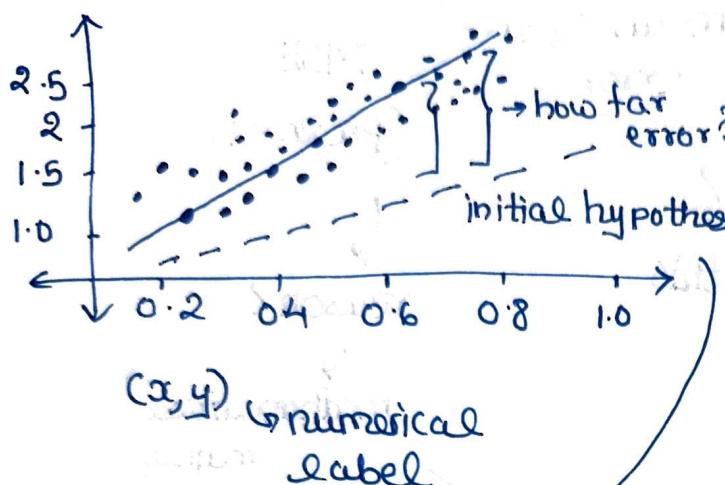


## \*Parametric Learning, Linear Regression



$y = ax + b$  } parametric learning

modelling } hypothesis

models  $\rightarrow$  specific  $a$  &  $b$  }

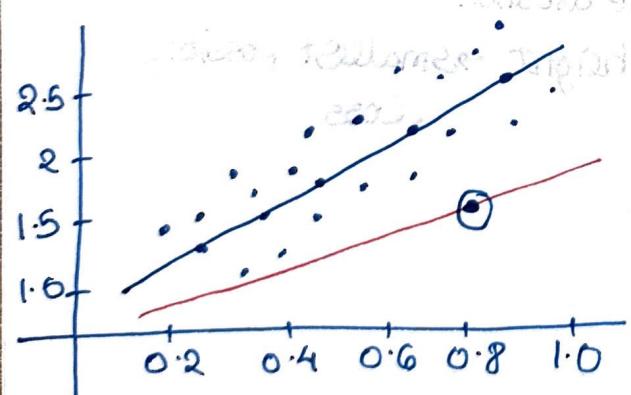
at one point we find the best possible hypothesis

Linear regression

$y = f(x)$   $\rightarrow$  possible hypothesis bigger not just one

we correct it.  $\hookrightarrow$  multiple iteration  $\rightarrow$  practice makes a man perfect.

## \*Loss function for Linear Regression:-



$y = ax + b$

way to quantify how good a model is?

