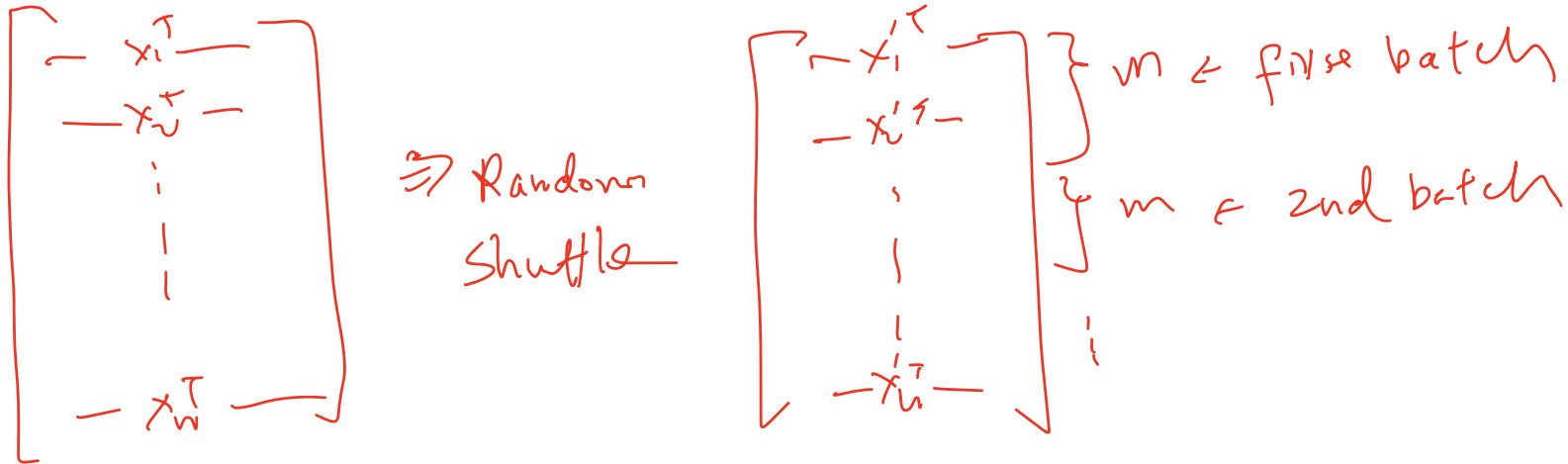


Linear Regression

Recap on Optimization

T/F: for mini-batch SGD, we should always use very large batch

Q: for mini-batch w/ size m , should we sample m points with replacement?



Objective

Learn the first regression algorithm (i.e., predict continuous variable)

