

## understanding the algorithms

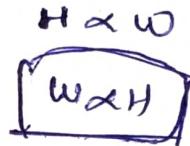
### ① Simple Linear Regression

Linear model in Regression  
(supervised learning)

Simple linear regression assumes that a linear relationship exists b/w the response variable & the explanatory variable, it models this relationship w/ a linear surface called a "hyperplane".

→ This will be clear with an ex.  
lets take the height & weight sample of few people

height	weight
5.6	60
5.7	62
5.8	63
5.9	62
6.1	64
6.2	75



Now, if I want to know what would be the weight of a person with 6.0 Height?

→ By observation or second like

$$H \propto W$$

→ to remove proportionality

we can multiply a constant

$$W = mH$$

→ But, it may also differ by some value, so,  $+ b$

$$W = mH + b$$

→ Now, this is the relation of height & weight.

constants

$$w = mH + b$$

given height

→ now, here if I know the value of 'm' & 'b'

then I can predict the weight.

Now, main aim is to find out

value of 'm' & 'b'.

→ Let's take a sample  $(5.6, 60), (5.7, 62)$ ,

$$60 = m(5.6) + b \quad \text{--- (1)}$$

$$62 = m(5.7) + b \quad \text{--- (2)}$$

$$m = 20, b = -52.$$

So, our equation

$$w = 20H - 52$$

Let's try with a "known" height 5.9.

$$w = 20(5.9) - 52 = 66 \text{ kg.}$$

Here we got '66', but we are expecting 62.

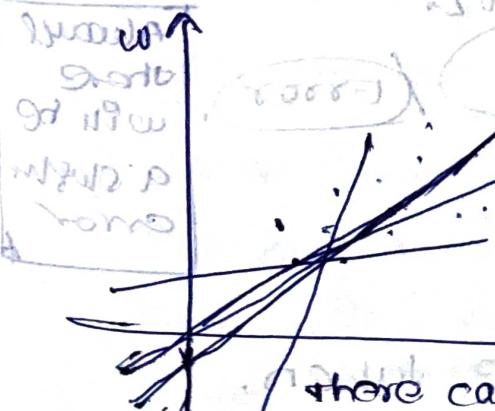
It may be that we change  $[m \& b]$  then we could get an exact value (approx).

Now, here we get the definition of ML as establishing a mathematical relationship or data, so that it can predict a new data.

→ So, when we say a model,

\* It is nothing but a mathematical relationship which can predict new data.

Plotting  $H$  &  $w$



Let's assume we have a best fit line.

Linear regression always tries to find the linear relation bw variable which is a straight line.

$$\text{Hence, } y = mx + c$$

there can be  
in lines.

→ we could find  $n$  diff lines.

But that one line which can predict all the value appears is the "Best Fit line".

$y = mx + c$  → and this best fit line is called model.

→ we have to find the ' $m$ ' & ' $c$ '.

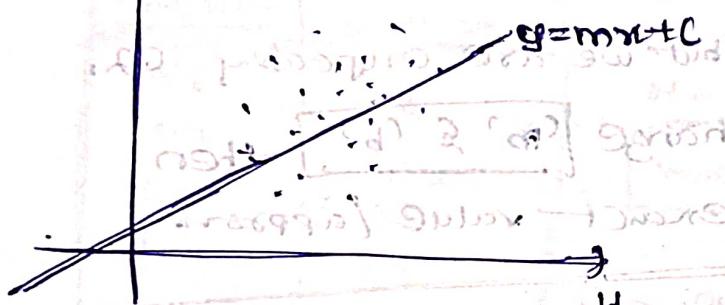
→  $m$ -slope,  $c$ -Intercept

\* Evaluating the fitness of the model with a

cost function

→ First, let us assume we have a best fit line.

But, now we have doubt.



How did we opt this line?  
Why this line?  
What about other lines?

→ Now, if it's the best fit line, then posture  
we want the weight or '0.0'.

→ It is nothing but the value of  $w$ -coordination

$y = mx + c$  with  $x$  coordinate 6.0

→ This line is  $y = 20H - 52$  and at  $x = 5.9$ ,

I got  $w \approx 66$ , but actually  $w = 62$ .

So, I can see there is a known

loss

Error

of  $10H = |(y - \hat{y})|$  Predicted.

Q) Why did we get this loss?

A) may be best fit is not taken,

Always  
there  
will be  
a slight  
error

## Gradient Descent

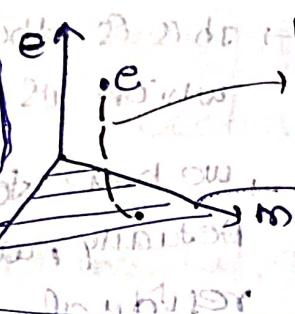
→ Now, do I improvise my model Strength  
I have to reduce the loss.

and this loss depends on the "c" taken.

