

# Lecture 13 Basic Optimization

Friday, October 14, 2016 1:20 PM

Goal:  $\min f(x) \leftarrow \text{objective}$

$x \in B \leftarrow \text{constraint}$

usually  $B \subseteq \mathbb{R}^n, f: \mathbb{R}^n \rightarrow \mathbb{R}$

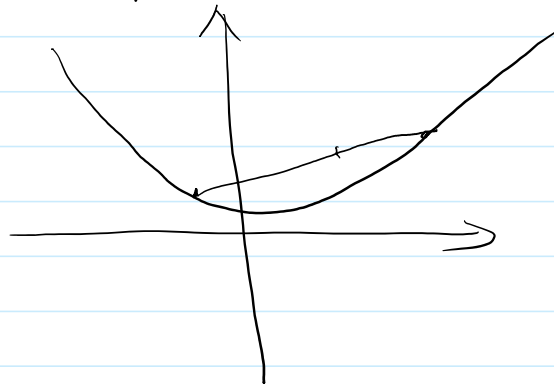
## - Convexity

- idea: if  $x, y$  are valid,  $\alpha x + (1-\alpha)y$  is also valid ( $\alpha \in [0, 1]$ )

- set  $B$  is convex, if  $\forall x, y \in B$   
 $\alpha x + (1-\alpha)y \in B \quad (\alpha \in [0, 1])$

- function  $f(x)$  is convex, if  $\forall x, y \in B$   
 $f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$   
( $\alpha \in [0, 1]$ )

- convex optimization: both  $f(x)$  and  $B$  are convex



## - Basic Algorithm: Gradient Descent

- recap: Gradient  $\nabla f(x) \in \mathbb{R}^n$ , 1st order derivative

$$(\nabla f(x))_i = \frac{\partial}{\partial x_i} f(x)$$

Hessian: 2nd order derivative  $\nabla^2 f(x) \in \mathbb{R}^{n \times n}$

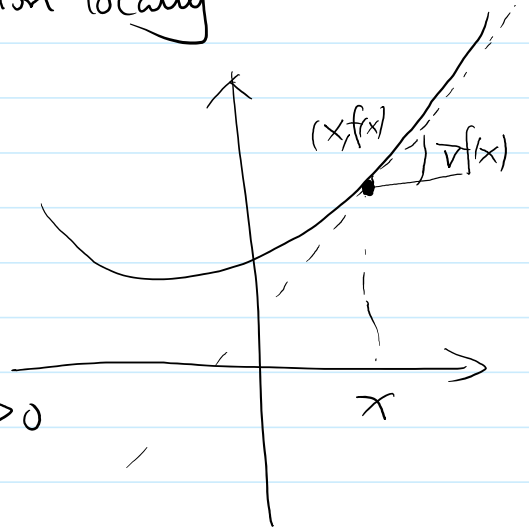
$$(\nabla^2 f(x))_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(x)$$

- General idea in optimization

approximate the objective function locally

- if  $f(x)$  is convex, then

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



→ intuition: if  $\langle \nabla f(x), y - x \rangle > 0$   
 $f(y)$  is always worse.

needs an upper bound to guarantee decrease.

- Lipschitz Gradient / "Smoothness"

Def:  $f(x)$  is  $L$ -smooth ( $L$ -Lipschitz Gradient)

$$\text{if } \forall x, y \quad \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$

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

$$\Leftrightarrow \langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2 \quad (*)$$

- Analyzing Gradient descent

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

$$\text{solution: } y = x - \frac{1}{L} \nabla f(x)$$

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

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

(need to show  $\|\nabla f(x)\|$  large to make progress)

$$x^{k+1} = x^k - \eta \nabla f(x^k) \quad (\eta \in (0, \frac{2}{L}])$$

Let  $r_k = \|x^k - x^*\|$ , first show never gets further

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

$$= r_k^2 - 2\eta \langle \nabla f(x^k), x^k - x^* \rangle + \eta^2 \|\nabla f(x^k)\|^2$$

$\downarrow (*) , \nabla f(x^*) = 0$   
 $\frac{1}{2L} \|\nabla f(x^k)\|^2$

$$\leq r_k^2 - \eta \left( \frac{2}{L} - \eta \right) \|\nabla f(x^k)\|^2 \leq r_k^2$$

$\Rightarrow$  always move closer!

