

Note Title

9/28/2017

$$\uparrow$$
  

$$- \quad \underline{\underline{D}} \quad \Rightarrow \quad \underset{\text{residual}}{e} = y - w^T x \quad w \in \mathbb{R}^D$$

$$\epsilon \sim N(0, \sigma^2)$$

A System  $L$  is linear if

$$L(ax + by) = aL(x) + bL(y)$$

example

$$L = A_{n \times n}$$

$$A_{n \times n} [aX + y] = aAX + Ay$$

right now  $X_1$  are  
feature

$$\begin{bmatrix} x_1 & x_2 & \dots & x_D & x_1^2 & x_2^2 & \dots & x_D^2 \end{bmatrix} = Y$$

$w_1 \quad w_2 \quad \dots \quad w_D \quad \quad \quad w_{2D} = W$

In general

$$Y = W^T X + \epsilon$$

^ We can modify features  $X \xrightarrow{\phi} \phi(X)$

$\phi \rightarrow$  polynomial

$\rightarrow$  kernel methods

$\rightarrow$  deep learning  
does feature engineering

$D \rightarrow 2D$

In general

We can model  
non-linear system  
as

$$N(Y | \phi(X), \sigma^2)$$

$$\rightarrow X^T X_2 \in \mathbb{R}$$

$$X_i \in \mathbb{R}^D$$

later we will see that we  
only need to know about interactions  
between feature  $X_i$

MLE estimation (Least Square Estimation)

let say we observed IID sample

$$D = \{(x_i, y_i)\}_{i=1}^N$$

$$P(D|\theta)$$

$$\hat{\theta} = (\hat{w}) = \arg \max_{\theta} \log P(D|\theta)$$

$$P(\theta|D)$$

$$= \arg \max_{\theta} \log \prod_{i=1}^N \mathcal{N}(y_i | w^T x_i, \sigma^2)$$

$$= \arg \min_{\theta} - \log \underbrace{\prod_{i=1}^N \frac{1}{(2\pi)^{1/2} \sigma} \exp\left(-\frac{1}{2} \frac{(y_i - w^T x_i)^2}{\sigma^2}\right)}_{LLL(\theta)}$$

(because lot of machine learning packages has minimization abstraction (API) built-in)

$$= \arg \min_{\theta} \left[ - \sum_{i=1}^N \log \frac{1}{(2\pi)^{1/2} \sigma} - \sum_{i=1}^N \frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right]$$

$$= \arg \min_{\theta} \sum_{i=1}^N \frac{1}{2\sigma^2} (y_i - w^T x_i)^2$$

$$= \arg \min_{w \in \mathbb{R}^D} \sum_{i=1}^N (y_i - w^T x_i)^2$$

↓  
RSS - Residual sum of square

$$= \arg \min_{w \in \mathbb{R}^D} \sum_{i=1}^N \|e_i\|^2$$

$$= \arg \min_{w \in \mathbb{R}^D} \|e\|^2$$

MSE (mean square error)

$$= \frac{RSS}{N}$$

$$e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix}$$

$$= (y - Xw)^T (y - Xw)$$

$$X = \begin{bmatrix} - & x_1^T & - \\ - & x_2^T & - \\ & \vdots & \\ - & x_i^T & - \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix} \in \mathbb{R}^D$$

$$y - Xw$$

$$= \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

side

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

$$C_{11} = A_{11}B_{11} + A_{12}B_{21}$$

For matrices block wise multiplication is also true.

$$= \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} x_1^T w_1 = w_1^T x_1 \\ x_2^T w_2 \\ \vdots \\ x_N^T w_N \end{bmatrix}$$

$$\begin{pmatrix} (x_1^T w_1) \\ \vdots \\ (x_N^T w_N) \end{pmatrix}_{1 \times 1}$$

$$= \begin{bmatrix} y_1 - x_1^T w_1 \\ y_2 - x_2^T w_2 \\ \vdots \\ y_N - x_N^T w_N \end{bmatrix}$$

$$\cancel{x}^T \cancel{x} = \|x\|_2^2$$

$$= x_1^2 + x_2^2 + x_3^2$$

$$\begin{aligned}
 & (y - Xw)^T (y - Xw) = \ell(w) \\
 & = \begin{bmatrix} y_1 - x_1^T w \\ y_2 - x_2^T w \\ \vdots \\ y_N - x_N^T w \end{bmatrix}^T \begin{bmatrix} y_1 - x_1^T w \\ y_2 - x_2^T w \\ \vdots \\ y_N - x_N^T w \end{bmatrix} \\
 & \underset{\text{argmin}}{w} = \sum_{i=1}^N (y_i - x_i^T w)^2 = \ell(w)
 \end{aligned}$$

