
 Taylor Thm in Exam Review $f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \int_0^1 \langle \nabla f(x+t(y-x)), y-x \rangle dt$

Additions:

- B_r Taylor, μ -str cvx w.r.t $\|\cdot\|_2 \Leftrightarrow \nabla^2 f(x) \succeq \mu I \forall x$
- smooth + non-smooth ~~\neq~~ smooth cvx + str cvx ~~\neq~~ str cvx
- Bregman Divergence: $D_\psi(x, y) = \psi(x) - \psi(y) - \langle \nabla \psi(y), x-y \rangle$
 (1st order approx error) $\forall x, y \in \mathbb{R}^d$
- 3-pt identity: $D_\psi(x, y) + D_\psi(z, y) + \langle \nabla \psi(z) - \nabla \psi(y), x-z \rangle + D_\psi(x, z)$
- $h(x) = \langle z, x \rangle + D_\psi(x, b)$; $y = \text{argmin}_{x \in X} h(x)$; $\forall x \in X$
 $h(x) \geq \langle z, y \rangle + D_\psi(y, b) + D_\psi(x, y)$

CS726 Cheat Sheet Final

IV. Constrained, Projection-Based Optimization

Def 1.1 (Normal Cone) X closed convex, $\forall x \in X$, $N_X(x)$ defined = $\{z \in \mathbb{R}^d : \langle z, y-x \rangle \leq 0, \forall y \in X\}$

Local solution: $(x^* \in X) - \nabla f(x^*) \in N_X(x^*)$, $\forall y \in X$

(Stationary) $\langle -\nabla f(x^*), y-x^* \rangle \geq 0 \Leftrightarrow \langle \nabla f(x^*), y-x^* \rangle \leq 0$

Thm 1.2 Given (P), $x^* \in X$ local sol $\Rightarrow -\nabla f(x^*) \in N_X(x^*)$; f convex $\Rightarrow -\nabla f(x^*) \in N_X(x^*)$ iff x^* global sol

f str. conv. $\Rightarrow \exists x^*$ s.t. $-\nabla f(x^*) \in N_X(x^*)$ unique sol

2. Def Euclidean Projection $P_X(x) = \operatorname{argmin}_{y \in X} \|y-x\|_2$

Lemma 2.1 X closed conv, $\forall x \in \mathbb{R}^d$, $P_X(x)$ exists & unique

$-(P_X(x)-x) \in N_X(x)$, $\forall y \in X: \langle P_X(x)-x, y-P_X(x) \rangle \geq 0$

Example 2.2 $X = \{x \in \mathbb{R}^d : x \geq 0\}$, $P_X(x) = \max\{x, 0\}$ element-wise

Example 2.3 $X = \{x \in \mathbb{R}^d : a \leq x \leq b\}$, $P_X(x) = \max\{a, \min\{x, b\}\}$

Example 2.4 $X = \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$, $P_X(x) = \frac{x}{\|x\|_2}$ if $\|x\|_2 > 1$

Example 2.5 ℓ_1 ball, prob simpl ex: compute $P_X(x)$ in $O(d \log d)$

Lemma 2.6 X closed conv, $\langle P_X(x)-y, x-y \rangle \geq 0 \forall y \in X$ equality iff $y = P_X(x)$

Lemma 2.7 X closed conv, P_X is nonexpansive $\forall x, y \in \mathbb{R}^d$
 $\|P_X(x)-P_X(y)\|_2 \leq \|x-y\|_2$

3. PGD $x_{k+1} = \operatorname{argmin}_{y \in X} \{f(x_k) + \langle \nabla f(x_k), y-x_k \rangle + \frac{1}{2} \|y-x_k\|_2^2\}$
 $= \operatorname{argmin}_{y \in X} \{ \frac{L}{2} \|y-x_k + \frac{1}{L} \nabla f(x_k)\|_2^2 \}$
 $= P_X(x_k - \frac{1}{L} \nabla f(x_k))$

$x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k) + P_X(x_k - \frac{1}{L} \nabla f(x_k))$

$= x_k - \frac{1}{L} L(x_k - P_X(x_k - \frac{1}{L} \nabla f(x_k)))$

Def 3.1 X closed conv, $\eta > 0$, G_η mapping defined by

$G_\eta(x) = \eta(x - P_X(x - \frac{1}{\eta} \nabla f(x)))$

PGD: $x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k)$

Lemma 3.2 X closed conv, \min L -smooth f , $\bar{x} = P_X(x - \frac{1}{L} \nabla f(x))$

$(G_\eta(x) = \eta(x - \bar{x}))$. If for some $\varepsilon > 0$, $\|G_\eta(x)\|_2 \leq \varepsilon$:

$-\nabla f(\bar{x}) \in N_X(\bar{x}) + B(0, \varepsilon(\frac{1}{\eta} + 1))$

$\Rightarrow (\forall y \in X) : \langle \nabla f(\bar{x}), y-\bar{x} \rangle \geq -\varepsilon(\frac{1}{\eta} + 1) \|y-\bar{x}\|_2$

$\Rightarrow (\forall y \in X \cap B(\bar{x}, 1)) : \langle \nabla f(\bar{x}), y-\bar{x} \rangle \geq -\varepsilon(\frac{1}{\eta} + 1)$

$G_\eta(x) = 0 \Rightarrow x = \bar{x}$ and $-\nabla f(x) \in N_X(\bar{x})$

Lemma 3.3 (Descent Lemma PGD) f L -smooth, X closed conv,

$\bar{x} = P_X(x - \frac{1}{L} \nabla f(x))$ ($x \in X$), then $G_\eta(x) = \eta(x - \bar{x})$.

If $\eta \geq L$, then $f(\bar{x}) \leq f(x) - \frac{1}{2\eta} \|G_\eta(x)\|_2^2$

Thm 3.4 (Convergence Nonconvex f PGD) f L -smooth, X closed conv
 $x_0 \in X$, $\eta \geq L$, $f_k \rightarrow \infty$. $\min_{x \in X} \|G_\eta(x)\|_2^2 \leq \frac{2\eta(f(x_0) - f^*)}{k+1}$
 $\Rightarrow \forall \varepsilon > 0$, $k \geq \lceil \frac{2\eta(f(x_0) - f^*)}{\varepsilon^2} \rceil$, $\exists \bar{x} \in \{x_k, \dots, x_{k+1}\}$
s.t. $-\nabla f(\bar{x}) \in N_X(\bar{x}) + B(0, \varepsilon)$

3.2 (Convergence Convex f PGD) $f(x_{k+1}) - f(x_k) \leq f(x_{k+1}) - f(x_k) + \langle \nabla f(x_k), x^* - x_{k+1} \rangle$
 $\leq \frac{\eta}{2} (\|x_{k+1} - x_k\|_2^2 - \|x^* - x_{k+1}\|_2^2 + \|x^* - x_k + x_{k+1} - x_k\|_2^2)$

$f(x_{k+1}) - f(x^*) \leq \eta \|x^* - x_k\|_2^2$, $\forall \varepsilon > 0$, $f(x_k) - f(x^*) \leq \varepsilon$

After $k = \lceil \frac{2\eta(f(x_0) - f^*)}{\varepsilon^2} \rceil$

V. Projection-free (Frank-Wolfe) Method for Constrained Convex Opt
produce sparse sol.

2. Basic FW $x_{k+1} = \operatorname{argmin}_{x \in X} \langle \nabla f(x_k), x \rangle$; $x_{k+1} = \frac{A_{k+1}}{A_k} x_k + \frac{a_k}{A_k} v_k$

(Analysis) $f(x^*) \geq \min_{x \in X} \{ \frac{1}{A_k} \sum_{i=1}^k a_i \langle \nabla f(x_i), x-x_i \rangle + \langle \nabla f(x_k), x-x_k \rangle \}$ ($U_k = f(x_{k+1})$)

$= \frac{1}{A_k} \sum_{i=1}^k a_i \langle \nabla f(x_i), x-x_i \rangle + \frac{1}{A_k} \sum_{i=1}^k a_i \min_{x \in X} \langle \nabla f(x_i), x-x_i \rangle$

$= \frac{1}{A_k} \sum_{i=1}^k a_i f(x_i) + \frac{1}{A_k} \sum_{i=1}^k a_i \langle \nabla f(x_i), v_i - x_i \rangle = L_k$

$A_k G_k - A_{k-1} G_{k-1} \leq \frac{a_k^2 L}{2 A_k} \|v_k - x_k\|_2^2 \leq \frac{a_k^2 L}{2 A_k} D^2$

$A_0 G_0 \leq \frac{a_0^2 L}{2 A_0} D^2$
 $f(x_{k+1}) - f(x^*) \leq G_k \leq \frac{L D^2}{2 A_k} \sum_{i=1}^k \frac{a_i^2}{A_i} \leq \frac{2 L D^2}{k+2} \sum_{i=1}^k \frac{a_i^2}{A_i}$ choose $a_i = \frac{1}{i^2}$

Lemma 3.1 f L -smooth, X closed conv, alg access the feasible set

X only via linear min oracle requires at least

$\min \{ \frac{d}{\varepsilon}, \frac{L D^2}{\varepsilon^2} \}$ iters to construct \hat{x} s.t. $f(\hat{x}) - \min_{x \in X} f(x) \leq \varepsilon$

for $\varepsilon > 0$, Lower bound applies even if f str. conv.

