

Lecture 9 - scribbles - convex analysis

Tuesday, February 20, 2018

14:31

today: convex optimization

Motivation: $\min_w \frac{\lambda \|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n f(x^{(i)}, y^{(i)}, w)$

convex analysis recap:

$f: \mathbb{R}^d \rightarrow \mathbb{R}$

f is convex \Leftrightarrow

$f(px + (1-p)y) \leq pf(x) + (1-p)f(y) \quad \forall x, y \in \mathbb{R}^d$

convex combination between x & y

$y + p(x - y)$

$y \xrightarrow{p} x$

if f is differentiable at x and convex

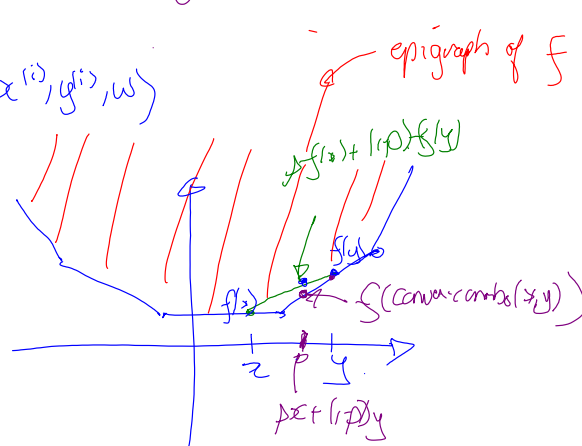
$\Rightarrow f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \quad \forall y$

(f convex)

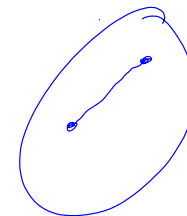
1. convexity 2. subdifferentiability 3. differentiability

subdifferentiable of f at x

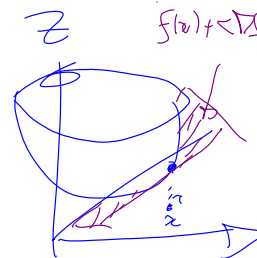
convex surrogate loss



epigraph of f



epigraph $f = \{(x, y) : y \geq f(x)\}$



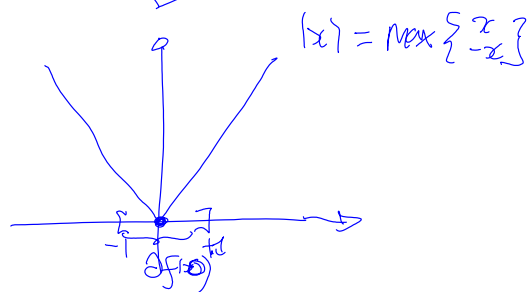
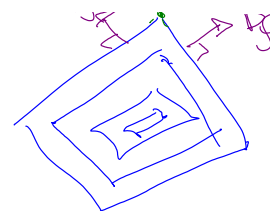
$f(x) + \langle \nabla f(x), y - x \rangle$



$v \in \partial f(x)$

subgradient: v of f at x : $v \in \partial f(x)$

$$\Leftrightarrow y \in \text{dom}(f) , f(y) \geq f(x) + \langle v, y-x \rangle$$



Some standard assumptions:

f is μ -strongly convex $\Leftrightarrow f(y) \geq f(x) + \langle \partial f(x), y-x \rangle + \frac{\mu}{2} \|x-y\|^2$

$\forall x, y \in \text{dom}(f)$ $\langle v, y-x \rangle$ for any $v \in \partial f(x)$

strong convexity constant

$$\text{dom}(f) \triangleq \{x : f(x) < \infty\}$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$$

$$f^{-1}([a, \infty]) \text{ is open for all } a$$

f is L -smooth i.e. f has a L -Lipschitz gradient for all x

$$\Leftrightarrow \|\nabla f(x) - \nabla f(y)\|_* \leq L \|x-y\|$$

$$(\|\cdot\|_p)^* = \|\cdot\|_q \text{ where } \frac{1}{p} + \frac{1}{q} = 1$$

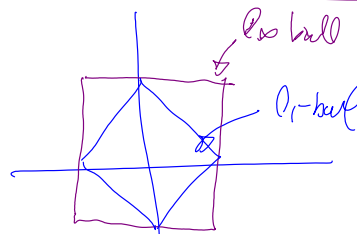
$$p=2 \Rightarrow q=2$$

$$p=1 \Rightarrow q=\infty$$

$$\|w\|_* \triangleq \sup_{\|v\| \leq 1} \langle w, v \rangle$$

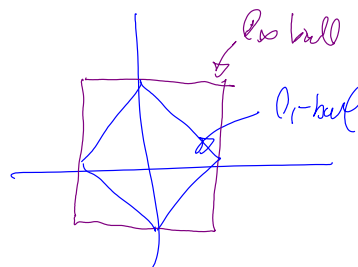
generalized C.S.

