

## \* Lecture 6

7.4.3.1  
(in Book)

- Matrix Differentiation

Consider  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$

$$f(\vec{x}) = (f_1(\vec{x}), \dots, f_m(\vec{x})), \quad f_i: \mathbb{R}^n \rightarrow \mathbb{R}, \quad 1 \leq i \leq m$$

Define the Jacobian of  $f$ .

$$J(f) = \frac{\partial}{\partial \vec{x}} f(\vec{x}) \quad \text{A } m \times n \text{ matrix} \quad 1 \leq i \leq m, 1 \leq j \leq n$$

$$\textcircled{1} \quad (J(f))_{ij} = \left( \frac{\partial f_i(\vec{x})}{\partial x_j} \right)_{i,j=1..n}^{\text{m,n}} \in \mathbb{R}^{m \times n}$$

E.g. Consider  $f(x, y, z) = (x^2 + yz) \in \mathbb{R}^2$

map

$\mathbb{R}^3 \rightarrow \mathbb{R}^2$

$$J(f) = \frac{\partial}{\partial \vec{x}} f(\vec{x}) = \begin{pmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} & \frac{\partial f_1}{\partial z} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} & \frac{\partial f_2}{\partial z} \end{pmatrix} \quad f_1(\vec{x}) = x^2 + yz, \quad f_2(\vec{x}) = xyz$$

$$= \begin{pmatrix} 2x & z & y \\ yz & xz & xy \end{pmatrix}$$

Remark: If  $f: \mathbb{R}^m \rightarrow \mathbb{R}$ , then  $\frac{\partial f}{\partial \vec{x}}(\vec{x}) = \nabla_{\vec{x}} f(\vec{x})^T$ .  
this is the gradient.

E.g.  $f(\vec{x}) = x^2 + y^2 + z^2$

$$f: \mathbb{R}^3 \rightarrow \mathbb{R} \quad \frac{\partial f}{\partial \vec{x}}(\vec{x}) = (2x, 2y, 2z) = \nabla_{\vec{x}} f(\vec{x})^T$$

E.g.  $A \in \mathbb{R}^{m \times n}$  be  $\mathbb{R}^m$  and define  $f(\vec{x}) = A\vec{x} + \vec{b}$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\frac{\partial f}{\partial \vec{x}}(\vec{x}) = \left( \frac{\partial f_i}{\partial x_j}(\vec{x}) \right)_{i,j=1..n}^{\text{m,n}}$$

$$= (A_{i,j})_{i,j=1..n}^{\text{m,n}} = A$$

$$f_i(\vec{x}) = A_{i,j} \vec{x} + b_i$$

i-th row of  $A$  times  $\vec{x}$

$$= (\sum_{j=1}^n A_{i,j} x_j) + b_i$$

row on Col vector.

Corollary:  $\frac{\partial}{\partial \vec{x}} \vec{x} = I_{n \times n}, \quad \vec{x} \in \mathbb{R}^n$

$$f(\vec{x}) = \vec{x}$$

$$= I \vec{x}$$

Rules  $\textcircled{1} \frac{\partial}{\partial \vec{x}} (A\vec{x} + \vec{b}) = A \quad (A, \vec{b} \text{ constant})$

$\textcircled{2} \frac{\partial}{\partial \vec{x}} (\vec{x}^T A) = A^T \quad (A \text{ constant})$

$A$  is a matrix compatible with  $\vec{x}^T A$ .

$\textcircled{3} \frac{\partial}{\partial \vec{x}} (\vec{x}^T A \vec{x}) = \vec{x}^T (A + A^T) \quad (A \text{ constant})$

$f$  &  $g$  have  
compatible dimensions,  
so the products  
make sense.

$\textcircled{4} \frac{\partial}{\partial \vec{x}} f(\vec{x})^T g(\vec{x}) = g(\vec{x})^T \frac{\partial}{\partial \vec{x}} f(\vec{x}) + f(\vec{x})^T \frac{\partial}{\partial \vec{x}} g(\vec{x})$   
similar to some product rules.