$$\hat{y} = a(0.8) + b$$

we maybe know the correct value is

$$X_i: \hat{y}_i = aX_i + b$$

↓  
predicted

$$y_L$$

correct label

Error  
in  
Correcting

$$e_i = |\hat{y}_i - y_L| = |ax_i + b - y_L|$$

error depends on  $a$  &  $b$

Loss fn  
 $L_{MAE}(a, b)$

predicted

$$\sum_{i=1}^n e_i = \frac{e_1 + \dots + e_n}{n}$$

Mean absolute error

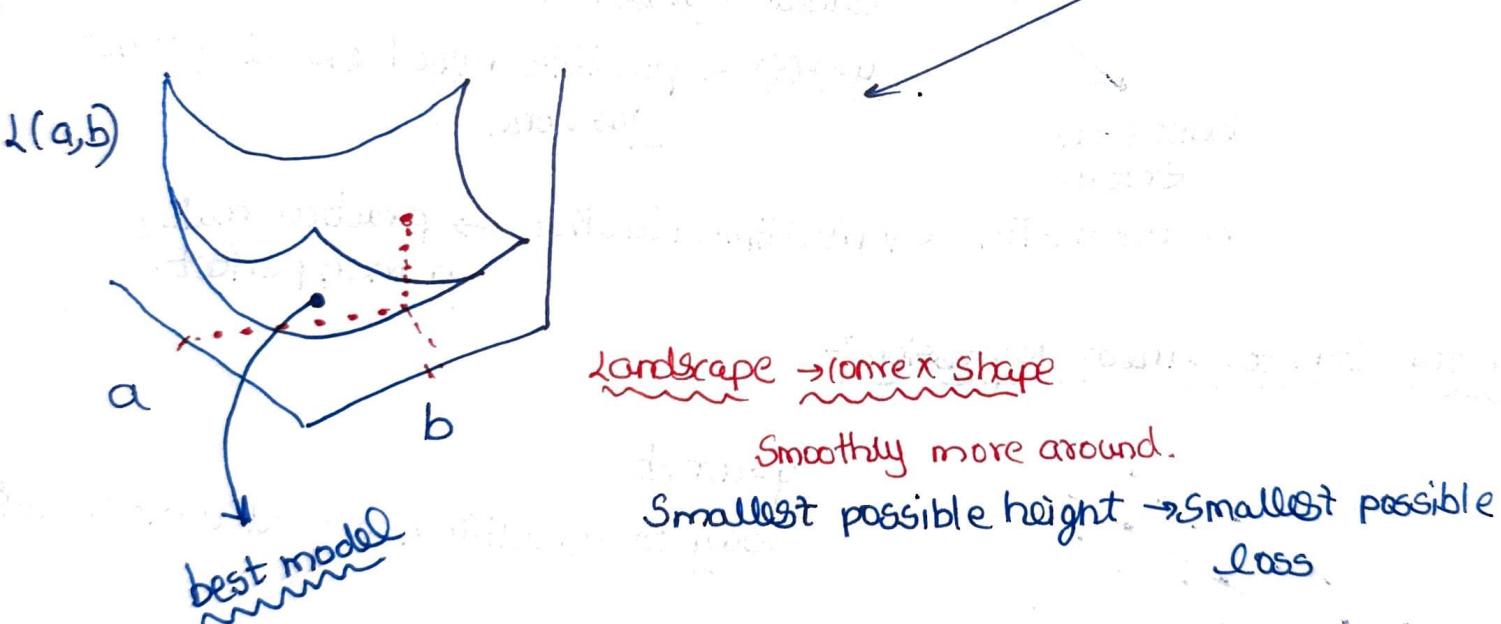
## Other type of loss function

$$L_{\text{MSE}}(a, b) = \frac{1}{n} \sum_{i=1}^n e_i^2 \quad \begin{array}{|l} \text{mean squared} \\ \text{error} \end{array}$$

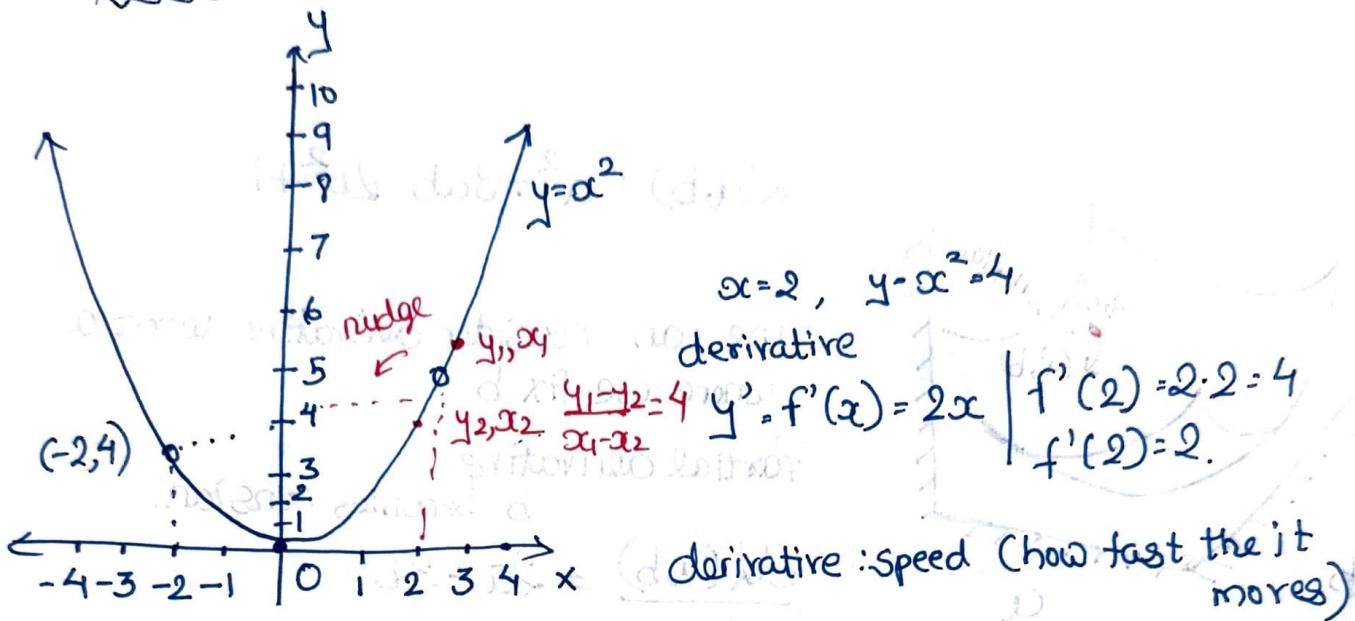
MSE  
(preferred)

