

Linear Regression

$$x \longrightarrow y$$

x is a vector

$x \in \mathbb{R}^d \rightarrow x$ is a vector of length d

y is a scalar

$$y \in \mathbb{R}$$

Prediction or Regression

Predict future income

Current GPA, # AI courses taken	$y =$	7000
	$y =$	<u>300000</u>

3.8	4
-----	---

$y ?$

Training data

GPA, # AI	Income
-	100000

3.7, 2	1000
3.9, 6	500000
2.4, 1	1,000,000

←

3.5, 2	?
--------	---

Functional models.

$$y = f(x)$$

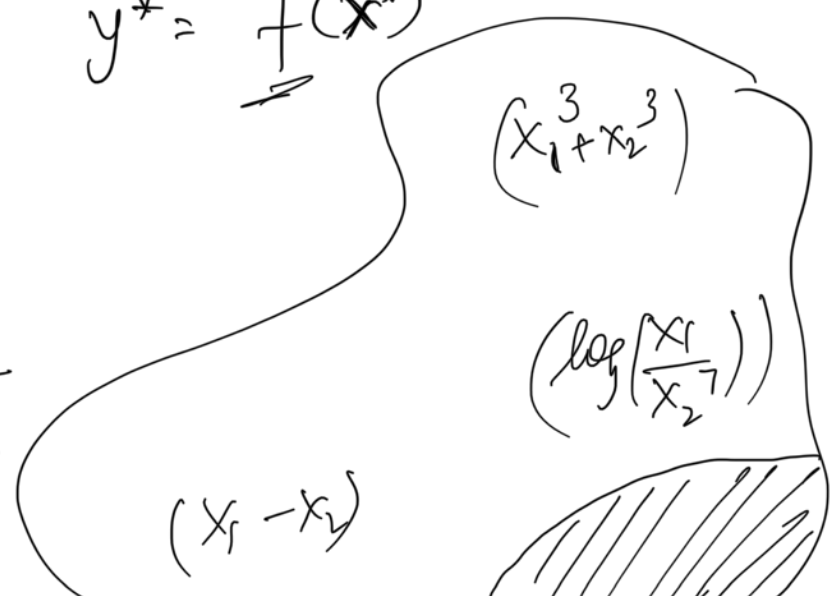
↑ some function.

