

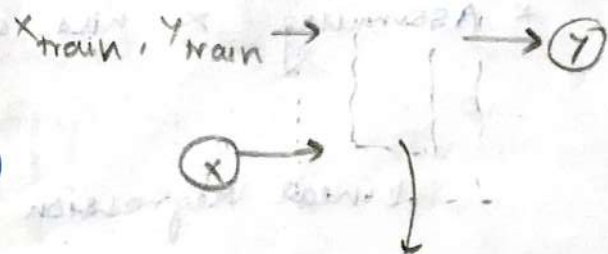
(\*\*) Linear Regression  $\rightarrow$  useful while predicting a continuous spectrum of output.

Now, given an  $x$ , we want the algo to predict;  
the corresponding  $y$

$x_{\text{train}}$  and  $y_{\text{train}}$  will be used  
by the machine, to come up  
(mapping) with a function (best possible)  
so that future prediction or  
test predictions can be as  
accurate as possible

(our main idea)

$\downarrow$   
fitting



a function b/w  
 $x$  and  $y$ , using  
which it will  
come up with  
values for future  
case of  $x$ .

(mapping)

## \* Function / Hypothesis

that this function fits  
our data / problem

(Linear regression) tries to put a linear function  
out there. i.e it assumes linear relation

b/w  $x$  and  $y$

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + b$$

output

dependent on  
each feature  
linearly

$m_n$  showcases the  
dependency of  $y$  on  
the  $n$ th feature

this dependency  
can be found  
out by varying  
 $m_1, m_2$  and  $m_3$ .

different features

The dependency of  $y$  on  
different features can be  
maintained using the values  
of  $m_1, m_2$  and  $m_3$

The 4 parameters  $m_1, m_2, m_3$  and  $b$   
are to be found out to reach  $y$  from  $x$ ,  
that is what the training data is  
for. (main objective is to reach the  
eqn. for that we need the 4 parameters)

\* Assuming  $x$  has only one feature (simplicity sake)

$x \leftrightarrow y$

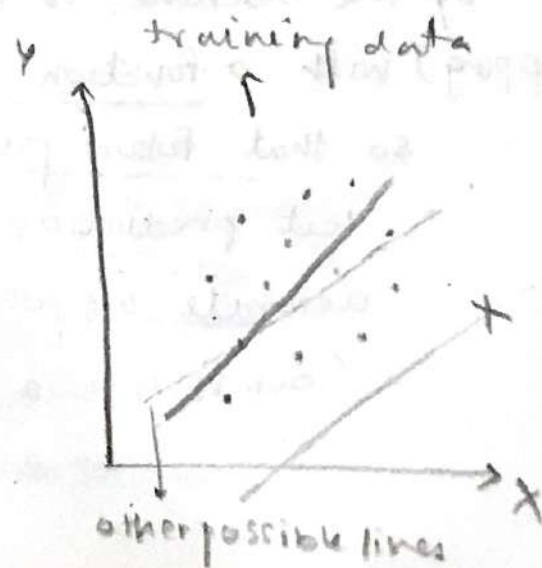
Linear Regression will try and  
fit it with one line, such  
that it is near to all points

i.e. combined error is

minimum

line selected to  
minimize error

how far we are  
from the line



\* To find the best fit line, we need to find  $m$  and  $c$  ( $c$  for 1 feature)

for  $n$  features:  $n+1$  parameters (including  $m_0, m_1, \dots, m_n, c$ )



$$y = \underbrace{m}_\text{slope} x + \underbrace{c}_\text{intercept}$$

\* given the best fit line is found for training set, we will just put it in the eqn. and find the

$y_{\text{pred}}$  for  $x_{\text{test}}$

i.e.  $y' = m x' + c$

(found using training data)



## ⑥ \* Coefficient of determination

↓  
To see how well the algo. is performing, we were just drawing a line and calculating how far the datapoints are from the line

↓  
(not the best approach)

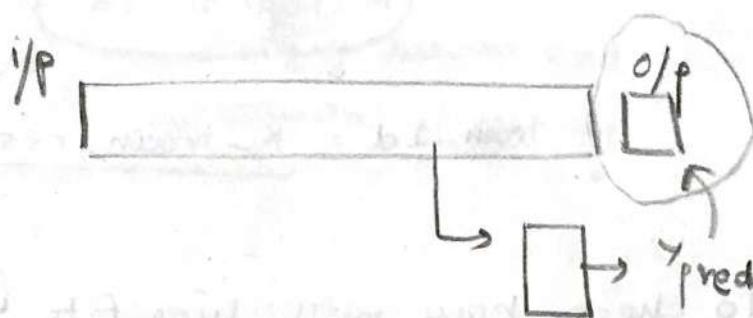
↓  
 $n$  features  $\rightarrow$   $n$  variables