Which one is better?  
We can maybe make this  
hyperparameter

mathematical  
mainly



## \* Gradient Descent (One variable)



Sign → Shows Speed / things

$\text{if } \alpha > 0 \rightarrow +$

$\alpha = 0 \rightarrow 0$

$\alpha < 0 \rightarrow -$

↓  
value of  $f(x)$  for values  
like  $f(2)$

Sign ( $f'(2)$ ) = +

its increasing  $f'(2) = 2 \cdot 2 = 2$

$f'(3) = 2 \cdot 3 = 6$

$f'(4) = 2 \cdot 4 = 8$

$f'(-2) = 2(-2) = -4$

Sign ( $f'(-2)$ ) = -

value is 4 but in negative

nudge

derivative

$$X := [X - \text{step\_size} \cdot f'(x)] \quad \text{step } 0.1$$

old X

$$f'(x) = 2x$$

$$X := 2 - 0.1(2 \cdot 2)$$

$$X := 2 - 0.1(4) = 1.6$$

} descended towards minimum

We need to be careful

If we take Step Size = 2

$$X = 2 - (2 \cdot 4) = -6$$

If we overnudge  
the we overshoot,

Correct direction but we overshoot.

→ This is called gradient descent

$X := 1.6 - 0.1(2 \cdot 1.6) \rightarrow$  descend towards minimum.

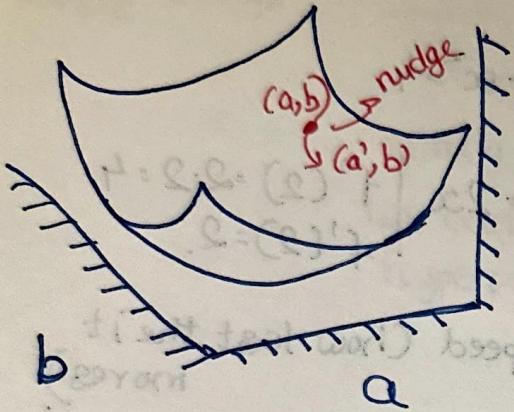
$$X := 1.6 - 0.32$$

$$X = 1.28$$

We can do iterations and go towards minimum.

## \* Gradient Descent (Multiple Variables)

$$L(a, b) = a^2 + 3ab + 2b^2 + 1$$



We can consider derivative w.r.t a where we fix b

Partial derivative

b becomes constant

$$\frac{dL(a, b)}{da} = 2a + 3b$$

$$+ = ((\Delta)^2) \text{ a.p.e}$$

$$\frac{dL(a, b)}{db} = \begin{aligned} & 3a + 2b(2) \\ & = 3a + 4b \end{aligned}$$

Now we gotta nudge it

new a old a  $\rightarrow ((\Delta)^2) \text{ a.p.e}$

$\rightarrow$  how fast a moves  
to fit to data

$$a := a - \text{step-size} \cdot \frac{dL(a, b)}{da}$$

$$b := b - \text{step-size} \cdot \frac{dL(a, b)}{db} \rightarrow \text{we have one for } b \text{ too}$$

$\rightarrow$  we do two updates.