If we learn $f()$,

then for a new x^* $x = (x_1, x_2)$

$$y^* = f(x^*)$$

a big bag of functions



Inductive bias

Monday Feb 8

$$x \rightarrow y$$

Functional models

$$y = f(x)$$

Probabilistic Models

$$P(x, y)$$

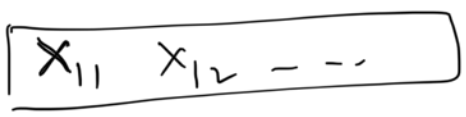
$$p(y|x) \leftarrow \text{Bayes Rule}$$

x_1



x, y, z

x_1, x_2



\mathbb{R}^D

$$x \cdot y = \sum_{i=1}^D x_i y_i$$

$$= x_1 y_1 + x_2 y_2 + \dots + x_D y_D$$

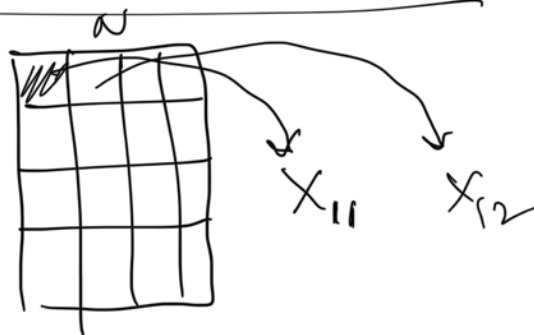
$$|x| = \sum_{i=1}^D |x_i|$$

$$\|x\|_2 = \sqrt{\sum_{i=1}^D x_i^2} \quad \text{l}_2 \text{ norm}$$

$$\|x\|_p = \left(\sum_{i=1}^D |x_i|^p \right)^{1/p}$$

Matrix

$$X \in \mathbb{R}^{M \times N}$$



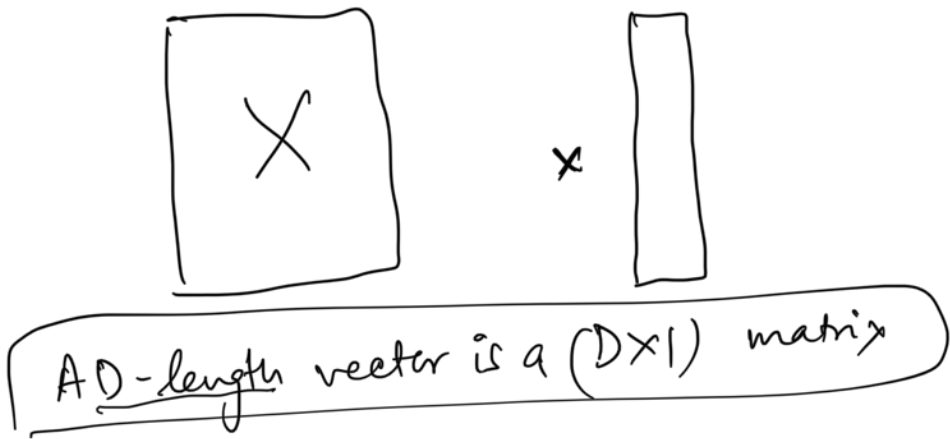
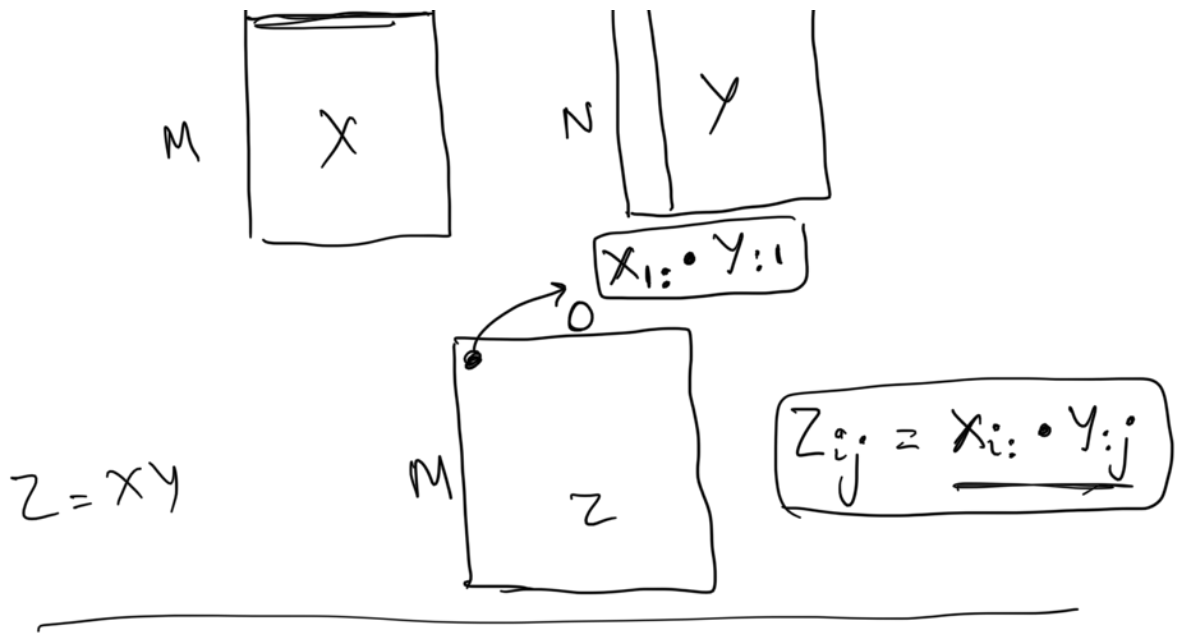
$$y = X^T \quad (N \times M)$$

