

2/18/15

①

REGRESSION

① Recall: Linear Regression

$$\text{Model } \hat{y} = w_0 + w_1 x_1 + \dots + w_d x_d$$

$$= \vec{w}^T \vec{x}$$

Least Squares Fitting / Ordinary Least Squares
 dataset: $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

$$\hat{w}_{OLS} = \underset{\vec{w}}{\text{argmin}} L(\vec{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \frac{1}{n} \|Y - Xw\|_2^2$$

$$= \frac{1}{n} [Y^T Y - 2Y^T Xw + w^T X^T Xw]$$

→ Detour: Matrix / Vector differentiation

Scalars: x, y

Vectors: \vec{x}, \vec{y}

Matrices: X, Y

	S	V	M
S	$\frac{\partial y}{\partial x}$	$\frac{\partial \vec{y}}{\partial \vec{x}}$	$\frac{\partial Y}{\partial X}$
V	$\frac{\partial y}{\partial \vec{x}}$	$\frac{\partial \vec{y}}{\partial \vec{x}}$	
M	$\frac{\partial y}{\partial X}$		

Convention: $\frac{\partial \vec{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \vdots \\ \frac{\partial y_m}{\partial x} \end{bmatrix}$ ↓ numerator = dim 1 / col-vector

$\frac{\partial y}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix}$ ← denominator = dim 2 / row-vector

$\frac{\partial \vec{y}}{\partial \vec{x}} = \begin{bmatrix} \ddots & \vdots & \ddots \\ \dots & \frac{\partial y_i}{\partial x_j} & \dots \\ \vdots & \vdots & \vdots \end{bmatrix}_{m \times n}$

Now,

$\frac{\partial (\vec{w}^T \vec{x})}{\partial \vec{w}} = \begin{bmatrix} \frac{\partial \vec{w}^T \vec{x}}{\partial w_0} & \dots & \frac{\partial \vec{w}^T \vec{x}}{\partial w_d} \end{bmatrix}$

$\frac{\partial \left[\sum_{i=0}^d w_i x_i \right]}{\partial w_0} = x_0 \quad (= 1 \text{ in this case})$

Thus $\frac{\partial (\vec{w}^T \vec{x})}{\partial \vec{w}} = [x_0 \ x_1 \ \dots \ x_d] = \vec{x}^T$
(very intuitive)

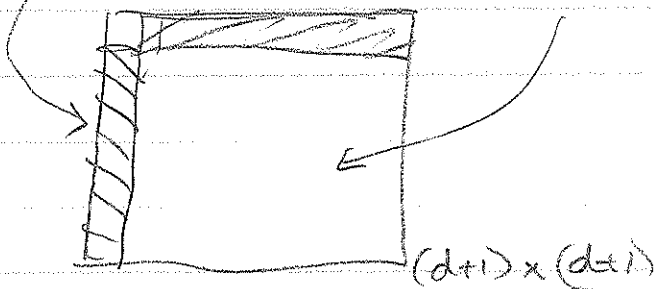
(2)

Also, $\frac{\partial (w^T A w)}{\partial \vec{w}} = \left[\frac{\partial (w^T A w)}{\partial w_0} \quad \dots \quad \frac{\partial (w^T A w)}{\partial w_d} \right]$

$$\frac{\partial [w^T A w]}{\partial w_0} = \frac{\partial \left[\sum_{i=0}^d \sum_{j=0}^d w_i a_{ij} w_j \right]}{\partial w_0}$$

$$= \frac{\partial \left[a_{000} w_0^2 + \sum_{i \neq 0} a_{ij} w_i w_0 + \sum_{j \neq 0} a_{ij} w_0 w_j \right]}{\partial w_0}$$

$$+ \sum_{i \neq 0} \sum_{j \neq 0} w_i a_{ij} w_j$$



$$= 2 a_{000} w_0 + \sum_{i \neq 0} a_{ij} w_i + \sum_{j \neq 0} a_{ij} w_j$$

$$= 2 a_{000} w_0 + 2 \sum_{i \neq 0} a_{ij} w_i$$

$$= 2 w^T A_0 \leftarrow \text{first column of } A$$

$$\Rightarrow \frac{\partial (w^T A w)}{\partial \vec{w}} = [2 w^T A_0 \quad 2 w^T A_1 \quad \dots \quad 2 w^T A_d]$$

$$= 2 w^T A \quad (\text{also intuitive!})$$

Just drop \vec{w}

← END DETOUR →

Putting it all together

$$\frac{\partial L(\vec{w})}{\partial \vec{w}} = \frac{1}{n} \left[-2 Y^T X + 2 W^T X^T X \right] = 0$$

$$\Rightarrow W^T X^T X = Y^T X$$

$$\Rightarrow (X^T X) w = X^T Y$$

$$\Rightarrow \boxed{\hat{w}_{OLS} = (X^T X)^{-1} X^T Y}$$

Single clean equation

Note 1: $(X^T X)^{-1} X^T$ is called the Moore-Penrose Pseudo Inverse of X

notation $X^+ = (X^T X)^{-1} X^T$

Why? behaves like an inverse

$$X^+ X = (X^T X)^{-1} X^T X = I$$

Recall $Y = Xw \Rightarrow w = X^{-1} Y$

$Y \approx \hat{Y} = Xw \Rightarrow w = X^+ \hat{Y}$ } see the resemblance?