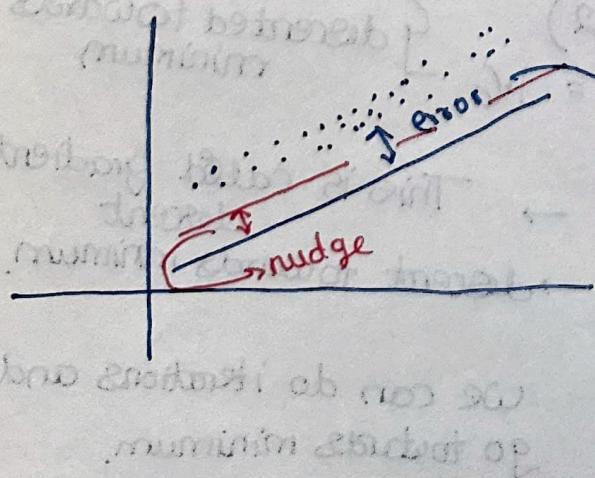
So we nudge it

and it moves

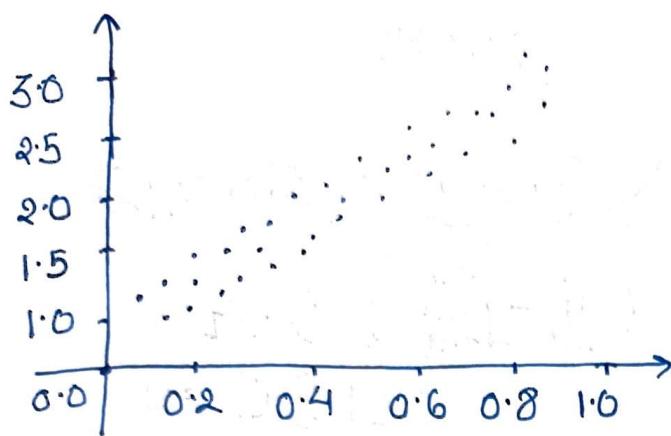
closer to our set of points

Gradient descent fits it better to

the data



## \* Gradient Descent for mSE



$$y = ax + b \rightarrow \text{intercept}$$

$\downarrow$   
slope coefficient.

Belonging to  $a$

$$L(a, b) = \sum_{i=1}^n e_i^2$$

$\rightarrow$  predicted

$$- \sum_{i=1}^n (\hat{y}_i - y_i)$$

$\rightarrow$  correct  $y_i$

labeled

$$= \sum_{i=1}^n (ax_i + b - y_i)^2$$

$$= \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\frac{\partial L}{\partial a} = \frac{1}{n} \frac{\partial e_i^2}{\partial a} + \dots + \frac{\partial e_n^2}{\partial a} = \frac{2}{n} (\hat{y}_1 - y_1) \cdot x_1 + \dots + 2(\hat{y}_n - y_n) x_n$$

Gradient Descent for  
mSE

$$\frac{\partial L}{\partial a} = \frac{2}{n} \sum_{j=1}^n (\hat{y}_j - y_j) \cdot x_j$$

$$\frac{\partial L}{\partial b} = \frac{2}{n} \sum_{j=1}^n (\hat{y}_j - y_j) \cdot 1$$

$$a = a - \text{step\_size} \frac{\partial L}{\partial a}$$

$$b = b - \text{step\_size} \frac{\partial L}{\partial b}$$

Example:-

$$x = \begin{pmatrix} 1 \\ -1 \\ 2 \end{pmatrix}, y = \begin{pmatrix} 1 \\ 2 \\ 5 \end{pmatrix}$$

$$y = x + 2$$

$$a=1, b=2$$

no of points ( $n$ ) = 3

$$(\hat{y}_1 - y_1)x_1 = (1-1) \cdot x_1 \\ = (1-1)(1) = 0$$

$$(\hat{y}_2 - y_2)x_2 = (-1-2)(-1) \\ = 3$$

$$(\hat{y}_3 - y_3)x_3 = (2-5)(2) \\ = -3(2) = -6.$$

$$\frac{dL}{da} = \frac{2}{n} \sum_{j=1}^n (\hat{y}_j - y_j) \cdot x_j$$

$$\frac{dL}{da} = \frac{2}{3} (0+3+(-6))$$

$$\frac{dL}{da} = \frac{2}{3} (-3) = -2$$

Step size ( $\rho$ ).