$$cX =$$

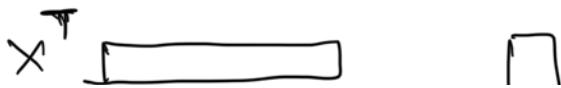


$$X \quad \underbrace{\quad}_N$$

$$y \quad \underbrace{\quad}_M$$



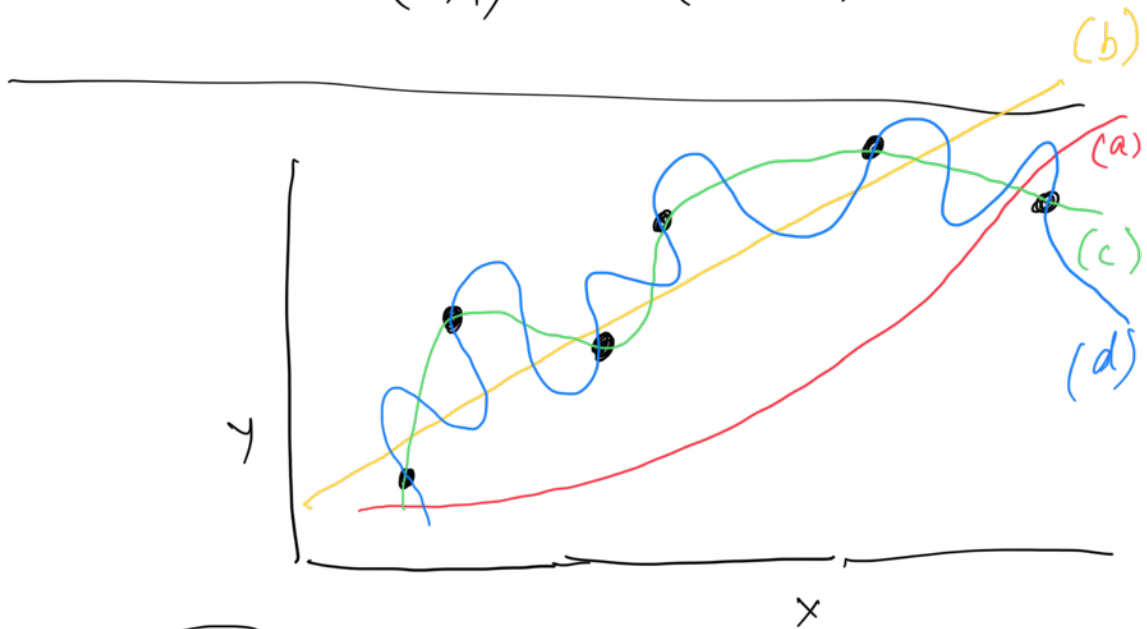
$$x \cdot y = x^T y$$



$(1 \times D)$

$\begin{bmatrix} | \\ | \\ | \\ | \\ | \end{bmatrix} y$
 $(D \times 1)$

$$Z = X^T y = (1 \times 1) = (X^T) \cdot y$$



error

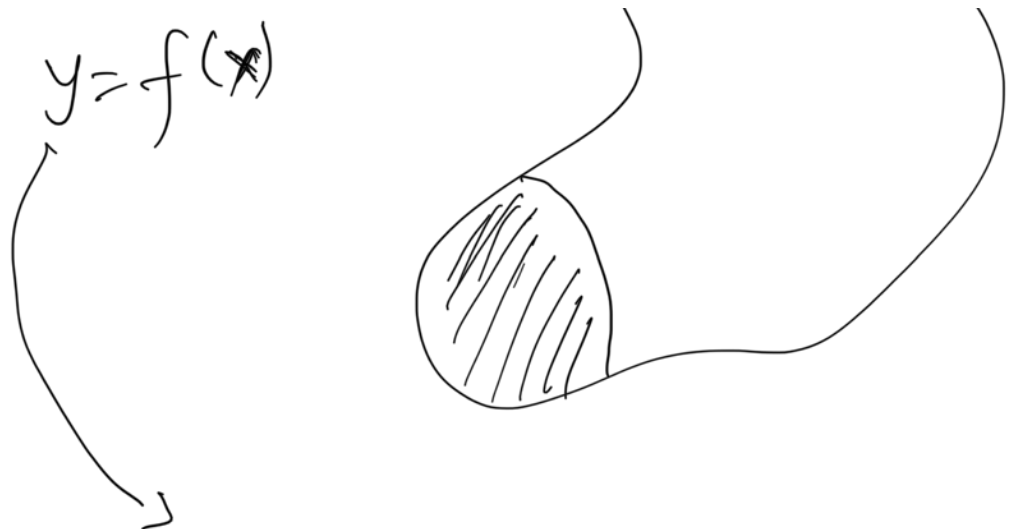
$$(b) \rightarrow y = \underline{m}x + \underline{c}$$

$$(c) \rightarrow y = \underline{a}x^3 + \underline{b}x^2 + \underline{c}x + \underline{d}$$

$$(d) \rightarrow y = \underline{a}x^{10} + \underline{b}x^9 + \dots + \dots$$

Linear Regression





$$y = f(x)$$

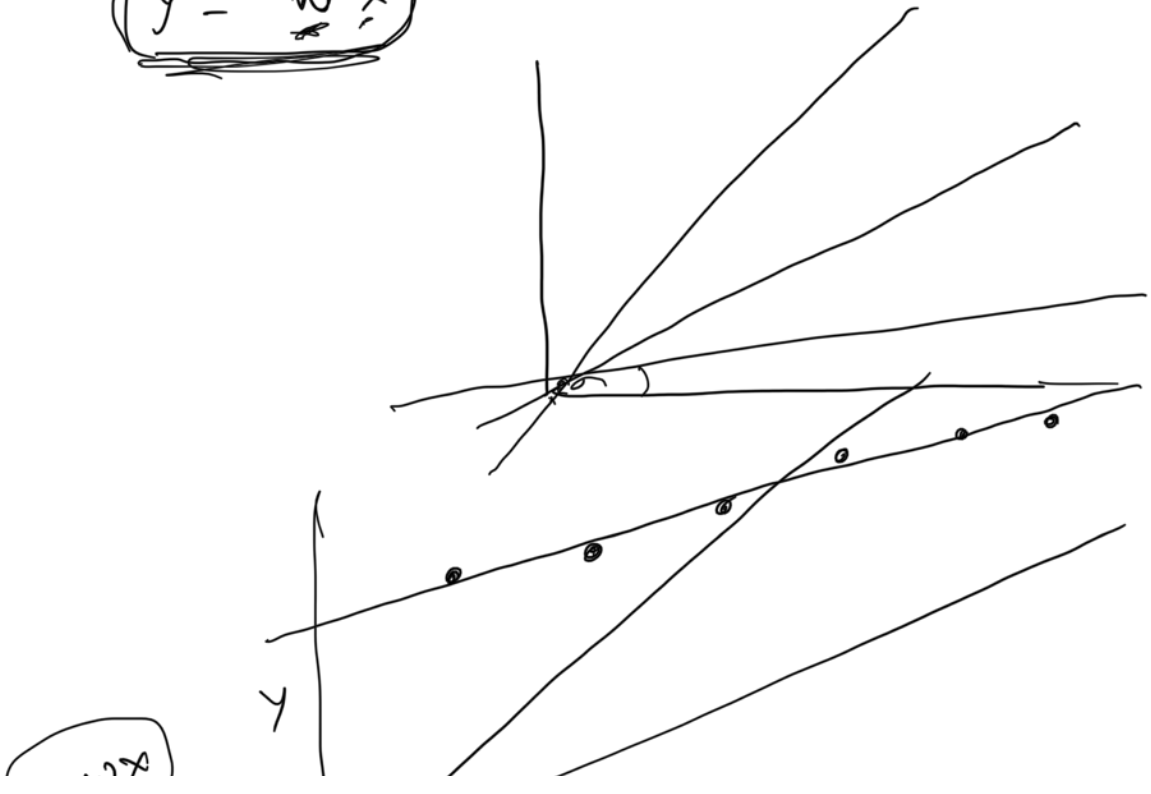
$$y = \underline{w^T x} \rightarrow D \times 1$$

weight vector
(D x 1)

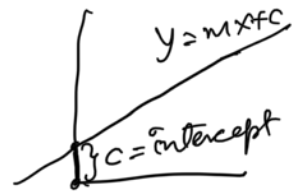
$$= w_1 x_1 + w_2 x_2 + \dots + w_D x_D$$