Outline for Today

1. Intro on Linear Regression

2. Normal equation for linear Regression

3. Interpretation of Linear Regression using MLE / MAP

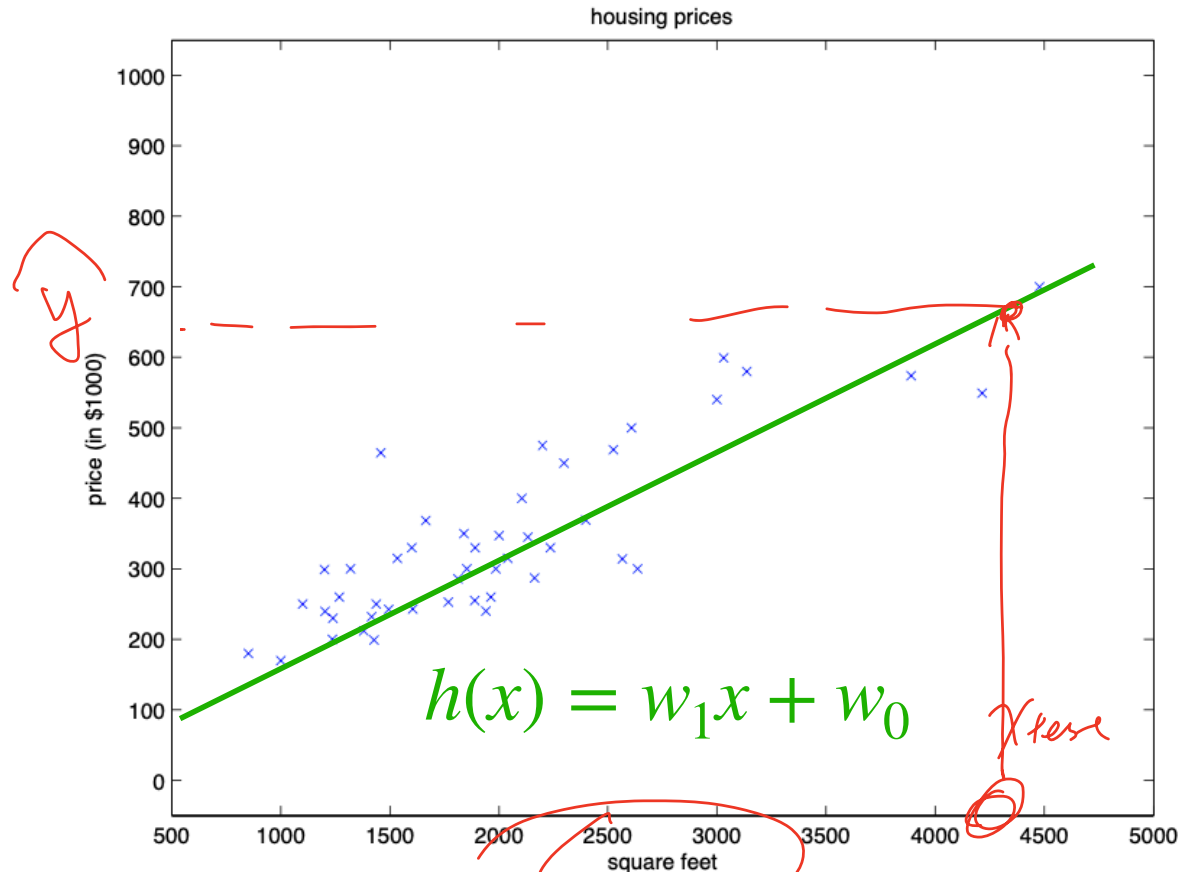
Ex: Predicting the house price

Dataset:

Living area (feet ²)	Price (1000\$)
2104	400
1600	330
2400	369
1416	232
3000	540
⋮	⋮
x	y

(Example from Stanford CS229)

Plot:



Ex: Predicting the house price (2d case)

Dataset:

Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮
$x[1]$	$x[2]$	y

Goal: finding the linear function

$$h(x) = w_1x[1] + w_2x[2] + w_0$$

that fits the data well

Ex: Predicting the house price (2d case)

Dataset:

Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮
$x[1]$	$x[2]$	y

As usual, we append 1 to the feature, i.e.,

$$x = \begin{bmatrix} x[1] \\ x[2] \\ 1 \end{bmatrix}$$

So the linear function can be written as:

$$h(x) = w^T x$$

w = $\begin{bmatrix} w(1) \\ w(2) \\ w_0 \end{bmatrix}$

Outline for Today

1. Intro on Linear Regression

2. Normal equation for linear Regression

3. Interpretation of Linear Regression using MLE / MAP

Mathematical formulation of linear regression

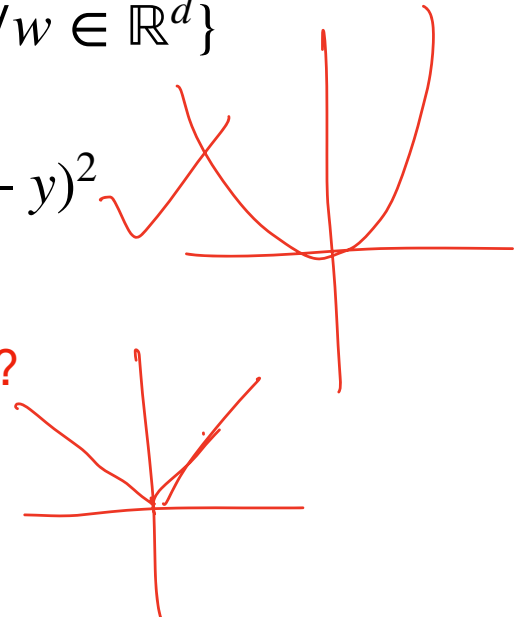
Input: dataset $\mathcal{D} = \{x_i, y_i\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$