$$a' = a - \rho \frac{dL}{da}$$

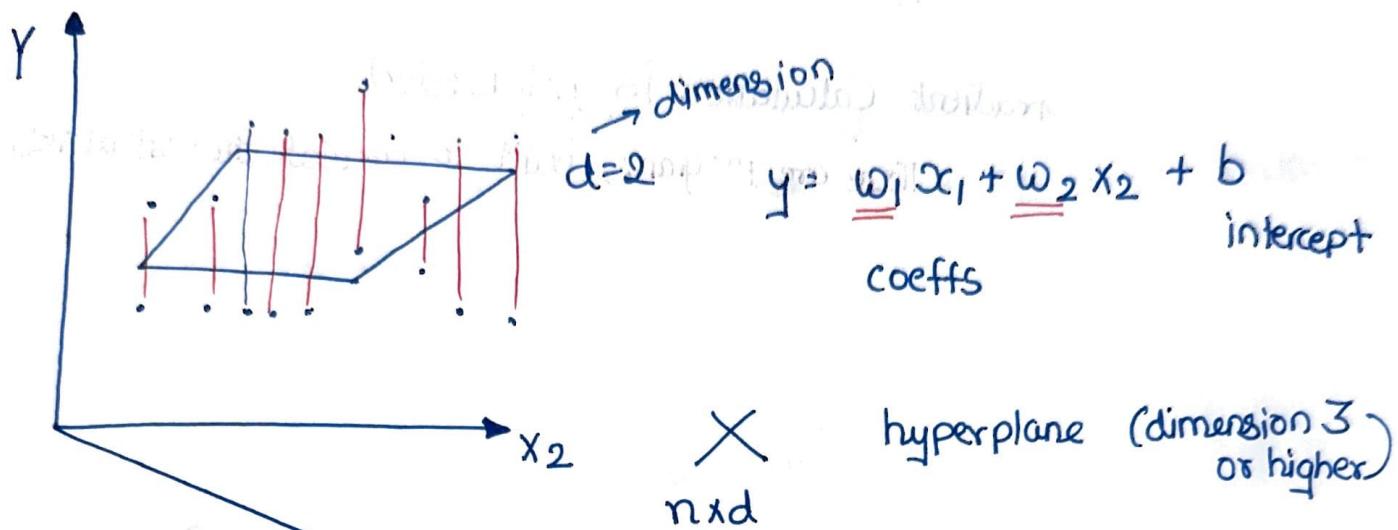
$$a = a - \rho(-2)$$

$$a = a - \rho(-2)$$

$$a = 1 - \rho(-2)$$

whatever step size ( $\rho$ ) is  
we get new  $a$ .

## \* Generalization in higher Dimensions



$y = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$

$y = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b$

$$X = [x_1, \dots, x_d] \quad Y = X \cdot w^T + b$$

w = [w<sub>1</sub>, ..., w<sub>d</sub>] dot product.

$$\begin{bmatrix} x_1 & \dots & x_d \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

matrix

update Rules?

update coeff w

$$w_i = w_i - \rho \sum_{j=1}^n (\hat{y}_j - y_j) X_{j,i}$$

new old (step size) n

for  $A \in \mathbb{R}^{d \times d}$

$w_i$  coordinate column i

$\rho$  learning rate

$\{\text{independently}\}$  nudge

$\{\text{add}\}$   $b$

$$b = b - \rho \sum_{j=1}^n (\hat{y}_j - y_j) \cdot 1$$

b bias

$\rho$  learning rate

$\{\text{add}\}$   $b$

P → every step we learn from our mistakes

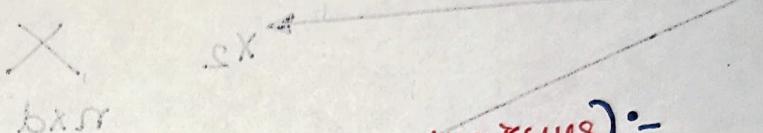
Gradient Calculation is not needed.