↓  
 $n$  dimensions

↓  
(impossible to plot)  $\leftarrow$  complexity of graphs  $\uparrow$

\* objective way to find performance of algo.

way to compare the prediction and output and come up with a score



↓  
i) important to compare different algorithms

ii) for  $n$  feature dataset

2	3	4	Truth (actual outputs)
↕	↕	↕	
1.8	2.6	3.5	Prediction

Scoring  $\rightarrow$  Coefficient of determination

$$1 - \frac{u}{v}$$

$$U = \sum (y_i^T - y_i^P)^2$$

$$V = \sum (y_i^T - \underbrace{y^T_{\text{mean}}})^2$$

$$\text{COD} = 1 - \frac{U}{V}$$

$$= 1 - \frac{\sum (y_i^T - y_i^P)^2}{\sum (y_i^T - y^T_{\text{mean}})^2}$$

\* Higher the score, better the algorithm. (b/w 0 and 1 preferably)

$U \rightarrow$  sum of (errors squared)

$V \rightarrow$  assuming the ~~error~~ <sup>prediction</sup> is mean for each datapoint, then sum of error squared.

i.e. we assume mean can be the worst value, we can predict.

for each  $x_i \rightarrow y^T_{\text{mean}} \rightarrow$  predicted value  $= y^T_{\text{mean}}$

and if all  $y_i^T = y^T_{\text{mean}}$ , then  $\text{COD} = 1 - \frac{0}{0} = 1 - 1 = 0$ .

LOL :D \* if we make predictions worse than mean, then COD can be -ve

\* if we predict all true values, then  $\text{COD} = 1 - \frac{0}{0}$   
 $\downarrow$   
 numerator = 0  $\rightarrow$  = 1 (upperband)

(64)

\* already implemented in sci-kit learn by  
the name of score

↓  
alg 1. score (X-test, Y-test)

↓  
takes out Y-prediction and  
compares with Y-test to  
give us the score

alg 2. score (X-train, Y-train)

\*) The score value is judged based on the data and  
problem at hand. In some cases 0.7 might be a great  
score and in some cases 0.9 may also not be  
enough.

\* algorithm may perform better on testing data, when  
compared with training data.

\* random split may split differently in the next case

↓  
(so the scores may also)  
vary.



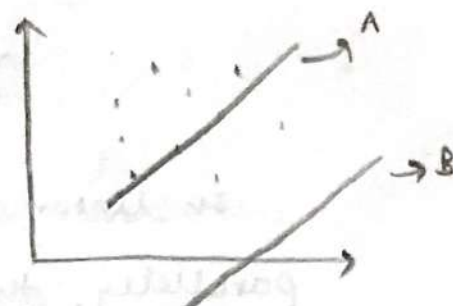
# \* Mathematical model behind linear regression

(65)

on how it finds the best suited line

↓  
an objective approach to find the best

↓  
(Error function)



$(x_i, y_i) \leftarrow i^{th} \text{ data point}$

input vector

$[x_1^i, x_2^i, x_3^i, \dots, x_n^i]$

↓  
n features  
related variables

↓  
the line predicted

$$\rightarrow m_1 x_1^i + m_2 x_2^i + m_3 x_3^i + \dots + m_n x_n^i + c$$

feature

$i \rightarrow$  not exponent, but  $i^{th}$  training example

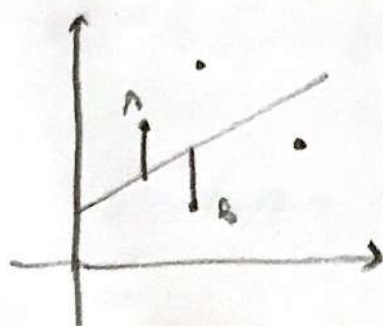
(Error :  
for certain  
training  
point)

$$y_i - (m_1 x_1^i + m_2 x_2^i + \dots + c)$$

Summation over entire dataset

Error for whole  
training set

$$: \sum_{i=1}^N (y_i - (m_1 x_1^i + m_2 x_2^i + \dots + c))$$

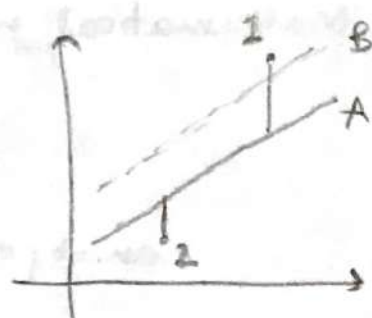


In this case, there will be a problem. that A and B will have diff. signs, hence they will cancel each other out.  
(not desirable).