Hypothesis: linear function $h(x) = w^\top x$

Hypothesis class: all possible linear functions $\{w^\top x, \forall w \in \mathbb{R}^d\}$

Loss function: squared loss $\ell(w^\top x, y) = (w^\top x - y)^2$

Q: can we use absolute loss, i.e., $|w^\top x - y|$?



Mathematical formulation of linear regression

$$\left(\frac{w^T x - y}{\quad} \right)^2$$

Formulating the optimization problem:

$$\arg \min_w \sum_{i=1}^n (w^T x_i - y_i)^2 \quad := \underline{\underline{l(w)}}$$

Convex

Q: how to solve this?

Linear regression solution

$$X = \begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_n \\ | & | & & | \end{bmatrix}_{d \times n}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$\arg \min_w \sum_{i=1}^n (w^T x_i - y_i)^2 \quad \checkmark$$

$$X^T w - Y$$

$$= \begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_n^T \end{bmatrix} w - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Let's compute the closed-form solution:

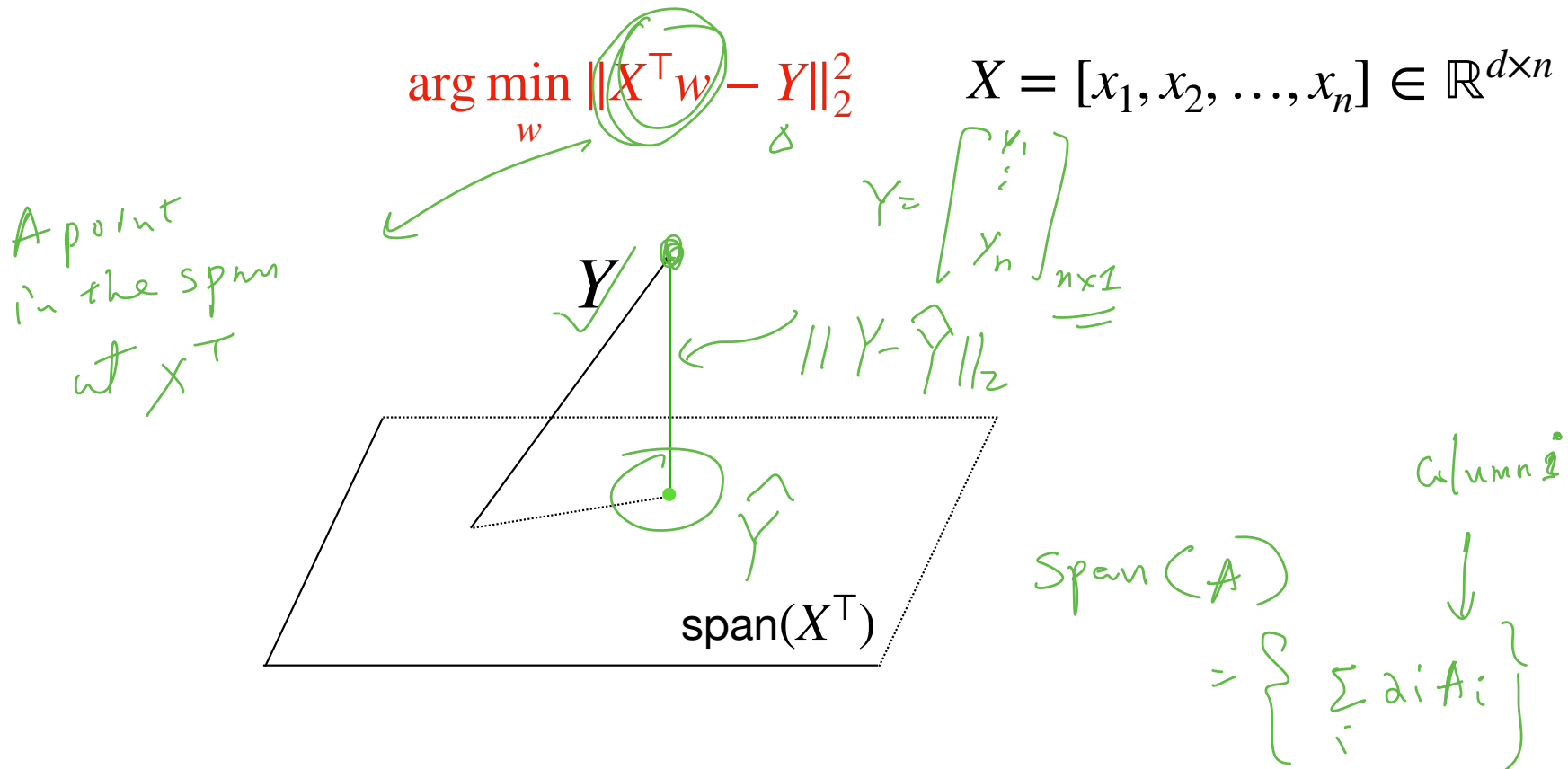
Define $X = [x_1, x_2, \dots, x_n] \in \mathbb{R}^{d \times n}$, $Y = [y_1, \dots, y_n]^T \in \mathbb{R}^n$