$$\ell(w) = (y - Xw)^T (y - Xw)$$

$$= (y^T - (Xw)^T) (y - Xw)$$

$$= (y^T - w^T X^T) (y - Xw)$$

$$\begin{aligned}
 & = y^T y - y^T \underbrace{Xw}_{= -2 y^T X w} - w^T X^T y + w^T X^T X w
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial \mathcal{L}(w)}{\partial w} &= 0 - 2y^T X + w^T (X^T X + (X^T X)^T) \\
 &= -2y^T X + w^T (X^T X + X^T (X^T)^T) \\
 &= -2y^T X + 2w^T X^T X
 \end{aligned}$$

let put it equal to 0

$$2w^T X^T X = +2y^T X$$

$$X^T X w = X^T y$$

$$w = (X^T X)^{-1} X^T y$$

$X_{N \times D}$  as far as column rank of  $X$  is  $D$

$$\frac{d\ell(w)}{dw} = \cancel{y^T X} + \underline{2w^T X^T X}$$

$$\frac{d\ell(w)}{dw^2} = \cancel{\text{scribble}} + 2 X^T X_{D \times D}$$

so  $w$  is minimizes if

$X^T X$  is positive

definite matrix ~~(PS)~~

$\equiv$   
positive definite matrix  $A_{m \times m}$

$\equiv$   $A$  is a positive definite matrix  
if for any non-zero vector  $v$



$$\begin{pmatrix} (V^T)_{1 \times m} A_{m \times m} V_{m \times 1} \end{pmatrix}_{1 \times 1} > 0$$

is  $X^T X$  positive definite.

$$\begin{aligned} & V^T X^T X V \\ &= (XV)^T XV = \|XV\|_2^2 > 0 \end{aligned}$$

side

$$X \in \mathbb{R}^n$$

$$A_{m \times n} X_{n \times 1}$$

$$(AX)_{m \times 1} \in \mathbb{R}^m$$

$$V = 0 \longrightarrow XV = 0$$

If  $X$  is full rank, then it is injective (1-1). Hence

$$XV \neq 0 \text{ for } V \neq 0 \\ \Rightarrow \|XV\|_2 \neq 0$$

Convex Function

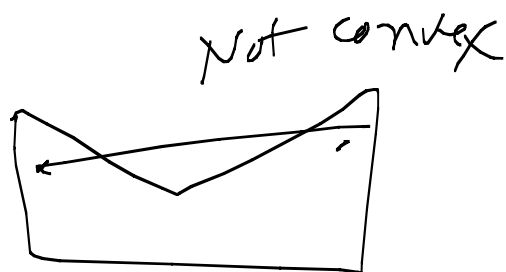
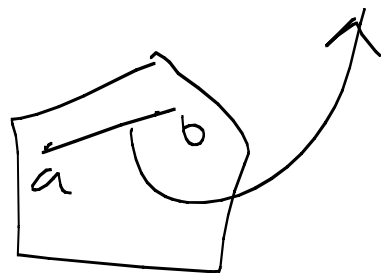
a function  $f(\theta)$

is called Convex \*

Convex Set  $A$  = A Set is Convex

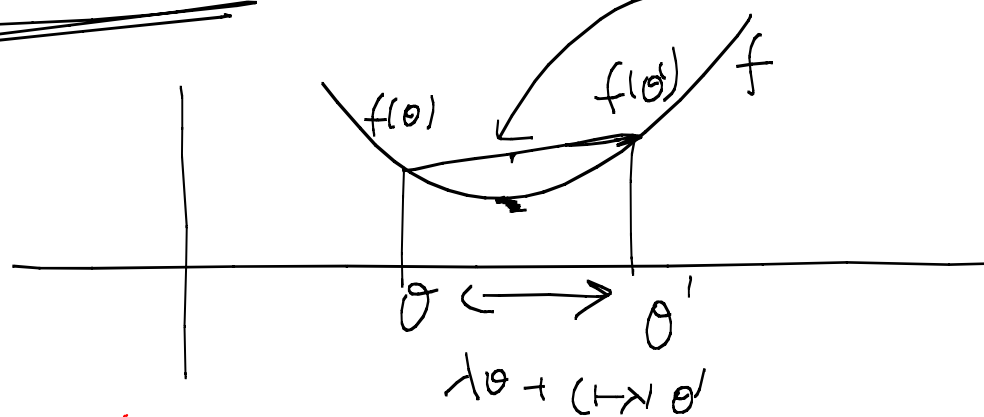
if for point  $a \in A, b \in A$

$$\lambda a + (1-\lambda)b \in A \quad \lambda \in [0, 1]$$



on a convex set  $S$  it for any  
 $\theta, \theta' \in S$

$$\underline{f(\lambda\theta + (1-\lambda)\theta')} \leq \underline{\lambda f(\theta) + (1-\lambda)f(\theta')}$$



i.e line joining  $f(\theta)$  and  $f(\theta')$  is always  
 above function between  $\theta$  and  $\theta'$ ,