squared loss

(3)

Note 2: $(X^T X)$ may not be invertible

Specifically,

$$(X^T X)_{(d+1) \times (d+1)} = \left[\begin{array}{c} \uparrow \\ \vec{x}_i \\ \downarrow \end{array} \right]_{(d+1) \times n} \left[\begin{array}{c} \leftarrow \vec{x}_i \rightarrow \\ \hline \end{array} \right]_{n \times (d+1)}$$

$$= \sum_{i=1}^n \underbrace{\vec{x}_i \vec{x}_i^T}_{\text{rank 1 matrix}}$$

n-summations

\Rightarrow if $n < (d+1)$

$\sum \vec{x}_i \vec{x}_i^T$ has rank $< (d+1)$

$\Rightarrow X^T X$ not invertible.

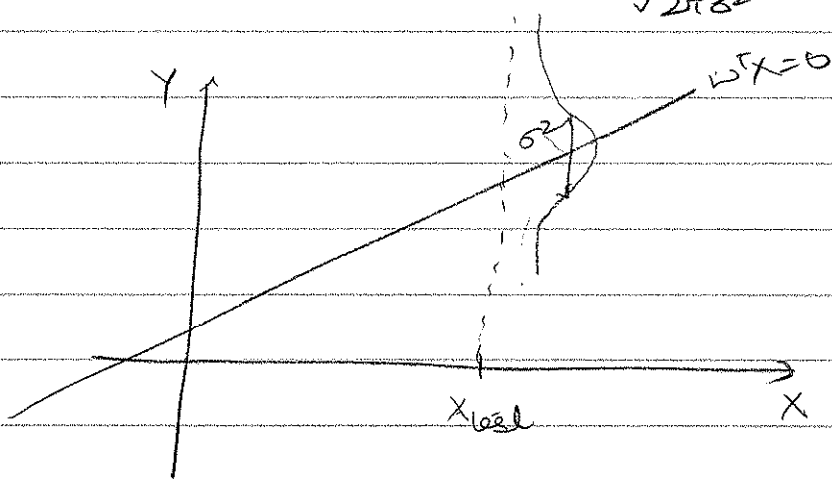
② Least Squares as MLE

Consider a probabilistic model

$$X \sim P(X) \quad [\text{Don't care}]$$

$$Y|X=x \sim N(W^T x, \sigma^2) \quad \leftarrow \text{fixed / don't care}$$

$$\Rightarrow p(y|X=x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-W^T x)^2}{2\sigma^2}}$$



$$e = y - W^T x \sim N(0, \sigma^2) \Leftrightarrow p(e) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-W^T x)^2}{2\sigma^2}}$$

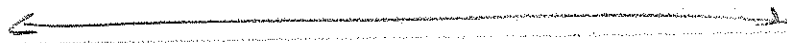
Now, given dataset $D = \{(\vec{x}_1, y_1) \dots (\vec{x}_n, y_n)\}$

$$\hat{w}_{MLE} = \underset{\vec{w}}{\text{argmax}} \log P(D | \vec{w})$$

$$= \underset{\vec{w}}{\text{argmax}} \sum_{i=1}^n \log P(e_i | w)$$

$$= \underset{\vec{w}}{\text{argmax}} \sum_{i=1}^n \left[-\frac{(y_i - W^T x_i)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right]$$

$$\begin{aligned} \Rightarrow \hat{w}_{MLE} &= \underset{w}{\operatorname{argmax}} \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T x_i)^2 \\ &= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (y_i - w^T x_i)^2 = \hat{w}_{OLS} \end{aligned}$$

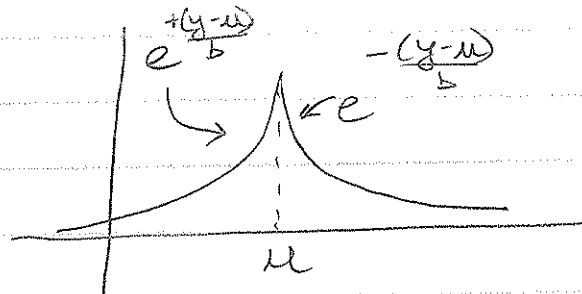


③ Robust Least Squares

$Y | X=x \sim \text{Laplacian}(w^T x, b)$ ↙ again fixed

Recall,

$$Y \sim \text{Lap}(\mu, b) \Leftrightarrow p(y|\mu, b) = \frac{1}{2b} e^{-\frac{|y-\mu|}{b}}$$



$$\hat{w}_{MLE} = \underset{w}{\operatorname{argmax}} \log P(D|w)$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \left[-\frac{|y_i - w^T x_i|}{b} - \log 2b \right]$$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n |y_i - w^T x_i|$$

So Laplace likelihood

↔ Least Absolute Error!

How do we solve this?

HW2!