For the flu example

$$y = w^T x$$



$y = w_1 x$



$$y = w_0 + w_1 x$$

↖
bias-term

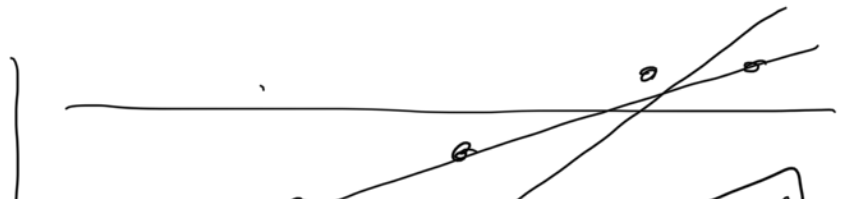
Given some data:

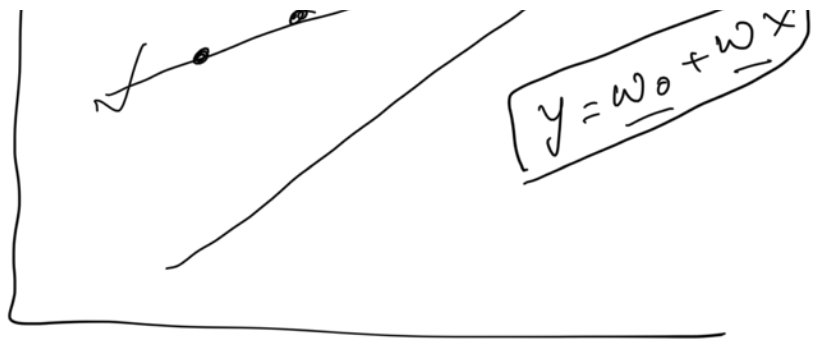
x	y
12	9
28	15
15	11
48	21
56	22

Find the best w_0, w_1 which "fits" the data best.

Wed, Feb 10

Resources on Piazza





Find w_0, w that does what?

<u>X</u>	<u>y</u>	<u>Pred. \bar{y}</u>
x_1	y_1	\bar{y}_1
x_2	y_2	\bar{y}_2
x_3	y_3	\bar{y}_3
\vdots	\vdots	\vdots
x_N	y_N	\bar{y}_N

