

Linear Regression

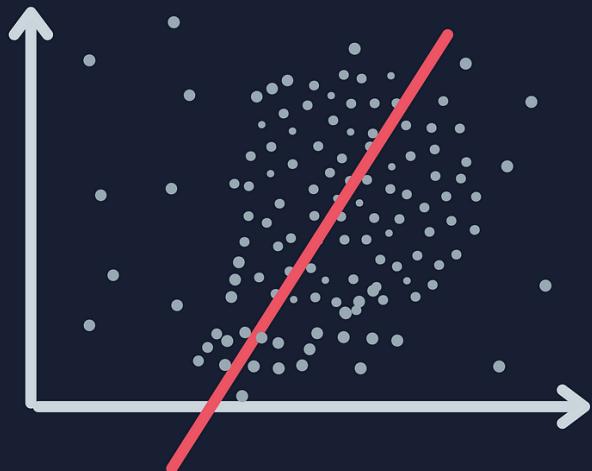


Table of Contents

1. What is Linear Regression
2. Understanding with an example
3. Evaluating the fitness of the model
4. Understanding Gradient descent
5. Understanding Loss Function
6. Measuring Model Strength
7. Another Approach for LR - OLS

understanding the algorithms

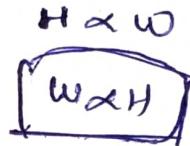
① Simple Linear Regression

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

$w = mH + b$ → 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).

$$\begin{aligned} 60 &= m(5.6) + b \quad \text{--- (1)} \\ 62 &= m(5.7) + b \quad \text{--- (2)} \end{aligned} \quad \left. \begin{array}{l} \text{By this we get} \\ m = 20, b = -52. \end{array} \right.$$

So, our equation $w = 20H - 52$

lets say with a known height 5.9.

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

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

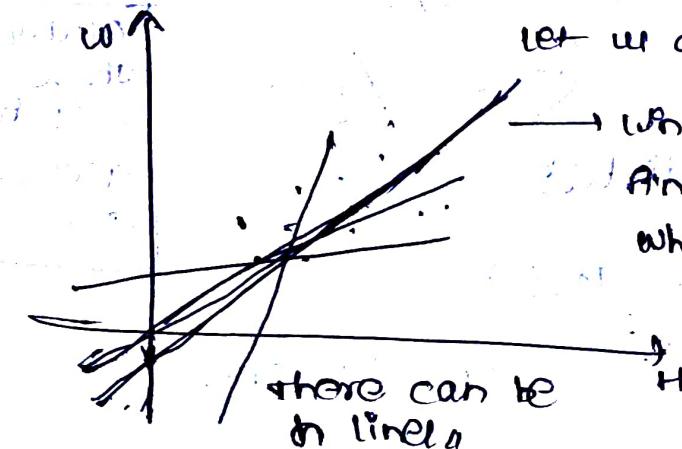
④ It may be or we change ['m' & 'b'] then we could get an exact value (approx.)

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,

④ It is nothing but a mathematical relationship which can predict new data.

Plotting 'H' & 'w'



Let us assume we have a best fit line.

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

Hence, $y = mx + c$

→ we could find n diff lines.

But that one line which can predict

all the values approx. is the "Best fit line"

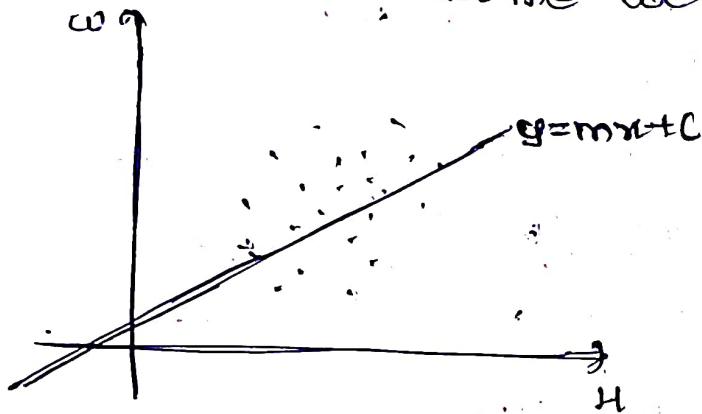
$y = mx + c$ → And this best fit line is called model
we have to find the 'm' & 'c'.

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
it? we want the weight or 'w'.

→ It is nothing but the value of y-coordinates
 $y = mx + c$ where worth or coordinate '6.0'

→ If this line is $w = 20H - 52$ and at $n = 5.9$,

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

So, I can see there is a known loss

$loss = |(y - \hat{y})|$ predicted

Q) why did we get this loss?

A may be best fit is not taken.