Let  $\Delta_k = f(x^k) - f(x^*)$ , then

$$\Delta_k \leq \underbrace{\langle \nabla f(x^k), x^k - x^* \rangle}_{\text{convexity}} \leq r_k \|\nabla f(x^k)\|$$

$$\leq r_0 \|\nabla f(x^k)\|$$

$$\Delta_{k+1} \leq \Delta_k - \underbrace{\eta \left( 1 - \frac{L\eta}{2} \right)}_{\omega} \|\nabla f(x^k)\|^2 \leq \Delta_k - \frac{\omega}{r_0^2} \Delta_k^2$$

$\uparrow$   
 same as (\*\*\*)

idea: if  $\Delta_k = \frac{r_0^2}{\dots}$ ,  $\Delta_{k+1} \leq \frac{r_0^2}{\dots}$

do induction carefully,

$$f(x^k) - f(x^*) \leq \frac{2L \|x^0 - x^*\|^2}{k+4} \quad \square$$

Strong convexity: better lowerbound

Def:  $f$  is  $\mu$ -strongly convex if  $\forall x, y$

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

Lemma: If  $f$  is  $L$ -smooth,  $\mu$ -strongly convex, then

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{\mu L}{\mu + L} \|x - y\|^2 + \frac{1}{\mu + L} \|\nabla f(x) - \nabla f(y)\|^2$$

Theorem: Choose  $\eta = \frac{2}{\mu + L}$ , then

$$\|x^k - x^*\| \leq \left( \frac{L(\mu - 1)}{L(\mu + 1)} \right)^k \|x^0 - x^*\|$$

## Lecture 14 Stochastic Gradient and Variance Reduction

Sunday, October 16, 2016

10:25 PM

### - Stochastic Gradient descent

#### - Least Squares

$$\min \|y - Ax\|^2$$
$$\frac{1}{2n} \sum_{i=1}^n (y_i - \langle a_i, x \rangle)^2, \quad a_i \in \mathbb{R}^d$$

For simplicity assume  $\|a_i\|=1$

$$y_i = \langle a_i, x^* \rangle + \varepsilon_i \quad \left( \begin{array}{l} \sum \varepsilon_i a_i = 0 \\ |\varepsilon_i| \leq \sigma \end{array} \right)$$

#### - Can rewrite

$$f(x) = f(x^*) + \frac{1}{2} (x - x^*)^T M (x - x^*)$$

where  $M = \frac{1}{n} \sum_{i=1}^n a_i a_i^T$

#### - SGD for least squares

$$f_i(x) = \frac{1}{2} (y_i - \langle a_i, x \rangle)^2$$

pick random  $i$

$$\begin{aligned} x^{t+1} &= x^t - \eta \nabla f_i(x) \\ &= x^t + \eta (y_i - \langle a_i, x^t \rangle) a_i \end{aligned}$$

#### - Analyzing SGD

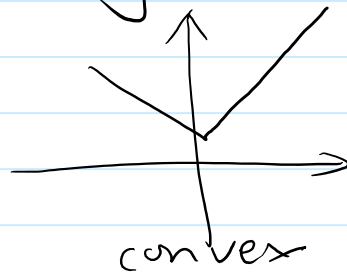
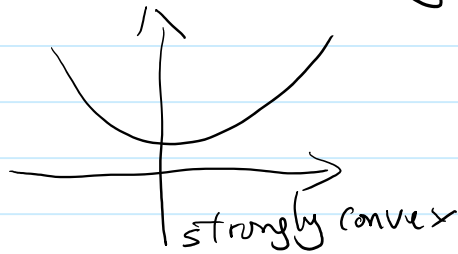
$$x^{t+1} = x^t - \eta \nabla f_i(x^t) = x^t - \eta (\nabla f(x^t) + \zeta_i)$$