There are programs built to calculate the derivatives

function

27300

( $\sum$  noise) enqrep  
Graph 10



\* Gradient Descent: The algorithm (efficiency & implementation issues) :-

function  $w = \text{Linear Regression-fit}(X, y)$   
Initialize randomly a d-dimensional vector  $w$  and a scalar  $b$

for  $j = 1$  to  $n$  {epochs} :

$$b = b - \rho \sum_{j=1}^n (b + \sum_i w_i x_{ij}) - y_j$$

predicted label      intercept

labels

every iteration all data points

} one update per loop  
sum of individual error

$$w_i = w_i - \rho \sum_{j=1}^n (b + \sum_k w_k x_{kj} - y_j) \cdot x_{ji}$$

return  $w, b$

Solved today

function  $w = \text{Linear Regression-Predict}(w, b, x)$  ← unseen point

$$return x w^T + b$$

Epoch: entire dataset →

random state

① initialization

issue → reproducibility (Get back to same values or reproduce a good fit)  
pseudo random number

Hyperparameter

② Learning Rate  $\rho$

③ In every epoch we consider entire

can change by iteration

dataset, individual error calculation, a lot work  $O(n \cdot d)$   
Intuitively doesn't make sense, do we need this much?

\* Scaling features for faster learning?

1) For each attribute  $i$ , calculate average

$$\bar{x}_{j,i} = \frac{1}{n} \sum_{j=1}^n x_{j,i}$$

2) for each attribute, calculate Std Dev

$$s_i = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_{j,i} - \bar{x}_{j,i})^2}$$

3) For each point  $j$  and attribute  $i$ , compute a new attribute.

$$x'_{j,i} = \frac{x_{j,i} - \bar{x}_{j,i}}{s_i}$$

$s_i$   $\rightarrow$  std dev

$$\text{mean}(x_i) = 0$$

$$\text{std}(x_i) = 1$$

$L(w_1, w_2, b)$ : Speed towards bottom/convergence

So we do scaling of features.

$$X = \begin{pmatrix} | & 0 \\ | & 0 \\ | & 0 \\ | & 0 \\ | & 0 \end{pmatrix}$$

in every column

$$X' = \begin{pmatrix} | & 0 \\ | & 0 \\ | & 0 \\ | & 0 \\ | & 0 \end{pmatrix}$$

for this column.

→ better landscape of the error

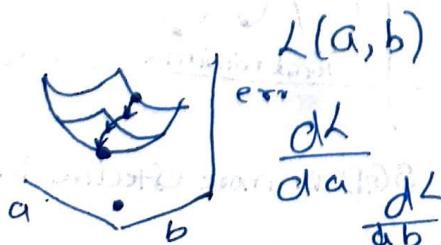
easier for gradient descent.

If we don't have info a column is not that important

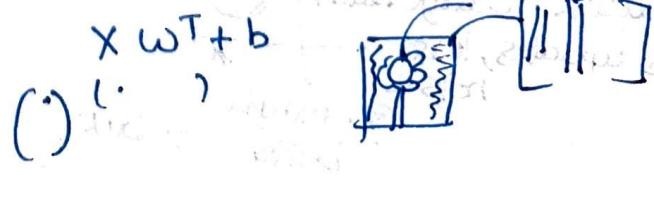
then scale it

make it more fair.

and no value



$$\frac{\partial L}{\partial a}, \frac{\partial L}{\partial b}$$



## Evaluating a regression model

$$y = w_0 + b$$

$$L(w, b) = \sum_{i=1}^n e_i^2$$

$X, Y$

$$X' = 2x, Y' = 2y$$

below all points by  
factor of 2

$$\text{Lnew}(w, b) = 4L(w, b)$$

} dependent on  
scale of data



\* MSE



$$MSE = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

↓ do not want  $(\hat{y}_i - y_i)$  being 0, when more spread.

( $\hat{y}_i$ ) of the data set.

Base line model :

$$\hat{y} = \text{mean}(y_i)$$

MSE base

average of existing points

no new calculations

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \text{mean}(y_i))^2$$

$$\text{Coefficient of Determination : } R^2 = 1 - \frac{MSE_{\text{model}}}{MSE_{\text{base}}} \quad (\text{new model})$$

(already existing)  
(mean of values)

Model not useful

} Consider  $MSE_{\text{model}} = 0$

} Perfect Case.