$\textcircled{5} \frac{\partial}{\partial \vec{x}} f(\vec{x})^T A g(\vec{x}) = g(\vec{x})^T A^T \frac{\partial}{\partial \vec{x}} f(\vec{x}) + f(\vec{x})^T A \frac{\partial}{\partial \vec{x}} g(\vec{x})$

Proofs/more rules in reference "Matrix Calculus" on Canvas.

E.g.  $\nabla \vec{y} (-2\vec{x}^T D \vec{y} + \vec{y}^T \vec{y}) = 0$

(recall using this for P)

$$\Rightarrow \frac{\partial}{\partial \vec{y}} (-2\vec{x}^T D \vec{y} + \vec{y}^T \vec{y}) = 0 \quad (\text{from a past lecture})$$

can we rule  $\textcircled{1}$  here, to get.

linearity  $\Rightarrow -2 \frac{\partial}{\partial \vec{y}} (\vec{x}^T D \vec{y}) + \frac{\partial}{\partial \vec{y}} (\vec{y}^T \vec{y}) = 0$

use rule  $\textcircled{5} \Rightarrow -2 \left( \vec{y}^T D \frac{\partial}{\partial \vec{y}} \vec{x} + \vec{x}^T D \frac{\partial}{\partial \vec{y}} \vec{y} \right) + \frac{\partial}{\partial \vec{y}} (\vec{y}^T \vec{y}) = 0$

$$\textcircled{3} \Rightarrow -2 \vec{x}^T D I + \vec{y}^T (I + I) = 0$$

$$\Rightarrow -2 \vec{x}^T D = -2 \vec{y}^T \quad I^T = I$$

$$\Rightarrow \vec{y} = D^T \vec{x}$$

Notice:

\*  $\frac{\partial}{\partial \vec{x}} (af(\vec{x}) + bg(\vec{x})) = a \frac{\partial}{\partial \vec{x}} f(\vec{x}) + b \frac{\partial}{\partial \vec{x}} g(\vec{x})$   
more general form. (linearity)

The Hessian matrix of  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is

$$H(f)(\vec{x}) = J(\nabla_{\vec{x}} f(\vec{x})) = \frac{\partial}{\partial \vec{x}} \frac{\partial}{\partial \vec{x}} f(\vec{x})$$

$$= \frac{\partial^2}{\partial \vec{x}^2} f(\vec{x})$$

$$= \left( \frac{\partial^2}{\partial x_i \partial x_j} f(\vec{x}) \right)_{i,j=1..n}^{n,n}$$

①

Remark: At points  $\vec{x}$  s.t. all 2nd partials of  $f$  are continuous,

$$\frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} f(\vec{x}) = \frac{\partial}{\partial x_j} \frac{\partial}{\partial x_i} f(\vec{x})$$

so  $H(f)(\vec{x})$  is symmetric, so  $H(f)(\vec{x})$  has an orthonormal eigenvector decomposition and has real eigenvalues.

②  $\hat{u}$  unit vector

$\hat{u}^T H(f)(\vec{x}) \hat{u}$  is the 2nd derivatives of  $f(\vec{x})$  in the direction  $\hat{u}$ .

Let's go back to gradient descent

(Find  $\underset{\vec{x}}{\operatorname{argmin}} f(\vec{x})$ )

set  $\vec{g} = \nabla_{\vec{x}} f(\vec{x})$

$$\hat{u} = \frac{\vec{x} - \vec{x}_0}{\|\vec{x} - \vec{x}_0\|} \text{ unit vector}$$

$$H = H(f)(\vec{x}_0)$$

apply the Taylor expansion

$$f(\vec{x}) \leq f(\vec{x}_0) + (\vec{x} - \vec{x}_0)^T \vec{g} + \frac{1}{2} (\vec{x} - \vec{x}_0)^T H(\vec{x} - \vec{x}_0)$$

$\approx D_{\hat{u}} f(\vec{x})$

$\approx$  2nd derivative of  $f$  in direction  $\hat{u}$

$$\nabla_{\vec{u}} f(\vec{v}) = \vec{u} \cdot \nabla_{\vec{v}} f(\vec{v})$$

$\downarrow$  (3.1)  $\frac{\partial}{\partial v_i}$

$$\text{dot } \vec{u}^T$$

④ tells how close  
 $\vec{x}_0$  &  $\vec{g}$  are depending  
 on what  $\vec{z}$  is.  
 (these?)