$$\sum_{i=1}^n (w^T x_i - y_i)^2 = \|X^T w - Y\|_2^2$$

$$\Rightarrow \arg \min_w \|X^T w - Y\|_2^2$$

$$\begin{aligned} &= \sum_{i=1}^n (x_i^T w - y_i)^2 \\ &= \|X^T w - Y\|_2^2 \end{aligned}$$

Linear regression solution



$$X = \begin{bmatrix} | & | & | \\ x_1 & x_2 & \dots & x_n \\ | & | & | \end{bmatrix}$$

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

$$XX^T \in \mathbb{R}^{d \times d}$$

Linear regression solution

$$\arg \min_w \|X^T w - Y\|_2^2$$

$$\forall \|X^T w - Y\|_2^2$$

$$\nabla_w \|X^T w - Y\|_2^2 = 2X(X^T w - Y)$$

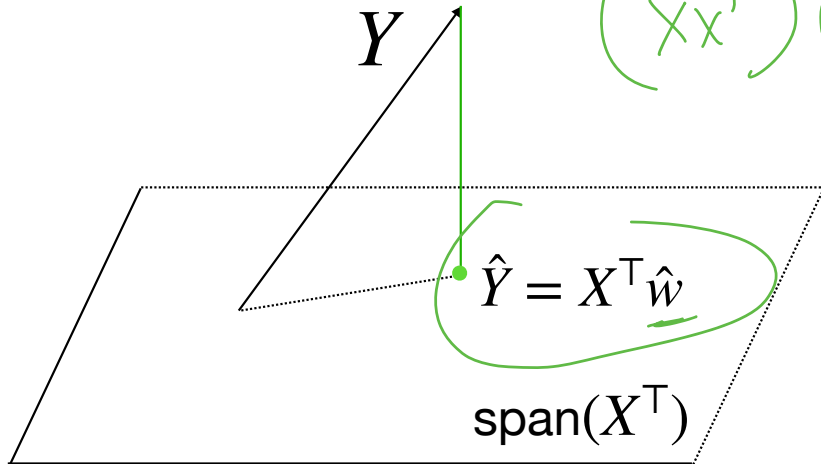
$$= 2X(X^T w - Y)$$

if XX^T is full rank, then $\hat{w} = (XX^T)^{-1}XY \in \mathbb{R}^d$

$$(XX^T)w - XY = 0 \Rightarrow w = (XX^T)^{-1}(XY)$$

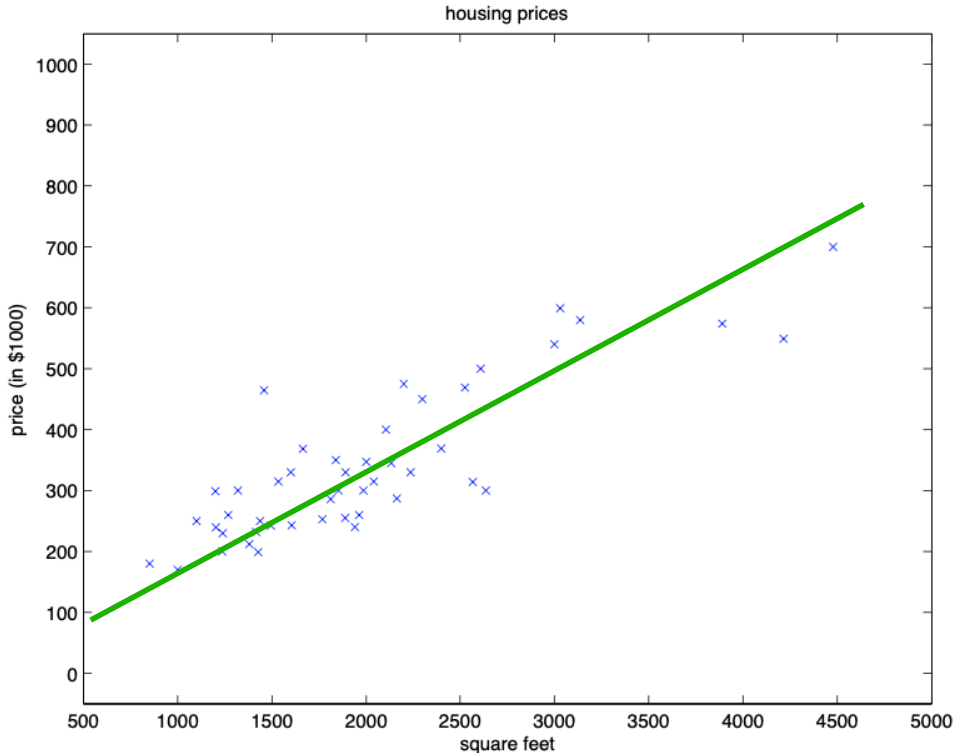
What if XX^T is not full rank?

(We will talk about regularization soon)



Prediction using linear regression

Once we learned \hat{w} , we can use it to make prediction on any new feature x



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

Given x_{test} , our prediction is:

$$\hat{y} = x_{test}^T \hat{w}$$

$$XY = \sum_{i=1}^n x_i y_i$$

$$= x_{test}^T (X^T X)^{-1} X^T Y$$

$$= \sum_i \left(x_{test}^T (X^T X)^{-1} x_i \right) \cdot y_i$$

$\beta_i \in \mathbb{R}$

Outline for Today

1. Intro on Linear Regression

2. Normal equation for linear Regression

3. Interpretation of Linear Regression using MLE / MAP

Derive Linear regression via Maximum Likelihood Estimation

Assume $P(y | x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - \overset{\text{mean}}{x^T w})^2 / \sigma^2\right)$, i.e., $y = w^T x + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$

$$Z = \sqrt{2\pi} \sigma$$

Let's maximize the log-likelihood of the data, i.e.,

$$\arg \max_w \sum_{i=1}^n \ln P(y_i | x_i; w)$$

$$= \arg \max_w \sum_{i=1}^n \left[-\frac{1}{2\sigma^2} (w^T x_i - y_i)^2 - \ln(Z) \right]$$