VI. Nonsmooth Convex Optimization

Def 2.1 f is subdifferentiable at x if $\exists g_x \in \mathbb{R}^d$ s.t. $\forall y \in \mathbb{R}^d$, $f(y) \geq f(x) + \langle g_x, y-x \rangle$, g_x : subgrad of f at x .

Set of all subgrads of f at x called subdifferential

Set of f at x ($\partial f(x)$). ∂ linear for every x .

Fact 2.2 Every conv lower semicont. fn is subdiff'ble on int dom f .

Thm 2.4 f conv lower semicont. fn. Then f M -Lip. cont. wrt $\|\cdot\|$ on int dom f iff $(\forall x \in \text{int dom } f) (\forall g_x \in \partial f(x)) : \|g_x\|_* \leq M$

Projected Subgrad Descent f M -Lip. cont. wrt $\|\cdot\|_2$ not descent

$x_{k+1} = \operatorname{argmin}_{y \in X} \{ a_k \langle g_x, y \rangle + \frac{1}{2} \|y - x_k\|_2^2 \}$ method

(Analysis) $f(x^*) \geq L_k := \frac{1}{A_k} \sum_{i=1}^k a_i \langle \nabla f(x_i), x^* - x_i \rangle$ $g_{x_i} \in \partial f(x_i)$

$x_k^{\text{out}} = \frac{1}{A_k} \sum_{i=1}^k a_i x_i$ $U_k = \frac{1}{A_k} \sum_{i=1}^k a_i f(x_i) \geq f(x_k^{\text{out}})$

$f(x_k^{\text{out}}) - f(x^*) \leq G_k = \frac{1}{A_k} \sum_{i=1}^k a_i \langle g_{x_i}, x^* - x_i \rangle$

$h_k(y) = a_k \langle g_{x_k}, y \rangle + \frac{1}{2} \|y - x_k\|_2^2$, $x_{k+1} = \operatorname{argmin}_{y \in X} h_k(y)$

h_k quadratic $\Rightarrow h_k(y) = h_k(x) + \langle \nabla h_k(x), y-x \rangle + \frac{1}{2} \|y-x\|_2^2$

(bound) $h(x^*) \geq h(x_{k+1}) + \frac{1}{2} \|x_{k+1} - x_k\|_2^2 \geq 0$ for x_{k+1}

$-a_k \langle g_{x_k}, x^* - x_k \rangle \leq -\frac{1}{2} \|x^* - x_k\|_2^2 + \frac{1}{2} \|x^* - x_k\|_2^2 + \frac{G_k M^2}{2}$

$A_k G_k = -\sum_{i=1}^k a_i \langle g_{x_i}, x^* - x_i \rangle \leq \frac{1}{2} \|x^* - x_0\|_2^2 + \sum_{i=1}^k \frac{a_i M^2}{2}$

$f(x_k^{\text{out}}) - f(x^*) \leq \frac{1}{2} \|x^* - x_0\|_2^2 + \sum_{i=1}^k \frac{a_i M^2}{2}$

Choose $a_k = \frac{1}{k+1}$, $f(x_k^{\text{out}}) - f(x^*) \leq \frac{M \|x^* - x_0\|_2^2}{k+1}$

Choose $a_k = \frac{1}{M \sqrt{k+1}}$, $f(x_k^{\text{out}}) - f(x^*) \leq \frac{M D}{\sqrt{k+1}}$

Choose $a_k = \frac{1}{M \sqrt{k+1}}$, $f(x_k^{\text{out}}) - f(x^*) = O(\frac{M D \log(k+1)}{\sqrt{k+1}})$

Choose $a_k = \frac{1}{\sqrt{k+1}}$, $\dots = O(\frac{\|x^* - x_0\|_2^2 + M^2}{\sqrt{k+1}} \log(k+1))$

VII. Stochastic Convex Optimization $\sqrt{k+1}$

f is M -Lip cont. $\forall x \in X$, estimate subgrad $\hat{g}(x, \xi)$ $\xi \sim \mathcal{P}$ i.i.d.

properties: $\mathbb{E} \xi \sim \mathcal{P}$ $\mathbb{E} \hat{g}(x, \xi) = g_x \exists g_x \in \partial f(x)$ (unbiased estimate)

$\mathbb{E} \|\hat{g}(x, \xi) - g_x\|_2^2 \leq \sigma^2 < \infty$ (bounded variance)