(66)

$$\sum_{i=1}^N (y_i - (mx_i + c))$$

$$\sum_{i=1}^N |y_i - (mx_i + c)|$$



In linear case, when we move 2 lines parallelly, then the sum of errors will remain the same, as sum of distances remains the same (example point 1 and 2 from A and B)

$$\sum_i (y_i - (mx_i + c))^2$$

→ for 1D.

error function/  
cost function

in this case for lines A and B, A will be suited more, since for B, the distance from point 1 is high.

A is more suited.

↓  
We are punishing more for more distance by squaring

for n dimension,

$$\sum_i (y_i - (m_1(x_1)^i + m_2(x_2)^i + m_3x_3^i + \dots + m_n x_n^i + c))^2$$

feature for  $i$ th data point



\* current calculations only for 1 feature

(67)

↓  
 { came up with an intuition  
 for n dimensions later }

$$\text{cost} = \sum_i (y_i - (mx_i + c))^2$$

↓  
 we need to find the line for which the  
 cost function is minimum.

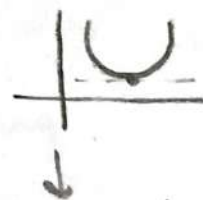
↓  
 we need m and c for that line

↓  
 find suitable m and c, so that  
 cost is minimized.

$$\text{cost}(m, c) = \sum_i (y_i - (mx_i + c))^2$$

⌘ quadratic  
 ↓  
 parabolic

$$\frac{\partial(\text{cost})}{\partial m} = 0 \quad \bigg| \quad \frac{\partial(\text{cost})}{\partial c} = 0$$



↓  
 corresponding  
 m and c for  
 cost function  
 to be minimum

↙ ↘  
 2 variables, 2 eqns

↓  
 m and c

$$\times \text{cost} = \sum_i (y_i - (mx_i + c))^2$$

$$\frac{\partial(\text{cost})}{\partial m} = \sum_i 2(y_i - (mx_i + c))(-x_i) = 0$$

$$= \sum_i y_i x_i - m \sum_i x_i^2 - c \sum_i x_i = 0$$

$$= \sum_i \frac{y_i x_i}{N} - m \sum_i \frac{x_i^2}{N} - c \sum_i \frac{x_i}{N} = 0 \quad (1)$$

$$\left( \frac{\sum x_i}{N} \right) = \text{mean of } x = x.\text{mean}()$$

(68)

$$\frac{\sum x_i^2}{N} = (x^* x) \cdot \text{mean}()$$

to get  
to mean

$$\frac{\sum x_i y_i}{N} = (x^* y) \cdot \text{mean}()$$

$$x \text{ cost} = \sum_i (y_i - (mx_i + c))^2$$

$$\frac{\partial \text{cost}}{\partial c} = \sum_i 2 (y_i - (mx_i + c)) = 0$$

$$= \sum_i y_i - m \sum_i x_i - \underbrace{nc}_{\substack{\text{for complete set} \\ \downarrow \\ \text{same across} \\ \text{all inputs}}} = 0$$

$$\Rightarrow \sum_i y_i - m \sum_i x_i = nc$$

$$\Rightarrow \sum_i \frac{y_i}{N} - m \sum_i \frac{x_i}{N} = c$$

$$c = y \cdot \text{mean}() - m^* x \cdot \text{mean}()$$

$$\frac{\partial \text{cost}}{\partial m} = x^* y \cdot \text{mean}() - m (x^* x) \cdot \text{mean}() - c^* x \cdot \text{mean}()$$

$$\Rightarrow x^* y \cdot \text{mean}() - m (x^* x \cdot \text{mean}()) - c^* x \cdot \text{mean}() = 0$$

$$\Rightarrow x^* y \cdot \text{mean}() - m (x^* x \cdot \text{mean}()) - (y \cdot \text{mean}() - m^* x \cdot \text{mean}())^* x \cdot \text{mean}() = 0$$

$$\Rightarrow (x \cdot y). \text{mean}() - m \cdot ((x^2 \cdot x). \text{mean}())$$

$$- (y \cdot \text{mean}() \cdot x \cdot \text{mean}())$$

$$- m^2 \cdot x \cdot \text{mean}() \cdot x \cdot \text{mean}() = 0$$

$$\Rightarrow \frac{(x \cdot y). \text{mean}() - (y \cdot \text{mean}() \cdot x \cdot \text{mean}())}{(x^2 \cdot x). \text{mean}() - x \cdot \text{mean}() \cdot x \cdot \text{mean}()}$$

$$(x^2 \cdot x). \text{mean}() - x \cdot \text{mean}() \cdot x \cdot \text{mean}()$$

$$\Rightarrow y \cdot \text{mean}() - m \cdot x \cdot \text{mean}() = c$$

These eqns give us the coefficients to find our best fitting line for our data.