$$= \arg \min_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

$$P(D | w) \stackrel{\text{i.i.d.}}{=} \prod_{i=1}^n P(x_i, y_i | w)$$

Bayes rule $i=1$

$$= \prod_{i=1}^n P(y_i | x_i; w) \cdot \underbrace{P(x_i | w)}_{P(x_i)}$$

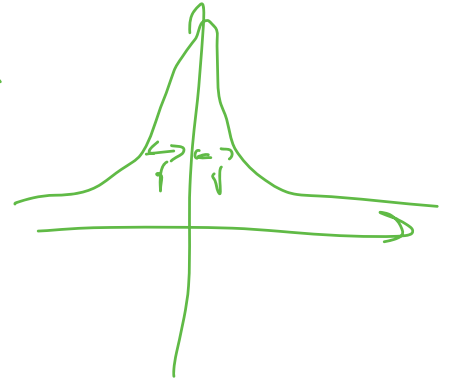
$$= \prod_{i=1}^n P(y_i | x_i; w) \cdot P(x_i)$$

Derive Linear regression via MAP

Assume $P(y | x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - x^\top w)^2 / \sigma^2\right)$, i.e., $y = w^\top x + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$

To use MAP, we need to define a prior over w , we use Gaussian as well here:

$$w \sim \mathcal{N}(0, r^2 I) \quad \text{prior}$$



$$p(w) = \frac{1}{Z} \exp\left(-\frac{1}{2} w^T w / r^2\right) \Leftrightarrow \mathcal{N}(0, r^2 I)$$

Derive Linear regression via MAP

$$w \sim \mathcal{N}(0, r^2 I) \quad P(y|x; w) = \frac{1}{Z} \exp\left(-\frac{1}{2}(y - x^T w)^2 / \sigma^2\right)$$

MAP:

$$P(w|D) \propto P(w) \cdot P(D|w)$$

Bayes Rule

$$\arg \max_w \ln P(w|D)$$

$$= \arg \max_w \ln P(w) + \ln P(D|w) \quad \text{same in MLE}$$

$$= \arg \max_w \left(\frac{w^T w}{2r^2} + \sum_{i=1}^n \frac{1}{2\sigma^2} (w^T x_i - y_i)^2 \right) \cdot \cancel{Z}$$

$$= \arg \min_w \frac{\sigma^2}{r^2} w^T w + \sum_{i=1}^n (w^T x_i - y_i)^2 \quad \left(= \arg \min_w \lambda \|w\|_2^2 + \sum_{i=1}^n (w^T x_i - y_i)^2 \right)$$

Ridge Linear Regression

Regularization

$\lambda \geq 0$

$$\arg \min_w \lambda \|w\|_2^2 + \sum_{i=1}^n (w^\top x_i - y_i)^2 \quad \checkmark \leftarrow \text{MAP}$$

In this case, we can derive a closed-form solution as well:

$$\hat{w} = \underline{\underline{(XX^\top + \lambda I)^{-1}XY}}$$

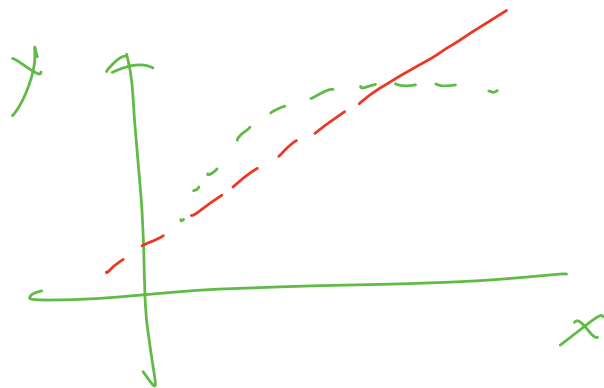
Note that it works even XX^\top is not full rank

PSD \checkmark $XX^\top + \lambda I$ is PD

$\lambda \rightarrow \infty$
then $\hat{w} = \vec{0}$

$\lambda \rightarrow 0$

$\hat{w} = (XX^\top)^{-1}XY$



Summary for today

1. Linear regression, Normal equation, and MLE / MAP interpretation
2. Your take-home question: what is the SGD update rule for Linear regression? Is the update rule intuitively explainable?
3. Next Tue: Support Vector Machine!