→ So, I am supposed to shift m & c such that I end up reducing the loss.

we have loss (e), slope (m), constant (c).

All I have to do is bring them down to my plane



should bring down & find the new coordinate

In this plane we have no error.

④ slowly, by examining each input m & c get changed and eventually loss decreased. this process is called learning.

so,

$$m_{\text{new}} = m_{\text{old}} - \eta \Delta m$$

$$c_{\text{new}} = c_{\text{old}} - \eta \Delta c$$

$m_{\text{new}}$  - changed one, new

$m_{\text{old}}$  - old m value.

$\eta$  - learning rate

$$\eta \rightarrow (0.0001 - 10) \text{ range.}$$

$\Delta m$  - Error in m / de

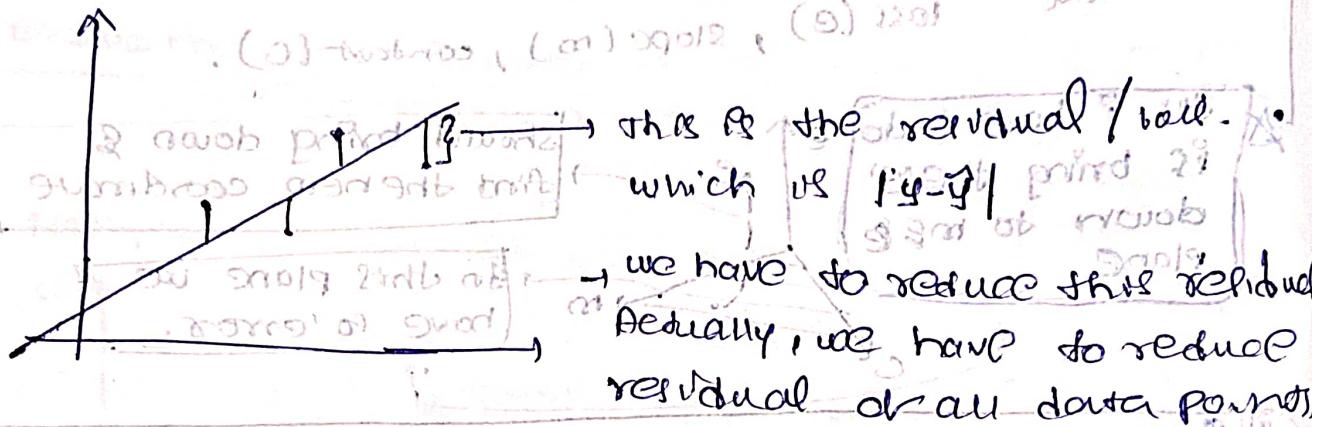
→ So, here we are taking a new m.

By subtracting a fraction of the error rate from old m that fraction is learning rate.

so, now I am supposed to choose the m & c, which give a least possible error. Then it will be the best fit line.

→ A cost function, also called a loss func., is used to define & measure the error of a model. The difference between the observed weights in training set are the difference between the weights predicted by model.

If the observed weights in training set are called "residuals", or training error / loss.



→ we can produce best weight predictor model by minimizing the sum of residuals.

$$(\text{Sum of residuals}) \quad E_n = \sum_{i=1}^n [y_i - (m_{\text{pred}} + c)]$$

so, to cancel the negative value, we take square

$$R(x) = \sum_{i=1}^n y_i^2 = \sum [y_i - (m_{\text{pred}} + c)]^2$$

sum of squares of residuals

$(e, m, b)$

Path it travel  
is gradient descent

→ so, we are supposed to shift the  $(e, m, b)$  to the origin  $(0, m, b)$ .

where  $\frac{de}{dm} = 0, \frac{de}{db} = 0$ .