### Lecture X

#### More on Gradient Descent.

keep going set  $\vec{x}_0 \in \mathbb{R}^n$ ,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

along the direction of the gradient,

unless we reach a stopping position.

then  $f(\vec{x}) \leq f(\vec{x}_0) + (\vec{x} - \vec{x}_0)^T \vec{g} + \frac{1}{2} (\vec{x} - \vec{x}_0)^T H (\vec{x} - \vec{x}_0)$

Gradient descent step

$$\textcircled{1} \quad f(\vec{x}_0 - \alpha \vec{g}) \leq f(\vec{x}_0) - \alpha \vec{g}^T \vec{g} + \frac{1}{2} \alpha^2 \vec{g}^T H \vec{g}$$

$\alpha$  is the stopping condition

case 1  $\vec{g}^T H \vec{g} > 0$  (guaranteed if  $H$  is positive definite).

we are trying to minimize  $f$  to get the largest possible drop

compute the optimal step size:

$$\frac{\partial}{\partial \alpha} f(\vec{x}_0 - \alpha \vec{g}) = -\vec{g}^T \vec{g} + \alpha \vec{g}^T H \vec{g} \Rightarrow$$

Hessian?

$$\alpha = \frac{\vec{g}^T \vec{g}}{\vec{g}^T H \vec{g}}$$

case 2  $\vec{g}^T H \vec{g} < 0$  (Previous case gives a local max... no good)

- set sth. wrong
- the problem is now meaningless.

when  $\tilde{g}^T H \tilde{g}$ . Taylor series says increasing & forever decreases  $f$   
 (probably not true ... Taylor series is only accurate locally)

probably looking @  $\nabla_x^2 f(\tilde{x}_0) = \tilde{0}$   
 sth. that's a saddle

- ① If  $H(f)(\tilde{x}_0)$  is positive definite, then  $\tilde{x}_0$  is a local min.
- ② If  $\dots$  negative  $\rightarrow$  local max.
- ③ If  $H(f)(\tilde{x}_0)$  has both positive and negative eigenvalues,  
 $\tilde{x}_0$  is a saddle point.
- (④.  $H(f)(\tilde{x}_0) = 0 (?)$  inconclusive, same as not saying anything ...)

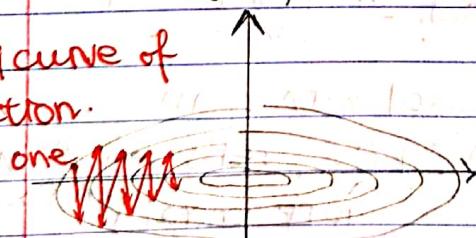
Remark: The condition number  $\max_{i,j} \frac{\lambda_i}{\lambda_j}$  of  $H$  indicates how well a gradient descent step will perform.

suppose  $m \geq 2$ ,  
 $H(f) = \begin{pmatrix} 10 & 0 \\ 0 & 1 \end{pmatrix}$

large  $\rightarrow$  bad  
 small (1)  $\rightarrow$  good.

the level curve of the function.

Idea: in one direction, stretch very far.



- horrible path,
- goes "forever"
- terrible performance when step size is fixed & large.

Fixes:

- ① Use Hessian info (algorithm centric)

② normalize dataset (data centric)

take a look @ mean & stdv for each variable; then perform some transformation that stdv of each is similar.

(more normalized :))

Q: when can we guarantee that gradient descent (or some variation) will converge?

Answer: In deep learning, almost never, BUT for some the function is so complicated that it never has the best properties (e.g. neural network).

### 3 Principal Component Analysis

special functions, we may get guarantees.

Defn: A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is  $L$ -Lipschitz continuous.

If  $\exists L > 0$ , s.t.  $\forall \vec{x}, \vec{y} \in \mathbb{R}^n$ ,

$$\|f(\vec{x}) - f(\vec{y})\|_2 \leq L \|\vec{x} - \vec{y}\|_2.$$

a positive constant  $L$  (bounds the possible change in function value based on size of change in inputs.)