$\zeta_i$  independent of  $x^t$ .  $E[\zeta_i] = 0$

$$\begin{aligned}
\text{Let } r_t &= \mathbb{E}[\|x^t - x^*\|^2] \\
r_{t+1}^2 &= r_t^2 - 2\mathbb{E}[\eta \langle \nabla f(x^t) + \zeta_i, x^t - x^* \rangle] \\
&\quad + \eta^2 \mathbb{E}[\|\nabla f(x^t) + \zeta_i\|^2] \\
&= r_t^2 - 2\eta \langle \nabla f(x^t), x^t - x^* \rangle \\
&\quad + \eta^2 \mathbb{E}[(y_i - \langle a_i, x^t \rangle)^2] \\
&= r_t^2 - 2\eta (x^t - x^*)^T M (x^t - x^*) + 2\eta^2 f(x^t) \\
&= r_t^2 - 2\eta (x^t - x^*)^T M (x^t - x^*) \\
&\quad + 2\eta^2 \left( f(x^*) + \frac{1}{2} (x^t - x^*)^T M (x^t - x^*) \right)
\end{aligned}$$

- Suppose  $\sigma_{\min}(M) = \mu$

(this is called Strong Convexity)



$$r_{t+1}^2 \leq \underbrace{r_t^2 - (2\eta - \eta^2) \mu r_t^2}_A + \underbrace{2\eta^2 f(x^*)}_B$$

we want term A to dominate term B!

$$\text{set } \eta \leftarrow \frac{\mu r_t^2}{f(x^*)} \text{ works}$$

in that case

$$r_{t+1}^2 \leq r_t^2 (1 - \eta) \\ \leq r_t^2 \left(1 - \frac{\mu r_t^2}{f(x^*)}\right)$$

again we solve the recursion and get

$$r_t^2 = \frac{f(x^*)}{\mu t} \quad (\text{if the initial point is close enough})$$

in the best case,  $M = \frac{1}{d}$   
(because  $\text{tr}(M) = \frac{1}{n} \sum \|a_i\|^2 = 1$ )

so we can hope to get reasonably close after  $O(d)$  iterations.

- System of linear equations

$$r_{t+1}^2 \leq r_t^2 - \underbrace{(2\eta - \eta^2)\mu r_t^2}_A + \underbrace{2\eta^2 f(x^*)}_B$$

we had to use a small  $\eta$  to let A dominate B.

what if  $f(x^*) = 0$ ? (this means  $y_i = \langle a_i, x^* \rangle$ )

then we can choose  $\eta = 1$  and

$$r_{t+1}^2 \leq (1 - \mu) r_t^2$$

$$\Rightarrow r_t^2 \leq (1 - \mu)^t r_0^2$$

$r_t^2$  decreases by a constant factor every  $\frac{1}{\mu}$  iterations!

- Variance Reduction

- idea: Previously, things did not work because

$$\nabla f_i(x^*) \neq 0$$

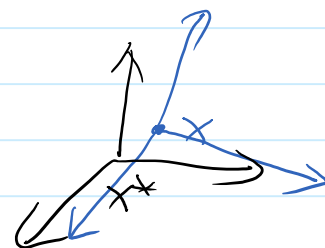
- if we choose a large step size will go away even if we are already at  $x^*$ !



- idea: if  $x$  is very close to  $x^*$

$$\nabla f_i(x) \approx \nabla f_i(x^*)$$

Fix  $\tilde{x}^0 = x$ , pick  $i$  randomly



$$\tilde{x}^{t+1} = \tilde{x}^t - \eta \underbrace{(\nabla f_i(\tilde{x}^t) - \nabla f_i(x) + \nabla f(x))}_{\text{variance reduction}} \quad \begin{matrix} \uparrow \\ \text{make sure} \\ \mathbb{E}[\tilde{L}] = \nabla f(\tilde{x}^t) \end{matrix}$$



# Lecture 15 Non-convex Optimization I Local Analysis

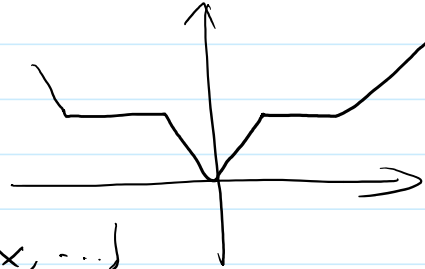
Sunday, October 23, 2016 10:59 PM

## - Non-convex optimization

### - what can a non-convex function look like?

- simpler case  
still has a unique  
minimum.

(quasi-convex, pseudo-convex, ...)



- complicated case  
multiple local optima..