As we can see residual is both a function of  $m$  &  $b$ , so differentiating partially with respect to  $m$  &  $b$  will give us:

$$\frac{\partial R}{\partial m} = \sum_{i=0}^n 2m_i(b + mx_i - y_i)$$

$$\frac{\partial R}{\partial b} = \sum_{i=0}^n 2(x_i(b + mx_i - y_i))$$

$m = 0$

so, we know for best line, residual should be min. minima of func. occurs where derivative = 0.

i.e.

$$\sum_{i=0}^n 2m_i(b + mx_i - y_i) = 0$$

$$\sum_{i=0}^n 2(x_i(b + mx_i - y_i)) = 0$$

$$\sum_{i=0}^n m_i + \sum_{i=0}^n m x_i^2 + - \sum_{i=0}^n y_i x_i = 0$$

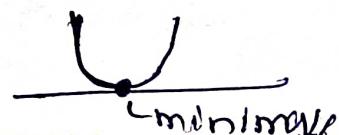
$$\sum_{i=0}^n c + \sum_{i=0}^n m x_i - \sum_{i=0}^n y_i = 0$$

$$\sum_{i=0}^n (c + mx_i) = \sum_{i=0}^n y_i$$

The same eq. can be written in matrix form as:

$$\begin{bmatrix} \sum_{i=0}^n m_i & \sum_{i=0}^n m^2 \\ \sum_{i=0}^n m & \sum_{i=0}^n m x_i \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix} = \begin{bmatrix} \sum_{i=0}^n m_i y_i \\ \sum_{i=0}^n y_i \end{bmatrix}$$

ideally, if we have an eq. or one dependent & one independent variable the minima will look like



dependent  
var  $y = m_0 + c$  independent var

→ the new values for slope & intercept  
are calculated.

repeat until convergence

$$m_{\text{new}} = m_{\text{old}} - \eta \left[ \sum_{i=1}^n (y_i - \hat{y}_i) x_{ij} \right] \quad \frac{\partial G}{\partial m}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \left[ \sum_{i=1}^n (y_i - \hat{y}_i) \right] \quad \frac{\partial G}{\partial b}$$

this is how LR works, Accuracy

\* measuring model strength RSE - residual sum of squares

the  $R^2$ -squared ( $R^2$ ) statistic provides a measure of fit.

It takes the form of proportion - (proportion explained).

$$R^2 = 1 - \frac{SSE}{TSS}$$

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

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

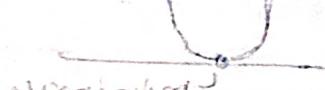
$$R^2 = 0.75$$

It says that our model.

fits 75% of data set.

Show that  $\hat{y}_i$  is the best linear fit.

new concentration



## Adjusted R<sup>2</sup> Statistic

Adjusted R<sup>2</sup> Statistic

→ As R<sup>2</sup> is just a linear eq., as we increase the number of independent variables in our eq., the R<sup>2</sup> increases as well.

But, that doesn't mean that the new independent variables have any correlation with the output variable.  
i.e., R<sup>2</sup> will increase, but it is not necessarily model yields better results.

→ To rectify this we use adjusted R<sup>2</sup> value which penalizes excessive use of such predictors which do not correlate with the output data.

$$R^2_{adj} = 1 - \frac{(1-R^2)(N-1)}{N-P-1}$$

P - no. of predictors  
N - total sample size.

If P=0,  $R^2_{adj} = R^2$ .

Note!

→ using training data to learn the values of the parameters for simple linear regression that produce the best fitting model is called ordinary least square (OLS) or linear least square.

↓↓↓↓↓

Another approach to find m & C.

Substitute for b

(to solve eqs)

Diff between two lines until a - b = 0

→ variance is a measure of how far apart  
the values all spread out

→ covariance is a measure of how much  
two variables change together. If the  
variables increase together, their cov is positive.  
o no relation, -ve → one increase & one decrease

$$\text{var} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$\text{cov} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

→ variance of explanatory variable.

→ covariance of the response, explanatory variables.

$$\beta = \frac{\text{cov}(x, y)}{\text{var}(x)}$$

$$y = \beta x + \alpha$$

$$\alpha = \bar{y} - \beta \bar{x}$$

Note:

- the independent variables are uncorrelated with the residual term, also known as "Exogeneity".  
This in layman terms generalizes that if no way should the error term be predicted given the value of independent variables.

→ The error terms have a constant variance, which ~~is true~~

### Homoscedasticity

→ The error terms are normally distributed.

→ No multicollinearity, i.e., no independent variables

Should be correlated with each other or

affect another. ~~so that they are not correlated with each other~~

### Model Confidence

⇒ If linear regression is a low bias / high variance model

(or) a high bias / low variance model?

↳ It's a "high bias / low variance" model.

→ Even after repeated sampling, the best-fit line won't stay roughly in the same pos. (low variance)

→ But, the avg'd the models created after repeated sampling won't do a great job in capturing the posted relation ship. (high bias)

→ Low variance is helpful when we don't have less training data.

```
import statsmodels.formula.api as smf
model_name = smf.OLS(formula=' ~ b1, b2, b3')
model_name.conf_int()
```

→ Conf. Int. are more than 100%

→ P-values < 0.05