For a given w, w_0

$$\bar{y}_1 = w_0 + w x_1$$

$$\bar{y}_i = w_0 + w x_i$$

$$\text{Error: } e_i = y_i - \bar{y}_i$$

$$J = \frac{1}{2} \sum_{i=1}^N e_i^2$$

$$J(w_0, w) = \frac{1}{2} \sum_{i=1}^N (y_i - (w_0 + w x_i))^2$$

Squared
loss function

$\frac{1}{2} \rightarrow$ just for mathematical
convenience

What if x is not 1-D

$\Rightarrow x \in \mathbb{R}^d$ where $d > 1$

$$y = w_0 + w^T x$$

$w \in \mathbb{R}^d$

Squared loss function

$$J(w) = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2$$

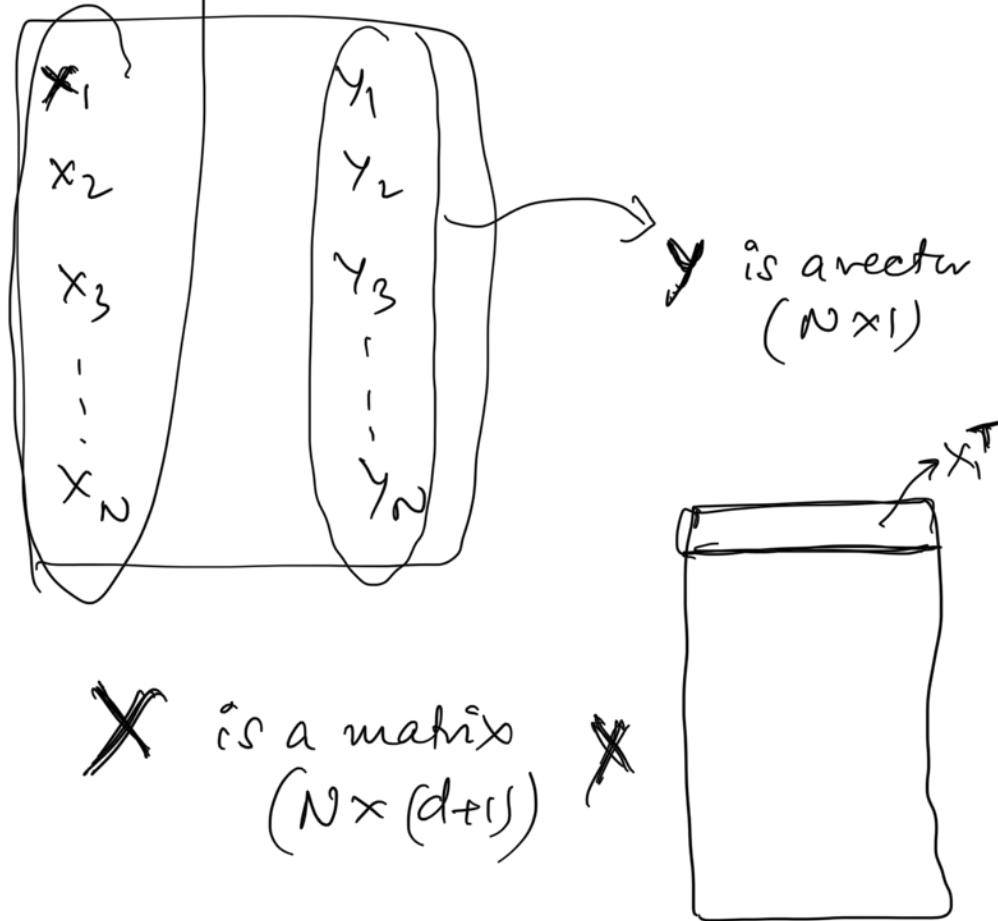
absorbed added a 1 to x_i

w_0 into w

Find w that minimizes $J(w)$

$$w^T x_i = \sum_{j=0}^d (w_j x_{ij})$$

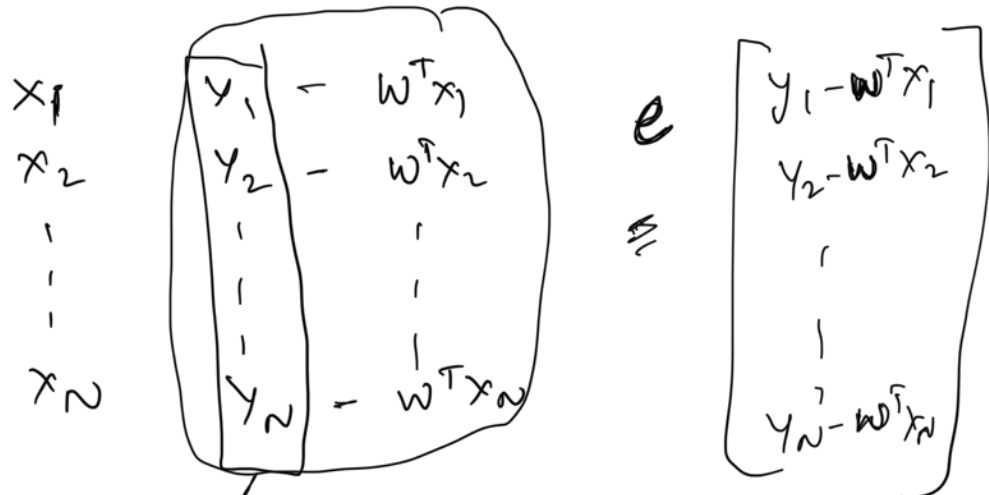
$w \cdot x_i$
 $w^T x_i$



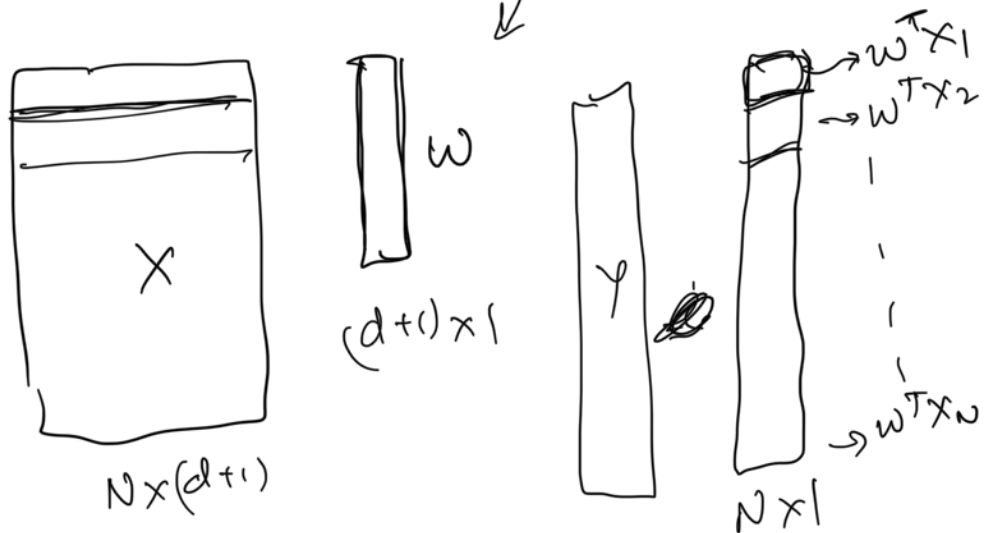
Training data: X, y
(N x (d+1)) (N x 1) target vector