$$\langle w, v \rangle \leq \|w\|_* \|v\|$$



$$p=2 \Rightarrow q=2$$

$$p=1 \Rightarrow q=\infty$$



Fundamental descent lemma:

when f is L -Lipschitz (though f is not nec. convex)

$$f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2$$

$$* f(\underbrace{x - \gamma \nabla f(x)}_{y_\gamma}) \leq f(x) - \gamma \langle \nabla f(x), \nabla f(x) \rangle + \frac{\gamma^2}{2} L \|\nabla f(x)\|^2$$

$$= f(x) - \underbrace{\left[\gamma \left(1 - \frac{\gamma L}{2} \right) \right]}_{> 0} \|\nabla f(x)\|^2$$

minimize RHS with respect to γ gives $\boxed{\gamma^* = \frac{1}{L}}$

$$> 0 \Leftrightarrow \boxed{0 < \gamma < \frac{2}{L}}$$

→ [think of 2nd order Taylor expansion]

$$\frac{1}{2} \int_{\gamma=0}^1 \langle y-x, \bar{H}(x+\gamma(y-x)) y-x \rangle d\gamma$$

Hessian $\bar{H}(x,y) \rightarrow$ some pt. between x & y

$$f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \frac{1}{2} \langle y-x, \bar{H}(x,y) y-x \rangle$$

to evaluate $\|H\| \leq L$ in absolute value

$$\lambda_{\min}(H) \|v\|^2 \leq v^T H v \leq \lambda_{\max}(H) \|v\|^2$$

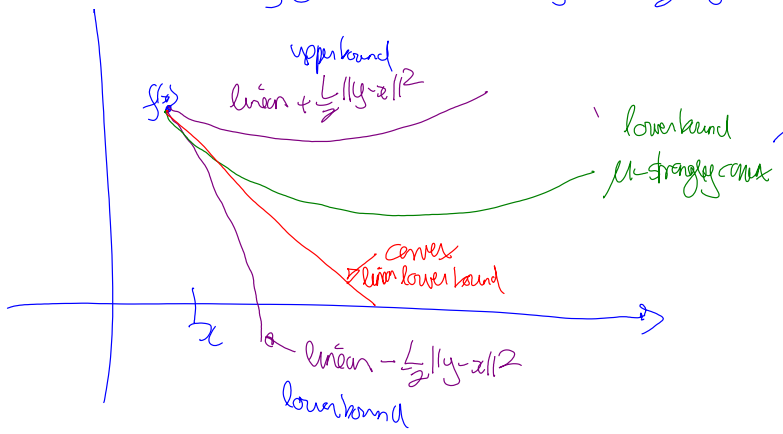
Or

"pro-way" is to use fundamental thm. of calculus

$$f(y) = f(x) + \int_0^1 \frac{d}{d\gamma} f(x + \gamma(y-x)) d\gamma$$

15h41

$$f(y) \leq f(x) + \langle \nabla f(x), y-x \rangle + \frac{L}{2} \|y-x\|^2 \leq \langle \nabla f(x + \delta(y-x)), y-x \rangle$$



$$f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2$$

f twice differentiable

$$L = \lambda_{\max}(H)$$

$$\mu = \lambda_{\min}(H)$$

$$f \text{ is } \mu\text{-strongly convex} \Leftrightarrow f - \frac{\mu}{2} \|\cdot\|^2 \text{ is convex}$$

* gradient descent: $x_{t+1} = x_t - \gamma \nabla f(x_t)$ $\gamma = \frac{1}{L}$

a) $f(x_t) - \min_x f(x) \leq O\left(\frac{L r_0^2}{t}\right)$ when f is convex & L -smooth

note: $\underbrace{\min_x f(x)}_{f^*}$ where $r_0 \geq \text{dist}(x_0, x^*)$ & $x^* = \arg\min_x f(x)$

* no guarantee on $\text{dist}(x_t, x^*)$ (in general) \rightarrow Nesterov lower bounds

"sublinear"