$$R^2 = 1$$

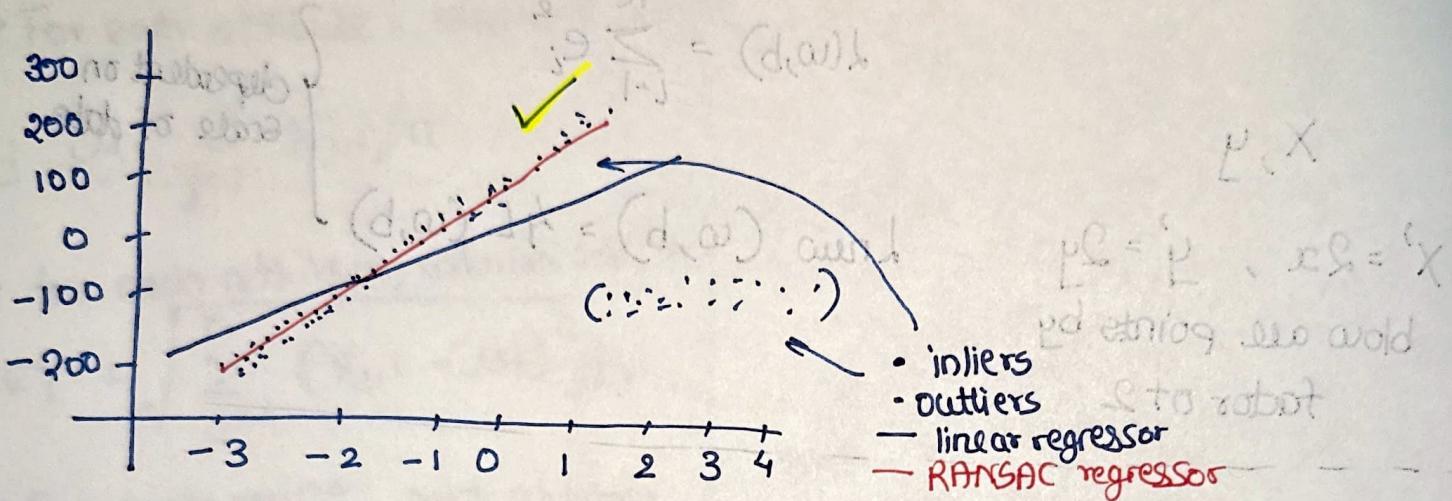
$$R^2 = 0$$

} if model Sucks  
 $\text{model} = \text{base}$

$MSE_{\text{model}} = MSE_{\text{base}}$   
worst-case

# Dealing with outliers RANSAC

\* random sample consensus



## Hyperparameters

$n'$  → number of sampled points

$\delta$  → inlier threshold

$n''$  → sufficient no of inliers

- Sample randomly  $n'$  points  $(x', y')$  from dataset  $(x, y)$
- fit a model  $f$  to  $(x', y')$
- Evaluate all points in  $(x, y)$  with  $f$  and keep only those that give an loss less than a threshold  $\delta$  These are inliers
- If number of inliers exceed than  $n'' \rightarrow$  model is good.
- If model is good then return new model fit on the inliers.

- Start with a model
- figure out threshold
- if there are enough below threshold
- add to inliers
- fit this new model.

