

Lecture 10 - subgradient method

Thursday, February 6, 2020 13:30

today:
 • fundamental descent lemma
 • subgradient method

Fundamental descent lemma:

when ∇f is L -Lipschitz (lemma holds even if f is not convex)

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

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

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

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

→ minimize RHS with respect to δ

gives $\boxed{\delta^* = \frac{1}{L}}$

$$f(y_\delta) \leq f(x) - \frac{1}{2L} \|\nabla f(x)\|^2$$

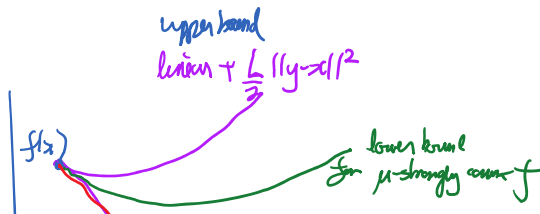
descent lemma:

think of 2nd order Taylor expansion

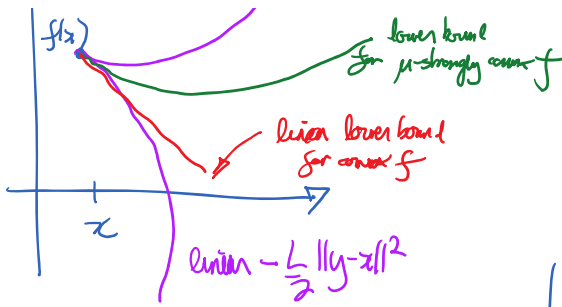
$$f(y) = f(x) + \langle \nabla f(x), y-x \rangle + \frac{1}{2} \int_{\delta=0}^1 \langle y-x, \underbrace{H(x + \delta(y-x))}_{\text{Hessian of } f} \rangle y-x \, d\delta$$

integral form of remainder
 top e-value of H in absolute value $\leq L$

$$\forall v \|Hv\| \leq \lambda_{\max}(H) \|v\|^2 \leq L \|y-x\|^2$$



when f is twice differentiable



when f is twice differentiable
 $L = \lambda_{\max}(H_f)$
 $\mu = \lambda_{\min}(H_f)$

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

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

a) when f is convex & L -smooth

$$f(x_t) - \underbrace{\min_x f(x)}_{f^*} \leq O\left(\frac{Lr_0^2}{t}\right)$$

where $r_0 \geq \text{dist}(x_0, X^*)$
 $X^* = \text{argmin}_x f(x)$

note: no guarantee on $\text{dist}(x_t, X^*)$
 (for general L -smooth convex f 's for $t \leq \text{dim}(x_t)$)

"sublinear rate"

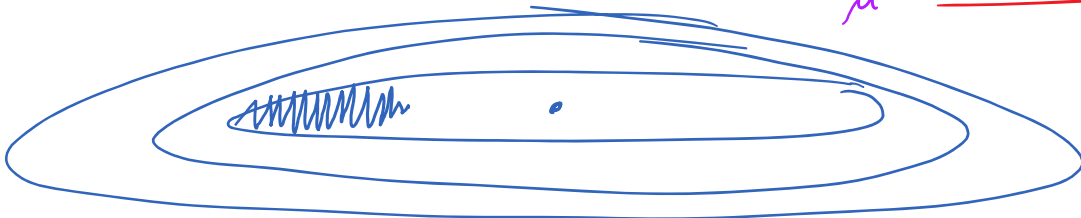
[see Nesterov book for proof]

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

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

"linear rate"

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



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

14h27

subgradient method

non-descent methods

f smooth \Rightarrow smooth sublevel sets



non-smooth f

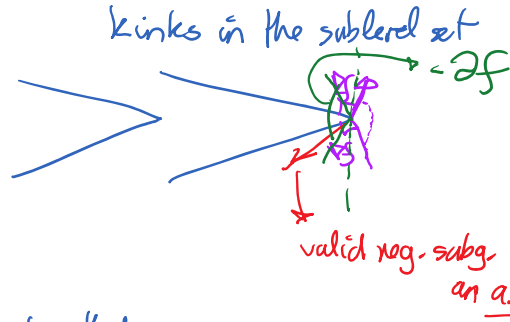
$$-\underbrace{f'(x_t)}_p$$

is not nec. a descent direction



∇f is a descent direction

$-\nabla f(x_k)$ is not nec. a descent direction
 ↑
 a subgradient



⊛ subgradient method is not a descent method

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

$x(\delta)$ gets closer to any \tilde{x} s.t. $f(\tilde{x}) \leq f(x_k)$ for small enough δ
 thus gets closer to any x^*