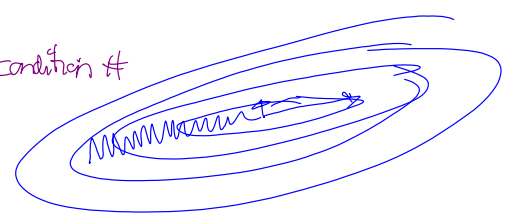
\rightarrow [see Nesterov book for $O(\frac{1}{t})$ rate]

b) if f is μ -strongly convex & L -smooth

"Quadratic rate"

$$f(x_t) - f(x^*) \leq O\left(\exp\left(-\frac{\mu}{L} t\right)\right)$$

$\frac{L}{\mu} \triangleq \text{condition \#}$

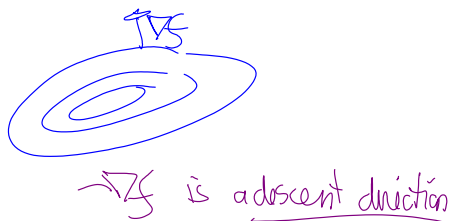


Newton's method $x_{t+1} = x_t - \gamma H(x_t)^{-1} \nabla f(x_t)$

Newton's method $x_{t+1} = x_t - \gamma H(x_t)^{-1} \nabla f(x_t)$

$\gamma \nabla_t \nabla f(x_t)$

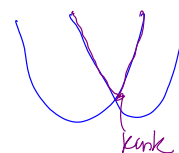
f smooth \Rightarrow smooth sublevel sets



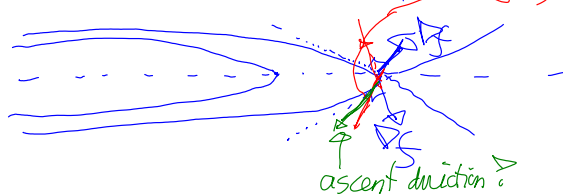
non-smooth f

$-\nabla f(x_t)$ is not necessarily a descent direction

subgradient



Kinks in the sublevel set $-2f(x)$



⊗ subgradient method is not a descent method

but $-\nabla f(x_t)$ is descent direction on $\|x(\gamma) - \tilde{x}\|^2$ for any \tilde{x} in sublevel set

(i.e. $\|x_t - \gamma \nabla f(x_t) - \tilde{x}\|^2 \leq \|x_t - \tilde{x}\|^2$ for γ small enough)

$$x(\gamma) = x_t - \gamma \nabla f(x_t)$$

16h43

thus get closer to x^*

[in non-smooth optimization, $f(x_t)$ can go up and down]

\Rightarrow combine multiple points x_t to get \hat{x}_T

argmin $f(x_t)$ (in batch setting)

\hat{x}_T

weights

or weighted average $\hat{x}_T = \sum_t P_t x_t$ ^{weights}
 for stochastic setting
 or too expensive to compute $f(x)$

* projection operator on closed convex set C $P_C(x) \triangleq \arg \min_{y \in C} \|x - y\|_2^2$

"Euclidean projection"

$P_C(\cdot)$ is a contraction i.e. $\|P_C(x) - P_C(y)\|_2 \leq \|x - y\|_2 \quad \forall x, y$

if $y \in C$, then $P_C(y) = y$

and thus $\|P_C(x) - y\|_2 \leq \|x - y\|_2 \quad \forall y \in C$

Stochastic subgradient method

Setup: want to solve $\min_{x \in C} f(x)$

where $f(x) \triangleq \mathbb{E}_{\xi} [h(x, \xi)]$

assumptions: 1) f & C are convex

2) projection on C is cheap

3) we have a stochastic oracle which gives a random $g(x, \xi)$

$$\text{s.t. } \mathbb{E}_{\xi} [g(x, \xi) | x] = f'(x)$$

↑ any subgradient of f at x

if C is differentiable in x and "well behaved" $\quad \uparrow$

[if f is differentiable in x and "well behaved"]

$$g(x, \xi) \triangleq \nabla_x h(x, \xi)$$

$$\text{then } \mathbb{E}_{\xi} [\nabla_x h(x, \xi)] = \nabla f(x)$$

EM example: $f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \rightarrow g(x^{(t)}, y^{(t)}, \overset{\text{parameter}}{\xi^{(t)}}) + \frac{\lambda \|x^{(t)}\|^2}{2}$

$$h(x, \xi) = f_{\xi}(x)$$

\uparrow
 $\{1, \dots, n\}$

at step t , sample $\xi_t \overset{\text{indep}}{\sim} \xi_1, \dots, \xi_n$

$$\text{use } g_t \triangleq g(x_t, \xi_t) \triangleq \nabla f_{\xi_t}(x_t)$$

$$\text{here } \mathbb{E}_{\xi} [\nabla f_{\xi}(x) | x] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x) = \nabla f(x)$$

$$4) \mathbb{E} \|g(x, \xi)\|^2 \leq B^2 \quad (\text{finite variance condition})$$

$$[\text{if } \|h'(x, \xi)\| \leq B]$$

\downarrow
is sufficient

algorithm:

$x_0 \in \mathbb{C}$ initialization

for $t=0, \dots, T-1$

let g_t be $g(x_t, \xi_t)$ (from oracle)

$$\text{let } x_{t+1} = \mathbb{P}_{\mathbb{C}} [x_t - \underset{\substack{\uparrow \\ \text{step size}}}{\eta_t} g_t]$$

end

output: $\hat{x}_T \triangleq \sum_{t=0}^{T-1} \beta_{T,t} x_t$

"weighted average"

where β_t are some
convex combo
coefficients
 $\sum_{t=0}^{T-1} \beta_t = 1$ $\beta_t \geq 0$

convergence proof:

important inequality

$$f(y) \geq f(x) + \langle f'(x), y-x \rangle + \frac{\mu}{2} \|y-x\|^2 \quad (\mu \geq 0)$$

$$\Rightarrow -\langle f'(x), x-y \rangle \leq -(f(x) - f(y) + \frac{\mu}{2} \|y-x\|^2) \quad (+)$$

use it on: $-\langle f'(x_t), x_t - \tilde{x} \rangle$

$$x_{t+1} = \mathcal{P}_C [x_t - \gamma_t g_t]$$

$$\|x_{t+1} - \tilde{x}\|^2 \leq \|x_t - \gamma_t g_t - \tilde{x}\|^2$$

any feasible point
i.e. $\tilde{x} \in C$

$$\stackrel{\mathcal{P}_C}{=} \|x_t - \tilde{x}\|^2 + \gamma_t^2 \|g_t\|^2 - 2\gamma_t \langle g_t, x_t - \tilde{x} \rangle$$

$$\mathbb{E}[\|x_{t+1} - \tilde{x}\|^2 | x_t] \leq \|x_t - \tilde{x}\|^2 + \gamma_t^2 \mathbb{E}[\|g_t\|^2 | x_t] - 2\gamma_t \underbrace{\langle \mathbb{E}[g_t | x_t], x_t - \tilde{x} \rangle}_{f'(x_t)}$$

(using (+)) \downarrow

\leq \downarrow

$$\leq \|x_t - \tilde{x}\|^2 - 2\gamma_t \left[f(x_t) - f(\tilde{x}) + \frac{\mu}{2} \|x_t - \tilde{x}\|^2 \right]$$

$$\mathbb{E}[\mathbb{E}[\cdot | x_t]] = \mathbb{E}[\cdot]$$

$$\mathbb{E}[\|x_{t+1} - \tilde{x}\|^2] \leq (1 - \mu \gamma_t) \mathbb{E}[\|x_t - \tilde{x}\|^2] - 2\gamma_t [\mathbb{E}[f(x_t)] - f(\tilde{x})] + \gamma_t^2 \mathbb{E}[\|g_t\|^2]$$

γ_t small enough, we have $\mathbb{E}[\|x_t - \tilde{x}\|^2]$ decreases
for any $\tilde{x} \in C$ s.t. $f(\tilde{x}) \leq \mathbb{E}[f(x_t)]$

* non-strongly setting: $\mu=0$

let \tilde{x} be some min point. x^* , i.e. $f(\tilde{x}) = \min_{x \in C} f(x)$

$$\text{let } r_t \triangleq \mathbb{E} \|x_t - x^*\|^2$$

$$\text{let } \boxed{\varepsilon_t \triangleq \mathbb{E}[f(x_t) - f(x^*)]} \quad \text{expected suboptimality error}$$

$$r_{t+1} \leq r_t - 2\gamma_t \varepsilon_t + \gamma_t^2 B^2$$

$$\Rightarrow 2\gamma_t [\varepsilon_t] \leq r_t - r_{t+1} + \gamma_t^2 B^2 \quad \forall t$$

* sum inequalities from $t=0$ to $t=T$

$$\Rightarrow 2 \underbrace{\sum_{t=0}^T \gamma_t \varepsilon_t}_{\downarrow} \leq \underbrace{r_0 - r_{T+1}}_{\text{telescoping sum } \sum_{t=0}^T (r_t - r_{t+1})} + \left(\sum_{t=0}^T \gamma_t^2 \right) B^2$$

$$2 \left(\sum_{t=0}^T \gamma_t \right) \min_{0 \leq t \leq T} \varepsilon_t \leq \downarrow$$

$$\Rightarrow \boxed{\min_{0 \leq t \leq T} \varepsilon_t \leq \frac{r_0 + \left(\sum_{t=0}^T \gamma_t^2 \right) B^2}{2 \sum_{t=0}^T \gamma_t}}$$

use $\gamma_t = \frac{r_0}{B\sqrt{T+1}}$ to minimize RHS

$$\Rightarrow \boxed{\min_{0 \leq t \leq T} \varepsilon_t \leq \frac{Br_0}{\sqrt{T+1}}}$$

* since f is convex,

$$f\left(\sum_t p_t x_t\right) \leq \sum_t p_t f(x_t)$$