## \* Feature Engineering

→ Meaningful features

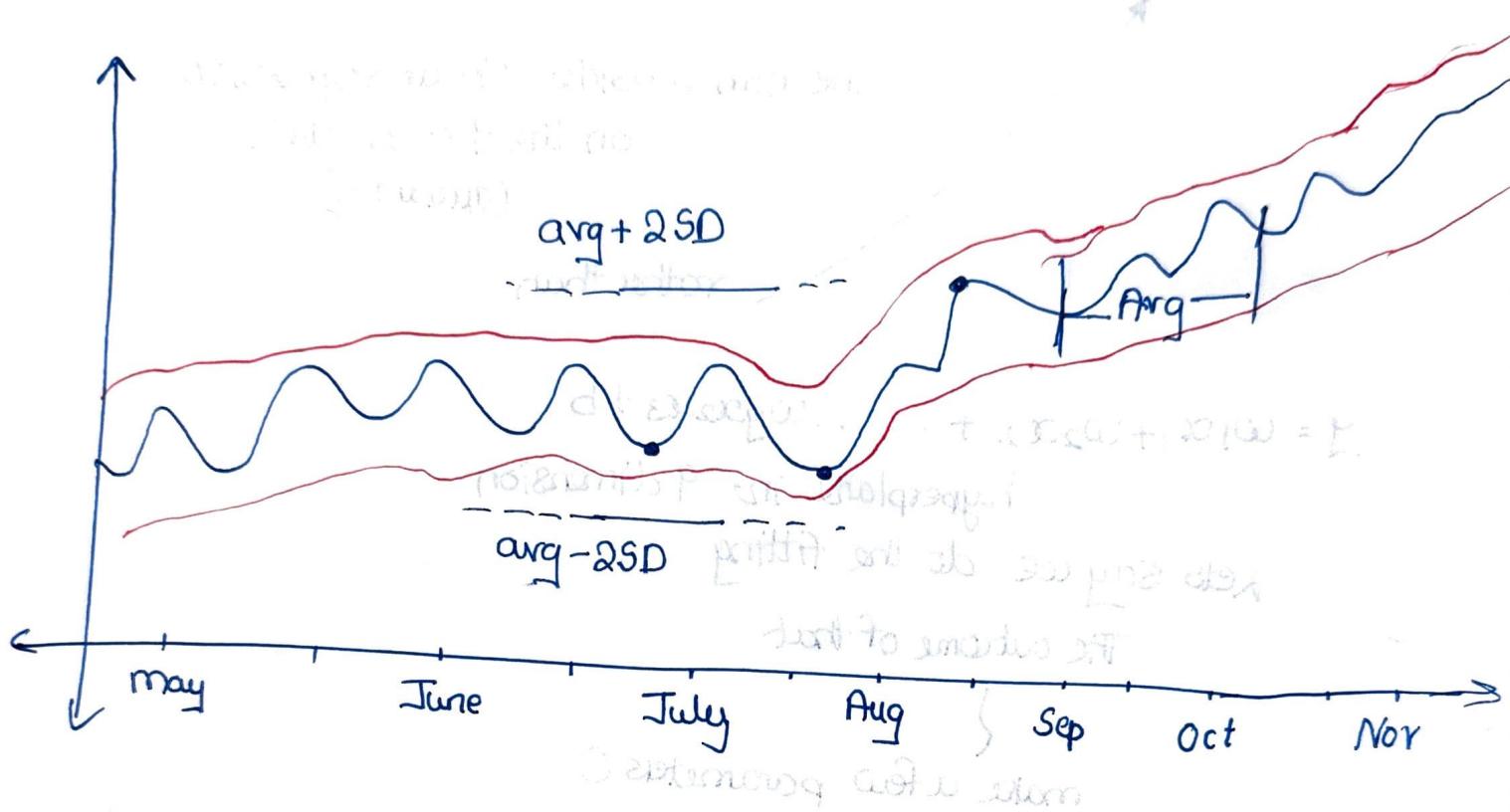
S&P 500

Open, High, Low, Close, Volume, Adj Close

Close price

Close price + 2SD

Close price - 2SD



$$\text{Definition of feature } X = [x_1, x_2, \dots, x_n]$$

For feature we will have  $\text{mean}(X)$  and  $\text{std dev}(X)$

$\text{std dev}(X)$

$2 \times \text{std dev}$

for every point plot  $\rightarrow$   $\text{Avg} + 2\text{SD}$

$\text{Avg} - 2\text{SD}$

$2 \text{std dev}(X)$

$$\sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$

at various point when it hits low

it seems to have tendency to bounce up as per SD

## Polynomial Regression :-

$$X = [x_1, x_2, x_3] \rightarrow [x_1, x_2, x_3, x_1^2, x_2^2, x_3^2, x_1x_2, x_2x_3, x_1x_3]$$

raw features    quadratic

we can consider linear regression  
on the 9 coeff class  
(quadratic)

~~add. p.s.~~  
rather than

$$y = w_1x_1 + w_2x_2 + \dots + w_9x_9 + b$$

hyperplane in 9 dimension

Let's say we do the fitting ~~add. p.s.~~

The outcome of that

make a few parameters 0

$$y = x_1^2 + x_2^2 \quad (\text{other coeff too small/0})$$

but this can is not linear  $\rightarrow$  it's quadratic

3D (raw space)  
not linear  
no hyperplane

it's just a surface.

We can also define cubic

Possible Applications :- amount of COVID cases

$x_1, x_2, \dots, x_n$  (following a increase)  $\rightarrow \underbrace{\log x_1, \log x_2, \dots, \log x_n}_{\text{log regression}}$

Covid Data  $R^2: 0.99$

actual for COVID from ~~other~~ sources  
So analysts thought it was made up.

(not necessarily true)