

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

$$\Rightarrow \epsilon = y - \hat{w}^T x \quad w \in \mathbb{R}^D$$

$$\mathbb{E} \sim N(0, \sigma^2)$$

A system L is linear if

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

∇

$$A_{\max} [ax + y]$$

$$A^T X + A y$$

$$\frac{\begin{bmatrix} x_1 & x_2 & \dots & x_D & x_1^2 & x_2^2 & \dots & x_D^2 \end{bmatrix}}{w_1 \ w_2 \ \dots \ w_D \quad w_{2D} = 1} \cdot \begin{bmatrix} x_1^2 \\ x_2^2 \\ \vdots \\ x_D^2 \end{bmatrix} = y$$

In general

$$y = w^T \phi + \epsilon$$

^ We can modify features $x \rightarrow \phi(x)$

$$D \rightarrow 2D$$

In general

We can model

non-linear system.

$$N(y | \phi(x), \sigma^2)$$

$\phi \rightarrow$ polynomial
 \rightarrow kernel methods
 \rightarrow deep learning
 does feature engineering

$\rightarrow (X_1^T \ X_2^T \ \dots \ X_D^T) \in \mathbb{R}^{1 \times D}$
 $X_i \in \mathbb{R}^D$
 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 $D = \{(x_i, y_i)\}_{i=1}^N$

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

$$P(\theta|D)$$

$$= \arg \max_{\theta} \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (y_i - w^T x_i)^2\right)$$

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

LLL(θ)

(because lot of machine learning

packages has minimization abstraction (API) built-in

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

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

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

$$= \arg \min_{w \in \mathbb{R}^D} \sum_{i=1}^N \|e_i\|^2 \quad \text{RSS - Residual Sum of Square}$$

$$= \arg \min_{w \in \mathbb{R}^D} \sum_{i=1}^N \|e_i\|^2 \quad \text{MSE (mean square error)}$$

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

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

$$= (y - Xw)^T (y - Xw) \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

$$X = \begin{bmatrix} - & - & x_1^T & - \\ - & - & x_2^T & - \\ - & - & x_i^T & - \end{bmatrix} \quad w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{pmatrix} \in \mathbb{R}^D$$

$$y - Xw$$

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

Side

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

For matrices block wise multiplication is also true.

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

$$\left(\begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} w \right)_{1 \times 1}$$

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

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

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

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

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

let put it equal to 0

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

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

$$w = \underbrace{(X^T X)^{-1}}_{\text{it's as far as column rank of } X \text{ is } D} X^T y$$

$$\frac{d\ell(w)}{dw} = \frac{\cancel{2w} - 2y^T x + 2w^T x^T x}{1}$$

$$\frac{d\ell(w)}{dw^2} = \cancel{0} + 2x^T x_{D \times D}$$

So w is minimizes ℓ

$x^T x$ is positive

definite matrix ~~(~~is~~)~~

positive definite matrix $A_{n \times n}$
 \Rightarrow A is a positive definite matrix
 \Rightarrow it for any non-zero vector v

$$\begin{pmatrix} (V^T A V)_{m \times m} & 0 \\ 0 & 0 \end{pmatrix}$$

is $X^T X$ positive definite.

$$V^T X^T X V$$

$$= (XV)^T XV = \|XV\|_2^2 > 0 \quad ?$$

side

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

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

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

$$V = 0$$



$$X V = 0$$

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

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

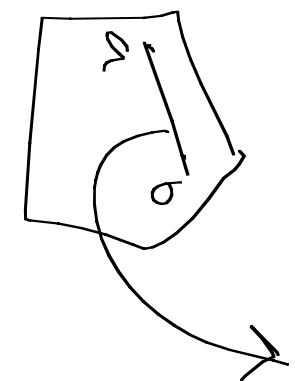
Convex Function a function $f(\theta)$

is called Convex *

Convex Set $A = A$ Set is convex

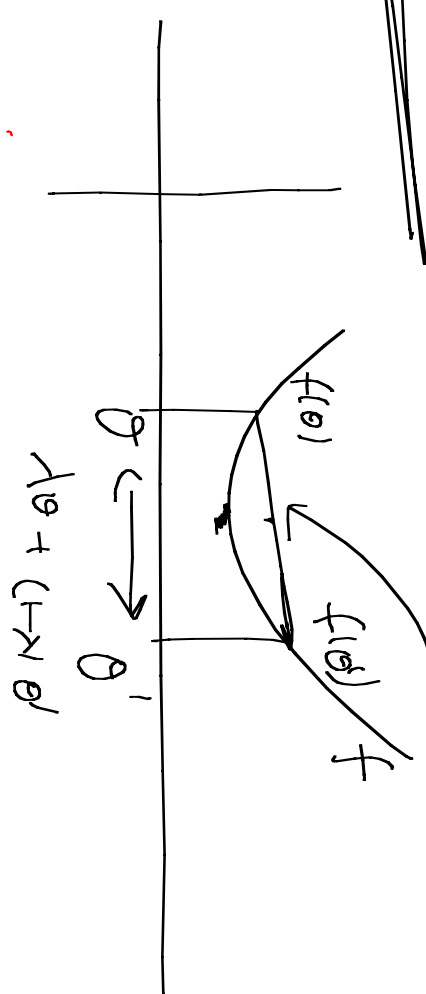
if for point $a \in A, b \in A$

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



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

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



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