S-PSGD (Stochastic Projected Subgrad Descent)

$x_{k+1} = \operatorname{argmin}_{y \in X} \{ a_k \langle \hat{g}(x_k, \xi_k), y \rangle + \frac{1}{2} \|y - x_k\|_2^2 \}$

$x_k^{\text{out}} = \frac{1}{A_k} \sum_{i=1}^k a_i x_i$, $U_k = \frac{1}{A_k} \sum_{i=1}^k a_i f(x_i) \geq f(x_k^{\text{out}})$

$L_k = \frac{1}{A_k} \sum_{i=1}^k a_i \langle \nabla f(x_i), x^* - x_i \rangle + \langle g_{x_i}, x^* - x_i \rangle$

$\Rightarrow f(x_k^{\text{out}}) - f(x^*) \leq G_k = -\frac{1}{A_k} \sum_{i=1}^k a_i \langle g_{x_i}, x^* - x_i \rangle$

$\Rightarrow -a_k \langle g_{x_k}, x^* - x_k \rangle = a_k \langle \hat{g}(x_k, \xi_k), x_k - x^* \rangle + a_k \langle g_{x_k} - \hat{g}(x_k, \xi_k), x_k - x^* \rangle$

$\Rightarrow a_k \langle \hat{g}(x_k, \xi_k), x_k - x^* \rangle \leq a_k \langle \hat{g}(x_k, \xi_k), x_k - x_{k+1} \rangle + a_k \langle g_{x_k} - \hat{g}(x_k, \xi_k), x_k - x_{k+1} \rangle$

$\mathbb{E} \langle \hat{g}(x_k, \xi_k), x_k - x^* \rangle \leq a_k \mathbb{E} \langle \hat{g}(x_k, \xi_k), x_k - x_{k+1} \rangle + a_k \mathbb{E} \langle g_{x_k} - \hat{g}(x_k, \xi_k), x_k - x_{k+1} \rangle$

$\Rightarrow \mathbb{E} \langle \hat{g}(x_k, \xi_k), x_k - x^* \rangle \leq a_k \mathbb{E} \langle \hat{g}(x_k, \xi_k), x_k - x_{k+1} \rangle + a_k \mathbb{E} \langle g_{x_k} - \hat{g}(x_k, \xi_k), x_k - x_{k+1} \rangle$

VIII. Second Order Methods

Basic Newton: $x_{k+1} = x_k - (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$

Local quadratic convergence of Newton's method:

$\{x_k\}$ seq converges to x^*

1) Q-linear, if $\exists r \in (0, 1)$ s.t. $\|x_{k+1} - x^*\| \leq r \|x_k - x^*\| \forall k$ large

2) Q-quad, if \exists const $M > 0$ s.t. $\|x_{k+1} - x^*\| \leq M \|x_k - x^*\|^2 \forall k$ large

Thm 3.5 Suppose f twice cont diff and $\nabla^2 f(x)$ Lip-cont. in nbhd of x^* ,

s.t. $\nabla f(x^*) = 0$ and $\nabla^2 f(x^*) \succ 0$. x_k evolve basic Newton:

(i) x_0 "sufficiently close" to x^* , $\{x_k\}$ converges to x^* ,

(ii) rate of convergence of $\{x_k\}$ quadratic

(iii) sequence of gradient norms $\{\|\nabla f(x_k)\|\}$ converges quadratically to 0.

Global converge $f(x_k) - f(x^*) \leq \epsilon$ after $\mathcal{C}(L_H, m, M, x_0, x^*) + \log(\frac{1}{\epsilon})$

(amped) $x_{k+1} = x_k - \alpha_k \nabla^2 f(x_k)^{-1} \nabla f(x_k)$ α_k : line search

For $\nabla^2 f(x_k) \succ 0$, consider Quasi-Newton:

$x_{k+1} = x_k - \alpha_k B_k^{-1} \nabla f(x_k)$ where $B_k \succeq \delta I$ for $\delta > 0$

$\alpha_k = -B_k^{-1} \nabla f(x_k)$ satisfies Weak Wolfe conditions: $\exists 0 < c_1 < c_2 < 1$

(WW1) $f(x_k + \alpha_k d_k) \leq f(x_k) + c_1 \alpha_k \langle \nabla f(x_k), d_k \rangle$

(WW2) $\langle \nabla f(x_k + \alpha_k d_k), d_k \rangle \geq c_2 \langle \nabla f(x_k), d_k \rangle$

(WW1) $f(x_k + \alpha_k p_k) \leq f(x_k) + c_1 \alpha_k \nabla f_k^T p_k$

(WW2) $\nabla f(x_k + \alpha_k p_k)^T p_k \geq c_2 \nabla f_k^T p_k$

$\Rightarrow f(x_{k+1}) \leq f(x_k) - c_1 \frac{1-c_2}{L} \frac{\langle \nabla f(x_k), p_k \rangle^2}{\|p_k\|^2}$ f L-smooth p_k direction

Thm 3.6 f twice cont. diff $x_{k+1} = x_k + \alpha_k d_k$, α_k satisfies WW1 & 2.