E.g.  $f(\vec{x}) = A\vec{x} + \vec{b}$

$$f(\vec{x}) - f(\vec{y}) = A\vec{x} + \vec{b} - A\vec{y} - \vec{b} = A(\vec{x} - \vec{y})$$

$$\therefore \max_{\text{here:}} \|f(\vec{x}) - f(\vec{y})\|_2 = \|A(\vec{x} - \vec{y})\|_2$$

$$\|A\|_2 = \max_{\vec{x} \neq \vec{0}} \frac{\|A\vec{x}\|_2}{\|\vec{x}\|_2} \leq \|A\|_2 \|\vec{x} - \vec{y}\|_2$$

so  $\forall \vec{x} \in \mathbb{R}^n$

$$\|A\|_2 \|\vec{x}\|_2 \geq \|A\vec{x}\|_2$$

so  $f$  is  $\|A\|_2$ -Lipschitz continuous.

this is the minimum value of  $L$  as well.

Defn:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if  $H(f)(\vec{x})$  is positive semi-definite for all  $\vec{x} \in \mathbb{R}^n$ .

" " strict convexity, " positive definite

" " "

extra ensured such as unique local minimum.

E.g. Consider  $f(\vec{x}) = \vec{b}^T \vec{x} + \vec{x}^T A \vec{x} + c$ . Let's see if this is

$$\frac{\partial f}{\partial \vec{x}}(\vec{x}) = \vec{b}^T + \vec{x}^T (A + A^T)$$

convex! to be general:

$A$  is a symmetric matrix.

$$= \vec{b}^T + 2\vec{x}^T A$$

$$\frac{\partial^2 f}{\partial \vec{x}^2}(\vec{x}) = 2A^T = 2A.$$

the rule says

so  $f$  is convex iff  $A$  is positive semi-definite.

E.g.  $f(\vec{x}) = \vec{x}^T \vec{x} = \|\vec{x}\|_2^2$  is convex.

$A = I$  here, so positive definite  
.. strictly convex.

can be Lipschitz CTS  
with many diff values of  $L$   
But the most in up is the smallest

usually.  
knowing the function  
satisfy some  
properties will follow

## \*Lecture 8

### Theorem //

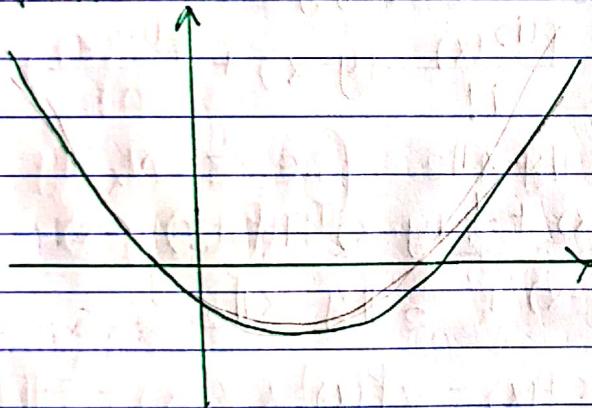
Remark: The definition I gave for convexity can be made more general.

Defn:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if  $\forall \vec{x}, \vec{y} \in \mathbb{R}^n \quad \forall 0 \leq t \leq 1$

$$f(t\vec{x} + (1-t)\vec{y}) \leq t f(\vec{x}) + (1-t) f(\vec{y})$$

Theorem: If  $f$  is twice differentiable on  $\mathbb{R}^n$  then  $f$  is convex  $\Leftrightarrow H(f)(\vec{x})$  positive semi-definite.  $\forall \vec{x} \in \mathbb{R}^n$

Eg.  $f: \mathbb{R} \rightarrow \mathbb{R}$



Q: How to see general defn here?