data matrix

$$J(w) = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2$$



$$e = (y - Xw)$$



$$\sum (y_i - w^T x_i)^2 = (y - Xw)^T (y - Xw)$$

$$e_i = y_i - w^T x_i$$

$$e = \begin{bmatrix} y_1 - w^T x_1 \\ \vdots \\ y_N - w^T x_N \end{bmatrix}$$

any vector p

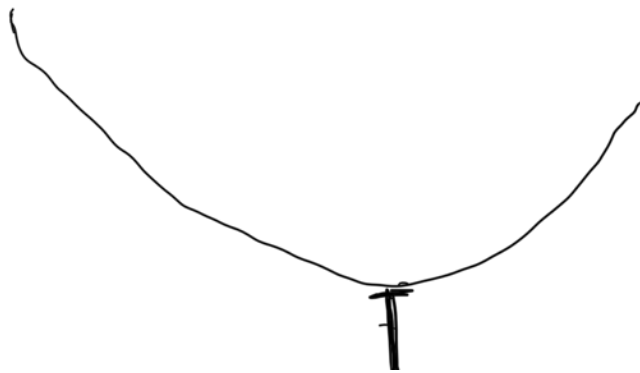
$$p^T p = \sum p_i^2$$

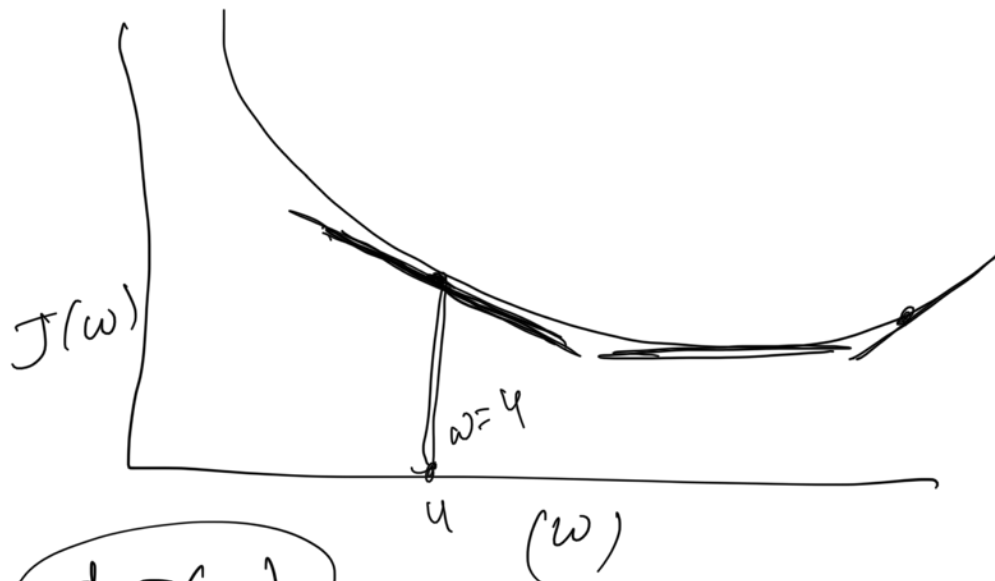
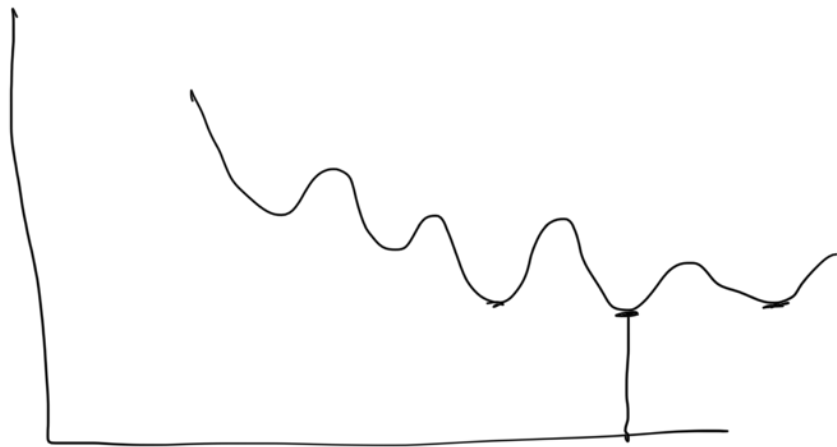
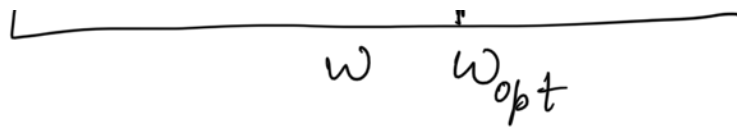
$$\sum (y_i - w^T x_i)^2 = (y - Xw)^T (y - Xw)$$

$$J(w) = \frac{1}{2} (y - Xw)^T (y - Xw)$$

How do we find w that minimizes $J(w)$

$J(w)$





$$\frac{dJ(w)}{dw}$$

e.g. $J(w) = 3w^3 + 4w$

$$\frac{dJ(w)}{dw} = 9w^2 + 4$$

at $w = 4$

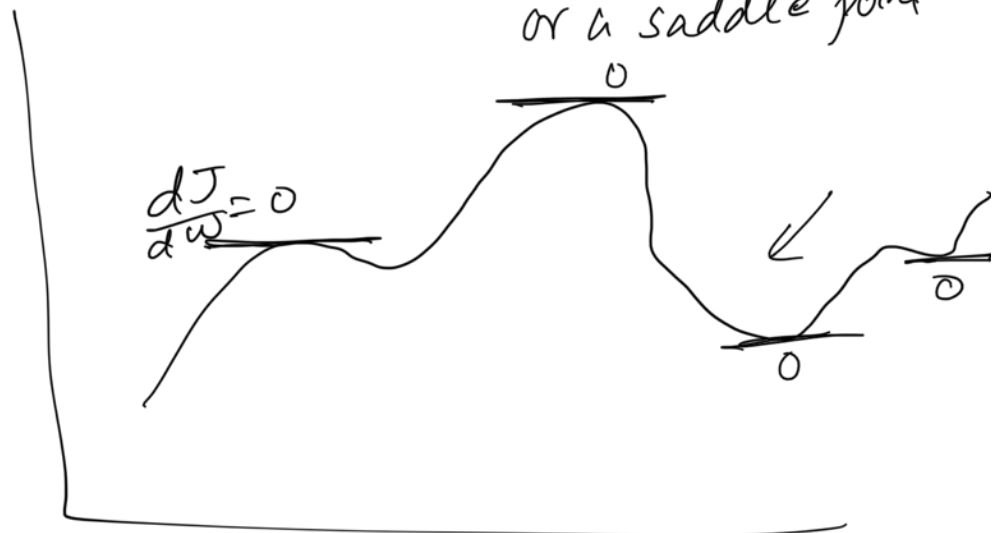
$$\frac{dJ(w)}{dw} = \underline{\underline{148}}$$