## - optimality conditions

- first order optimality condition

$$\nabla f(x) = 0$$

- such points are called critical points.
- for (strongly) convex function,  $\nabla f(x) = 0 \Rightarrow x$  is optimal.

- second order condition

$$\nabla^2 f(x) \succeq 0$$

e.g.  $f(x) = x_1^2 + x_2^2$

$$\nabla^2 f(x) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \succeq 0.$$

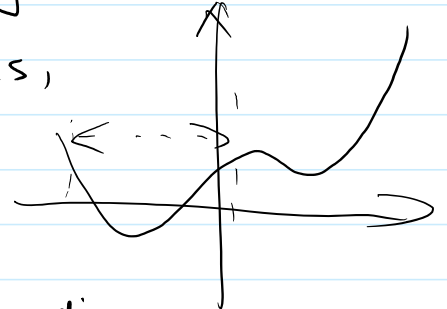
- saddle points

$\nabla^2 f(x)$  is not positive semidefinite or negative semidefinite.

e.g.  $f(x) = x_1^2 - x_2^2$   
 $\nabla^2 f(x) = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$

- Local convergence vs global convergence.

- when multiple local minima exists,  
 can hope to converge to global  
 minimum with good initialization.



- local structure of a non-convex function can  
 behave like a convex function!

- Approximate Gradient Descent.

- idea: after a good initialization, maybe the function  
 is very similar to convex.

- how to measure "similarity" to convex?

- Consider gradient descent, initial  $z^0$ , goal  $z^*$

$z^*$  is the optimum for convex function  $f(z)$

however, only has non-convex function  $g(z)$

hope:  $g(z)$  close to  $f(z)$

$$z^{t+1} = z^t - \eta \nabla g(z^t)$$

- Def:  $g$  is  $(\alpha, \beta, \varepsilon)$ -correlated if

$$\langle \nabla g(z^t), z^t - z^* \rangle \geq \alpha \|z^t - z^*\|^2 + \beta \|\nabla g(z^t)\|^2 - \varepsilon$$

note: if  $f$  is  $\mu$ -strongly convex and  $L$ -smooth

$$\langle \nabla f(z^t), z^t - z^* \rangle \geq \frac{\mu L}{\mu + 1} \|z^t - z^*\|^2 + \frac{1}{\mu + 1} \|\nabla f(z^t)\|^2$$

$$\langle \nabla f(\tilde{z}^t), \tilde{z}^t - \tilde{z}^* \rangle \geq \frac{\mu L}{\mu + L} \|\tilde{z}^t - \tilde{z}^*\|^2 + \frac{1}{\mu + L} \|\nabla g(\tilde{z}^t)\|^2$$

$(\frac{\mu L}{\mu + L}, \frac{1}{\mu + L}, 0)$  - correlated.

Theorem:  $\|\tilde{z}^{t+1} - \tilde{z}^*\|^2 \leq (1 - 2\alpha\eta) \|\tilde{z}^t - \tilde{z}^*\|^2 + 2\eta\varepsilon$ ,  
 (when  $\eta \leq 2\beta$ ), in particular

$$\|\tilde{z}^t - \tilde{z}^*\|^2 \leq (1 - 2\alpha\eta)^t \|\tilde{z}^0 - \tilde{z}^*\|^2 + \frac{\varepsilon}{2}$$

Proof:  $\|\tilde{z}^{t+1} - \tilde{z}^*\|^2 = \|\tilde{z}^t - \tilde{z}^*\|^2 - 2\eta \langle \nabla g(\tilde{z}^t), \tilde{z}^t - \tilde{z}^* \rangle + \eta^2 \|\nabla g(\tilde{z}^t)\|^2$

$$= \|\tilde{z}^t - \tilde{z}^*\|^2 - \eta (2 \langle \nabla g(\tilde{z}^t), \tilde{z}^t - \tilde{z}^* \rangle - \eta \|\nabla g(\tilde{z}^t)\|^2)$$

$$\leq \|\tilde{z}^t - \tilde{z}^*\|^2 - \eta (2\alpha \|\tilde{z}^t - \tilde{z}^*\|^2 + (2\beta - \eta) \|\nabla g(\tilde{z}^t)\|^2 - 2\varepsilon)$$

$$\leq (1 - 2\alpha\eta) \|\tilde{z}^t - \tilde{z}^*\|^2 + 2\eta\varepsilon$$