* in non-smooth optimization, $f(x_k)$ can go up & down

to stabilize

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

argmin $f(x_t)$ [in batch setting]
 $\{x_t\}_{t=1}^T$

weighted average $\hat{x}_T = \sum_t p_t x_t$ [for stochastic setting]
 convex weights or when too expensive to compute $f(x_t)$

* projection operator on a closed convex set C

$$P_C(x) \triangleq \operatorname{argmin}_{y \in C} \|x - y\|_2^2$$

"Euclidean projection" of x on C

$P_C(\cdot)$ is non-expansive i.e. $\|P_C(x) - P_C(y)\|_2 \leq \|x - y\|_2 \forall x, y$

• if $y \in C$, then $P_C(y) = y$

can 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 $g(x, \xi)$ for random ξ

st. $\mathbb{E}_{\xi} [g(x, \xi) | x] = f'(x)$
some subgradient of f at x

[example:

a) f is differentiable in x & "well behaved"

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

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

b) ERM example: $f(x) = \int \sum_{i=1}^n f_i(x)$ eg $f_i(x) = \mathcal{L}(x^{(i)}, y^{(i)}; "x") + \frac{\Delta}{2} \|x^{(i)}\|^2$
parameter \downarrow

$h(x, \xi) \triangleq f_{\xi}(x)$
 $\hookrightarrow \xi \in \{1, \dots, n\}$

at step t , sample $i_t \stackrel{\text{unif}}{\sim} \{1, \dots, n\}$

use $g_t \triangleq g(x_t, i_t) \triangleq f'_{i_t}(x_t)$

here $\mathbb{E}_{\xi} [f'_{\xi}(x) | x] = \int \sum_{i=1}^n f_i'(x) = f'(x)$

4) $\mathbb{E} \|g(x, \xi)\|_2^2 \leq B^2$ (finite variance condition)

this replaces the Lipschitz gradient assumption

(if f is Lipschitz instead)

[sufficient condition: $\|h(x,y)\| \leq B \|x-y\|$]

stochastic subgradient method
algorithm

$x_0 \in C$ initialization

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

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

let $x_{t+1} = P_C [x_t - \alpha_t g_t]$

↑
step size

end
output

$$\hat{x}_T \triangleq \sum_{t=0}^T P_{T,t} x_t$$

"weighted average"

where $P_{T,t}$ are some
convex comb. coeffs

$$\sum_t P_{T,t} = 1 \quad P_{T,t} \geq 0$$

convergence proof:

important inequality

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

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

→ use this on $-\langle f'(x_t), x_t - x^* \rangle$

$x_{t+1} = P_C(x_t - \alpha_t g_t)$ by def.

$$\|x_{t+1} - \tilde{x}\|^2 \stackrel{\text{by } P_C}{\leq} \|x_t - \alpha_t g_t - \tilde{x}\|^2$$

any feasible pt.
 $\tilde{x} \in C$

$$= \|x_t - \tilde{x}\|^2 + \alpha_t^2 \|g_t\|^2 - 2\alpha_t \langle g_t, x_t - \tilde{x} \rangle \quad [\text{valid } \forall \tilde{x} \in C]$$

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

$$\mathbb{E}[\|x_{t+1} - \tilde{x}\|^2 | x_t] \leq \|x_t - \tilde{x}\|^2 + \delta_t^2 \mathbb{E}[\|g_t\|^2 | x_t] - 2\delta_t \langle \mathbb{E}[g_t | x_t], x_t - \tilde{x} \rangle$$

\downarrow $\downarrow \sqrt{B^2}$ $\downarrow f'(x_t)$

$$\stackrel{\text{using } (f)}{\leq} \|x_t - \tilde{x}\|^2 + \delta_t^2 B^2 - 2\delta_t [f(x_t) - f(\tilde{x}) + \frac{\mu}{2} \|x_t - \tilde{x}\|^2]$$

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

(also true for $\mu=0$)

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