Q1 $dJ(w)$

if $\frac{dJ(w)}{dw} = 0$ at a given w

that means $w \rightarrow$ point of minima
or maxima

or a saddle point



If $J(w)$ is convex

then $\frac{dJ(w)}{dw} = 0$

solution for this
will give us
 w for which $J(w)$ is minimum

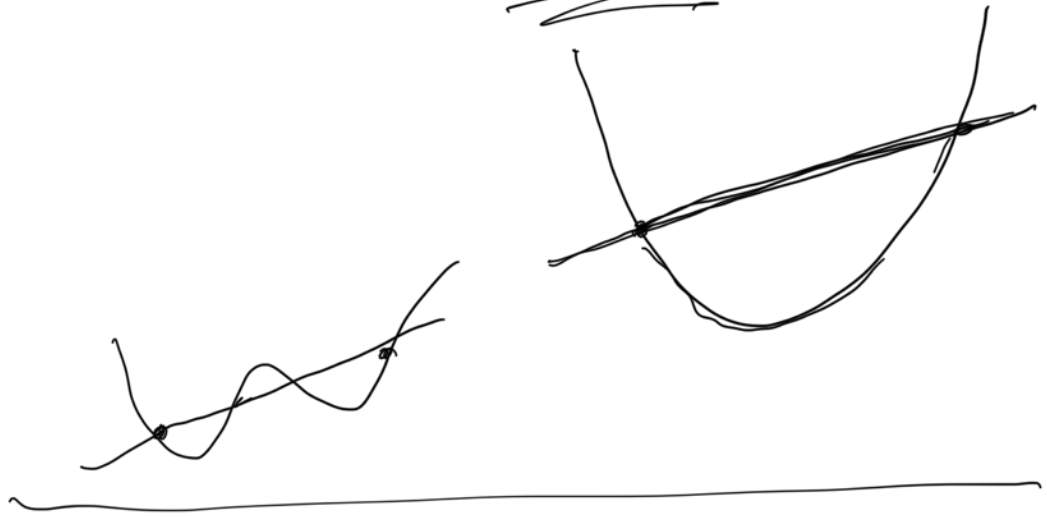
Feb 12

$$f(w) = \underline{7w^3 - 4w^2 + 8}$$

↑
scalar

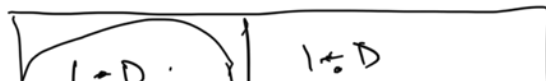
$$\frac{d}{dw} f(w) = f'(w) = \underline{21w^2 - 8w}$$

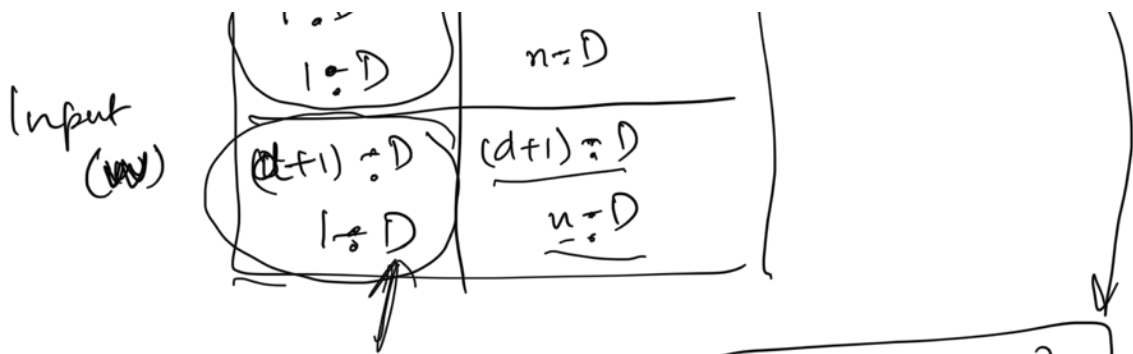
$$\frac{d^2}{dw^2} f(w) = \frac{d}{dw} \left(\frac{df(w)}{dw} \right)$$



Squared loss function for LR
is convex

$$J(w) = \frac{1}{2} (y - Xw)^T (y - Xw)$$





$$\frac{1}{2} \sum (y_i - w^T x_i)^2$$

$$\frac{d J(w)}{d w}$$

$$= \nabla J(w)$$

↑
gradient

$$f(w) = 3w_0^2 + 17w_1$$

$$w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\frac{\partial f(w)}{\partial w_0} = 6w_0 + 8$$

$$\frac{\partial f(w)}{\partial w_1} = 17$$

$$\frac{d f(w)}{d w} = \begin{bmatrix} \frac{\partial f(w)}{\partial w_0} \\ \frac{\partial f(w)}{\partial w_1} \end{bmatrix}$$

Gradient

$$\nabla f(w)$$

$$\frac{d^2}{d w^2} f(w)$$

$$\frac{d}{d w} \nabla f(w)$$

$$f(w) =$$

$$w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

$$\frac{d}{d w} f(w)$$

$$\begin{bmatrix} 7w_0^2 + 8w_1 \\ w_0 w_1 + 9w_1^3 \end{bmatrix}$$

$$= H \quad \left| \quad d w = \begin{bmatrix} 14w_0 & 8 \\ w_1 & w_0 + 27w^2 \end{bmatrix} \right.$$