Always
there
will be
a slight
error

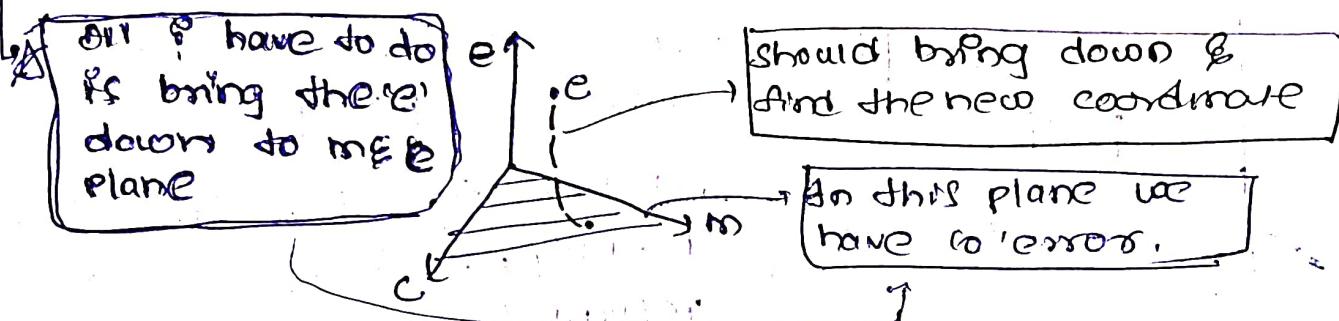
Gradient Descent

→ Now, to improve my model strength,
I have to reduce the loss.

and this loss depends on the "c" taken.

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

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



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_{new} - changed one, new

m_{old} - old m value.

η - learning rate

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

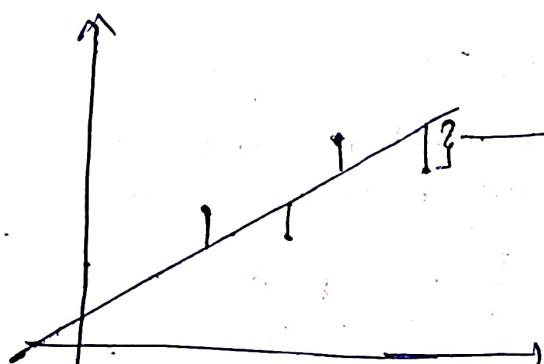
Δm - Error in m / m_{de}

→ So, here we are taking a new m .

By subtracting a fraction of the error rate from old 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 diff. b/w the weights predicted by model & the observed weights in training set are called "residuals", or training error / loss.



→ This is the residual / loss which is $|y_i - \hat{y}_i|$

→ we have to reduce this residual
 Actually, we have to reduce residual of all data points.

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

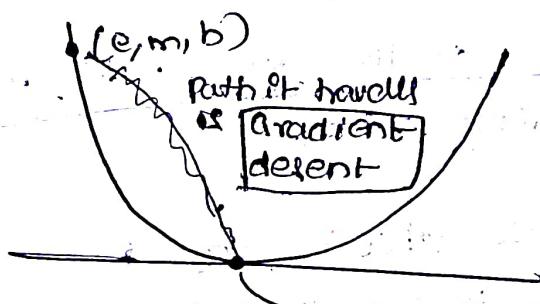
$$(\text{Sum of residuals}) \quad \sum_{i=1}^n |y_i - (m x_i + c)|$$

so, to cancel the negative value, we take square

$$R(x) = \sum_{i=1}^n y_i^2 = \sum_{i=1}^n [y_i - (m x_i + c)]^2$$

sum of squares
of residuals

Now, it is a quadrature in 3D.



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

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

As we can see residual is both a func. of m & b , so differentiating partially w.r.t m & b will give us:

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

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

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

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

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

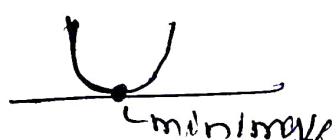
$$\boxed{\sum_{i=0}^n \alpha_i c + \sum_{i=0}^n \alpha_i m x_i^2 + -\sum_{i=0}^n \alpha_i y_i x_i = 0}$$

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

The same eq. can be written in matrix form as

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

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



repeat for $y = mx + c$ and update var

→ the new values for slope & intercept are calculated.

repeat until converge()

{

$$m_{\text{new}} = m_{\text{old}} - \eta \left[\sum_{i=1}^n (h_0 \cdot x_i - y_i) x_i \right] \quad \frac{\partial E}{\partial m}$$

$$c_{\text{new}} = c_{\text{old}} - \eta \left[\sum_{i=1}^n (h_0 \cdot x_i - y_i) \right] \quad \frac{\partial E}{\partial c}$$

Accuracy this is how LR works.

* measuring model strength RSS - residual sum of squares
 TSS - total sum of square

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

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

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

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{Product of } \text{Residual change}$$

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

BT $R^2 = 0.75$

It says that our model fits 75% of data set.

Adjusted R² Statistic

→ As R² is just a linear eq., as we increase the no. of independent variables in our eq., the R² increases as well.

But, that doesn't mean that the new independent variables have any correlation with the output variable.

i.e., R² will increase, but it is not necessarily model yields better results.

→ To rectify this, we use adjusted R² value which penalises excessive use of such features which don't 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 squares (OLS) or linear least square.

Another approach to find m & C.

(to solve α s)

- variance is a measure of how far a set of 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

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

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

→ variance of explanatory variable.

→ covariance of the response, explanatory variables.

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

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

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

thus cov α s zero

Note :-

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

the error terms have a constant variance.

Homoscedasticity

- the error terms are normally distributed.
- no multicollinearity, i.e., no independent variables should be correlated with each other or affect another.

model confidence

or is linear regression 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 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 p'st relationship.

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

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