Theorem: suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is twice differentiable convex, and has  $L$ -Lipschitz continuous gradient, and  $f$  has a global min  $x^* \in \mathbb{R}^n$ . Then if  $x_k$  is the  $k^{th}$  step of gradient descent, with fixed step size  $\alpha < 1/L$ , then

$$f(x_k) - f(x^*) \leq \frac{\|x_0 - x^*\|^2}{2L}$$

proof:  $\nabla f$  L-Lipschitz  $\Rightarrow H(f)(\vec{x}) - L\mathbf{I}$  negative semi definite

proof:  $f_{\vec{x}, \vec{u}}(t) = f(\vec{x} + t\vec{u}) \quad t \in \mathbb{R}$

$$\text{then } f'_{\vec{x}, \vec{u}}(0) = \vec{u}^T \nabla f(\vec{x})$$

$$\vec{u}^T (H(f)(\vec{x}) - L\mathbf{I}) \vec{u} = \vec{u}^T H(f)(\vec{x}) \vec{u} - \vec{u}^T L\mathbf{I} \vec{u}$$

$$= \vec{u}^T \vec{u} - L \vec{u}^T \vec{u}$$

$$= f''_{\vec{x}, \vec{u}}(0) - L = \lim_{h \rightarrow 0} \frac{f'_{\vec{x}, \vec{u}}(h) - f'_{\vec{x}, \vec{u}}(0)}{h} - L$$

$$\begin{aligned}
 &= \lim_{h \rightarrow 0} \frac{\hat{u} \cdot [\nabla_{\vec{x}} f(\vec{x} + h\hat{u}) - \nabla_{\vec{x}} f(\vec{x})]}{h} - L \\
 &\leq \lim_{h \rightarrow 0} \frac{\|\nabla_{\vec{x}} f(\vec{x} + h\hat{u}) - \nabla_{\vec{x}} f(\vec{x})\|_2}{h} - L \\
 &\leq \lim_{h \rightarrow 0} \frac{L \|h\hat{u}\|_2}{h} - L = 0 \quad \frac{\partial L}{\partial \hat{u}}
 \end{aligned}$$

$$\therefore \forall x, y, z \in \mathbb{R}^n, (x-y)^T (H(f)(z) - L) (x-y) \leq 0$$

Let  $x, y \in \mathbb{R}^n$ . By Taylor's Remainder Theorem,

(in  $\mathbb{R}$ , TRT says  $\exists z \in [x, y]$ )

$$f(y) = \sum_{i=0}^{n-1} \frac{f^{(i)}(x)}{i!} (y-x) + \frac{f^{(n)}(z)}{n!} (y-x)$$

$$\begin{aligned}
 &\exists z \text{ s.t. } \|z - x\|_2 \leq \|y - x\|_2 \quad (\text{i.e. } z \in N_x(\|y - x\|_2)) \\
 f(y) &= f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T H(f)(z) (y-x) \\
 &\leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y - x\|_2
 \end{aligned}$$

$$\begin{aligned}
 \text{Now, } f(x - \alpha \nabla f(x)) &\leq f(x) - \nabla f(x)^T \alpha \nabla f(x) + \frac{L}{2} \|\alpha \nabla f(x)\|_2^2 \\
 &= f(x) - \left(1 - \frac{\alpha L}{2}\right) \alpha \|\nabla f(x)\|_2^2
 \end{aligned}$$

$$0 < \alpha < \frac{1}{L}$$

$$\text{Since } 0 < \alpha < \frac{1}{L} \iff 0 < \frac{\alpha L}{2} < \frac{1}{2}$$