If the sequence $\{x_k\}$ converges to x^* s.t. $\nabla f(x^*) = 0$ and

$\nabla^2 f(x^*) \succ 0$ and if satisfies $\lim_{k \rightarrow \infty} \frac{\|\nabla f(x_k) + \nabla^2 f(x_k) d_k\|}{\|d_k\|} = 0$

(super-linear) then

(i) $\alpha_k = 1$ is admissible $\forall k > k_0$ and (reverse: 3.7) (QN $\alpha_k = 1$)

(ii) if $\alpha_k = 1 \forall k > k_0$, $\{x_k\}$ converges to x^* super-linearly

Proof of Thm 3.5 Goal: if $\|x_0 - x^*\|$ suff. small, then

1) $\exists M > 0$ constant s.t. $\forall k: \|x_{k+1} - x^*\| \leq M \|x_k - x^*\|^2$

2) $\exists M' > 0$ also const s.t. $\forall k: \|\nabla f(x_{k+1})\| \leq M' \|\nabla f(x_k)\|^2$

1) $x_{k+1} = x_k - (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$, x_k in nbhd of x^* , $\nabla^2 f$ LH-Lip

$$\|x_{k+1} - x^*\| = \|x_k - x^* - (\nabla^2 f(x_k))^{-1} \nabla f(x_k)\|$$

$$= \|(\nabla^2 f(x_k))^{-1} (\nabla^2 f(x_k)(x_k - x^*) - \nabla f(x_k))\|$$

$$= \|(\nabla^2 f(x_k))^{-1} (\nabla^2 f(x_k)(x_k - x^*) - (\nabla f(x_k) - \nabla f(x^*)))\|$$

$$\leq \|(\nabla^2 f(x_k))^{-1}\| \cdot \|\nabla^2 f(x_k)(x_k - x^*) - (\nabla f(x_k) - \nabla f(x^*))\|$$

$$\|x_{k+1} - x^*\| \leq \|(\nabla^2 f(x_k))^{-1}\| \cdot \int_0^1 \|\nabla^2 f(x_k) - \nabla^2 f(x_k + t(x^* - x_k))\| \|x_k - x^*\| dt$$

$$\leq \|(\nabla^2 f(x_k))^{-1}\| \cdot \int_0^1 L_H \cdot t \|x^* - x_k\| dt$$

$$\leq \|(\nabla^2 f(x_k))^{-1}\| \cdot \frac{L_H}{2} \|x^* - x_k\|^2 = M \|x^* - x_k\|^2$$

2) $\nabla f(x_k) = \nabla^2 f(x_k) (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$

$$x_{k+1} - x_k = -(\nabla^2 f(x_k))^{-1} \nabla f(x_k)$$

$$\nabla^2 f(x_k)(x_{k+1} - x_k) = -\nabla f(x_k)$$

$$\|\nabla f(x_{k+1})\| = \|\nabla f(x_{k+1}) - \nabla f(x_k) + \nabla f(x_k)\|$$

$$= \left\| \int_0^1 \nabla^2 f(x_k + t(x_{k+1} - x_k))(x_{k+1} - x_k) dt - \nabla^2 f(x_k)(x_{k+1} - x_k) \right\|$$

$$\leq \int_0^1 \|\nabla^2 f(x_k + t(x_{k+1} - x_k)) - \nabla^2 f(x_k)\| \|x_{k+1} - x_k\| dt$$

$$\leq \frac{L_H}{2} \|x_{k+1} - x_k\|^2 \leq \frac{L_H}{2} \|(\nabla^2 f(x_k))^{-1}\|^2 \|\nabla f(x_k)\|^2$$

$$\leq 2L_H \|(\nabla^2 f(x_k))^{-1}\|^2 \|\nabla f(x_k)\|^2$$

$$= M' \|\nabla f(x_k)\|^2 \quad M' = 2L_H \|(\nabla^2 f(x_k))^{-1}\|^2 \text{ const}$$

secant equation: $B_{k+1} s_k = y_k$

$$s_k = x_{k+1} - x_k = \alpha_k p_k$$

$$y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$$

curvature: $s_k^T y_k > 0$