\swarrow Hessian of $f(w)$ \swarrow Jacobian

$$J(w)$$

Calculate $\nabla J(w) \equiv \frac{d}{dw} J(w)$

and then solve for w : $\nabla J(w) = 0$

$$J(w) = \frac{1}{2} (y - Xw)^T (y - Xw)$$

$$J(w) = \frac{1}{2} [(y^T - Xw)^T] (y - Xw)$$

$$= \frac{1}{2} [(y^T - w^T X^T) (y - Xw)]$$

$$= \frac{1}{2} \left[(y^T y) - \frac{y^T Xw}{\quad} - \frac{w^T X^T y}{\quad} + w^T X^T X w \right]$$

$$= \frac{1}{2} [y^T y - 2 w^T X^T y$$

$$(A+B)^T = A^T + B^T$$

$$(AB)^T = B^T A^T$$

$$a^T b = b^T a$$

$$\frac{y^T Xw}{2} \quad (w^T X)$$

$$+ w^T x^T x w \quad \left. \begin{array}{l} \rightarrow y^T z \\ \underline{z^T y} \end{array} \right\}$$

$$\begin{aligned} \frac{d}{dw} J(w) &= \frac{1}{2} \left[\frac{d}{dw} \underline{y^T y} - 2 \frac{d}{dw} \underline{w^T x^T y} \right. \\ &\quad \left. + \frac{d}{dw} \underline{w^T x^T x w} \right] \\ &= \frac{1}{2} \left[-2 x^T y + 2 x^T x w \right] \\ &= -x^T y + x^T x w \end{aligned}$$

$$\text{Set } \frac{d}{dw} J(w) = 0 \quad \left[\nabla f(x) = 0 \right]$$

$$-x^T y + x^T x w = 0$$

$$x^T x w = x^T y$$

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

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

$X - N \times (d+1)$

$$\boxed{A^{-1} A = I}$$

$$\boxed{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}$$

$$y - N \times 1$$

$$X^T X = (d+1) \times (d+1)$$

Scalability

$$A_{d \times d}$$

$$A^{-1}$$

$$O(d^3)$$

$$\frac{1}{\kappa} \rightarrow \kappa \rightarrow 0$$

A^{-1} is not easy
to calculate

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

$$A \rightarrow z$$

Mon, Feb 15

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

w

Evaluation Metric

$$(y_i - \hat{y}_i)^2$$

$$+ \begin{pmatrix} y_2 - \hat{y}_2 \end{pmatrix}^2 + \begin{pmatrix} y_3 - \hat{y}_3 \end{pmatrix}^2 + \begin{pmatrix} y_4 - \hat{y}_4 \end{pmatrix}^2 + \begin{pmatrix} y_5 - \hat{y}_5 \end{pmatrix}^2$$

Root Mean Squared Error

(RMSE) $\sqrt{\frac{1}{N} \sum (y_i - \hat{y}_i)^2}$ $\hat{y}_i = \omega^T x_i$

unit = cm

Squared Loss

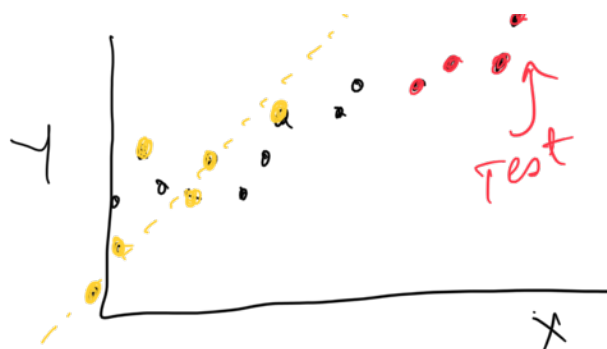
$\frac{1}{2} \sum (y_i - \hat{y}_i)^2$ $\hat{y}_i = \omega^T x_i$

What data to measure error on?

Training data?

Generalization Error

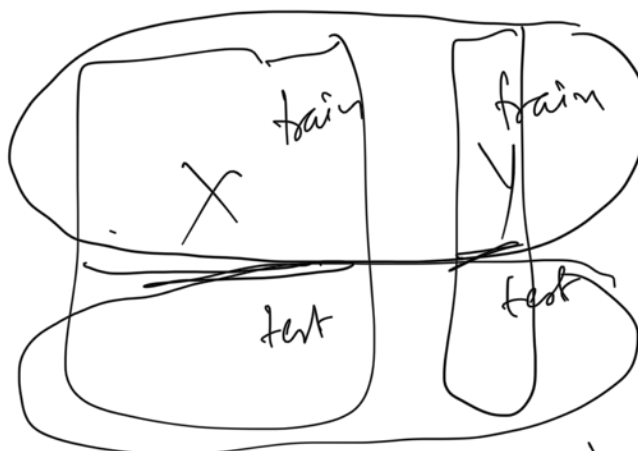
Testing error



$$\text{MAE: } \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$$

Scaling \longleftrightarrow

$$\text{Price} = 0.108 \text{ CRIM} + 0.0464 \text{ ZN} + \dots$$



GD - LR

Choose some w (init)

$$\text{Calculate } J(w) = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2$$

Calculate ∇J

$$w \leftarrow w - \eta \nabla J$$

Choose some w (init)

for $i = 1$ to N :

while not converged

$$J(w) = \frac{1}{2} (y_i - w^T x_i)^2$$

∇J

$$w \leftarrow w - \eta \nabla J$$

Stochastic Gradient
Descent

SGD