$$\iff 0 > -\frac{\alpha L}{2} > -\frac{1}{2}$$

$$\iff 1 > \boxed{1 - \frac{\alpha L}{2} > 1 - \frac{1}{2}}$$

$$\frac{1}{2}$$

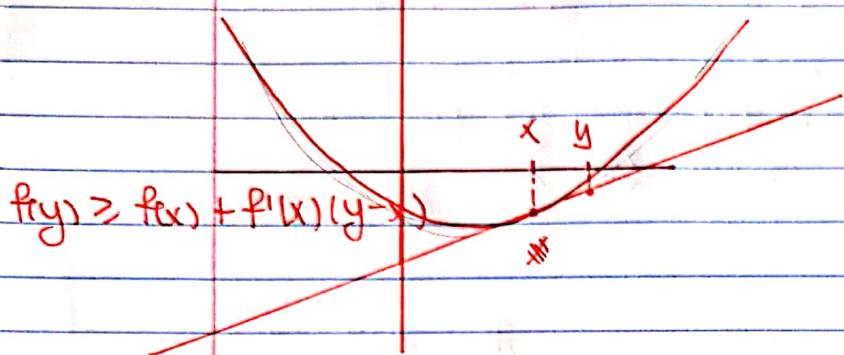
$$f(x - \alpha \nabla f(x)) \leq f(x) - \frac{\alpha}{2} \|\nabla f(x)\|_2^2 \quad \begin{cases} \Rightarrow \text{sequence} \\ (f(x_i))_{i \geq 0} \end{cases}$$

weakly decreasing

$$f \text{ convex} \Rightarrow f(x^*) \geq f(x) + \nabla_x f(x)^T (x^* - x)$$

note:  $f: \mathbb{R} \rightarrow \mathbb{R}$

$$\Rightarrow f(x) \leq f(x^*) + \nabla_x f(x)^T (x - x^*)$$



$$f(y) \geq f(x) + \nabla_x f(x)(y - x)$$

$$\begin{aligned} f(x_{i+1}) &= f(x_i - \alpha \nabla_x f(x_i)) \leq f(x_i) - \frac{\alpha}{2} \|\nabla_x f(x_i)\|_2^2 \\ &\leq f(x^*) + \nabla_x f(x_i)^T (x - x^*) \\ &\quad - \frac{\alpha}{2} \|\nabla_x f(x_i)\|_2^2 \end{aligned}$$

$$x_{i+1} = x_i - \alpha \nabla_x f(x_i) \Rightarrow \nabla_x f(x_i) = x_i - x_{i+1}$$

$$\begin{aligned} f(x_{i+1}) &\leq f(x^*) + \nabla_x f(x_i)^T (x_i - x^*) - \frac{\alpha}{2} \|\nabla_x f(x_i)\|_2^2 \\ &= f(x^*) + \frac{1}{2} (x_i - x_{i+1})^T (x_i - x^*) - \frac{\|x_i - x_{i+1}\|_2^2}{2\alpha} \\ &= f(x^*) + \frac{1}{2} (x_i - x_{i+1})^T (x_i - x^*) - \frac{1}{2\alpha} (x_i - x_{i+1})^T (x_i - x_{i+1}) \end{aligned}$$

complete the square

$$\begin{aligned} &= f(x^*) - \frac{1}{2\alpha} [(x_i - x_{i+1})^T (x_i - x_{i+1}) - 2(x_i - x_{i+1})^T (x_i - x^*) \\ &\quad + (x_i - x^*)^T (x_i - x^*) - (x_i - x^*)^T (x_i - x^*)] \end{aligned}$$

$$= f(x^*) - \frac{1}{2\alpha} [((x_i - x_{i+1}) - (x_i - x^*))^T ((x_i - x_{i+1}) - (x_i - x^*)) - 1]$$

$$= f(x^*) - \frac{1}{2\alpha} [\|x^* - x_{i+1}\|_2^2 - \|x_i - x^*\|_2^2]$$

$$\sum_{i=0}^k (f(x_i) - f(x^*)) = \sum_{i=0}^k \frac{1}{2\alpha} [\dots]$$

$$= \frac{1}{2\alpha} [\|x^* - x_0\|_2^2 - \|x^* - x_k\|_2^2]$$

$$\leq \frac{1}{2\alpha} \|x^* - x_0\|_2^2$$

$(f(x_i))$  decreasing  $\Rightarrow f(x_k) - f(x^*) \leq f(x_i) - f(x^*) \quad \forall i < k$

$$\begin{aligned}\therefore k(f(x_k) - f(x^*)) &= \sum_{i=1}^k (f(x_k) - f(x^*)) \\ &\leq \sum_{i=1}^k (f(x_i) - f(x^*)) \\ &\leq \frac{1}{2\alpha} \|x_0 - x^*\|_2^2\end{aligned}$$