

51 Gradient Descent from Scratch:

A first-order iterative optimization algorithm for finding a local minimum of a differentiable function.

→ considering m const. (0.01)

$$\beta_{\text{new}} = \beta_{\text{old}} - \underset{\substack{\uparrow \\ \text{learning} \\ \text{rate}}}{\eta} \cdot \underset{\substack{\uparrow \\ \frac{dL}{d\beta}}}{(\text{slope})}$$

no. of iterations = epochs

↑
when to stop.

$\eta \cdot \frac{dL}{d\beta} \rightarrow$ step size.

→ Now m & β are both unknown.

(1) initialize random values of m & β .
say, $m = 1$ $\beta = 0$

(2) decide epochs & learning rate.
say epochs = 100 $lr = 0.01$

→ Apply same formula for both m & β .

But now, L will depend on both m & β . Therefore, in both equation, slope will be partial derivative of $L(m, \beta)$

$$\beta_{\text{new}} = \beta_{\text{old}} - \eta \frac{dL}{d\beta}$$

$$m_{\text{new}} = m_{\text{old}} - \eta \frac{dL}{dm}$$

$$\frac{dL}{d\beta} = \frac{d}{d\beta} \left(\sum (y_i - m x_i - \beta)^2 \right)$$

$$\frac{dL}{d\beta} = -\frac{2}{n} \sum (y_i - m x_i - \beta) \quad \sum (y_i - \hat{y}_i)$$

and

$$\frac{dL}{dm} = -\frac{2}{n} \sum (y_i - m x_i - \beta) x_i$$

→ right learning rate is necessary to get optimum epochs, otherwise model will take too long to process.

52

Variations / variants of GD:

① Batch

↓

considering whole data as 1 batch

and the new values of m & β_0 are updated after batch processing

Stochastic

↓

updating after every row.

mini-batch

↓

considering whole data into some no. of mini-batches & updating every time each mini-batch gets processed.

→ for multiple linear regression with n inputs.

$$\frac{dL}{d\beta_0} = -\frac{2}{n} \sum (y_i - \hat{y}_i) \quad \frac{dL}{d\beta_j} = -\frac{2}{n} \sum (y_i - \hat{y}_i) x_{ij}$$

doing this summation in 1 go, we need to use dot product.

→ Normal gradient descent is
Batch gradient descent

- problem with batch GD:
 - As higher dim data with high epochs to run, there will be large no. of derivative calculation increasing the model time.
- Hardware problem due to memory overflow.

② Stochastic GD: (faster)

- In the epoch we update 1 row the basis per update hogya.

↳ minimize no. of epochs.

- I row ke use ke karan memory overflow nhi hoga.
- random row selection & update.
- Slightly in precise prediction due to randomness.

- Due to randomness, after reaching near to soln, it fluctuates. So to reduce this fluctuations, we use learning schedules varying learning rate with epochs

③

Mini-Batch Gradient Descent:

→ Creating batches in rows.

Say $n = 1000$

batches $\rightarrow 100$ — size = 10

So, every epoch we do 100 updates.

* Session 1-2 on Regression Analysis

↳ The process of studying the relationship between X and y using stats and with the help of computed coefficients.

• steps :-

- (1) Identify dependent and independent variables.
- (2) collect and prepare the data.
- (3) visualize the data
- ✓(4) Check assumptions
- (5) fit the LR model
- ✓(6) Interpret the model : Analyzing estimated regression coefficient, their std errors, t-values, p-values to determine statistical significance of relationship between dependent & independent variables.
- (7) Validate the model
- (8) Report results.

$$\begin{array}{ccccc} \text{dependent} & \nearrow & Y = f^*(X_1, X_2, \dots) & + & \epsilon \\ \text{variables.} & & \uparrow & & \uparrow \\ & & \text{independent} & & \text{irreducible error} \\ & & \text{variables} & & \end{array}$$

$f(X_1, X_2, X_3, \dots) \rightarrow \text{true fn.}$

But there is deviation in f^n due to sample (as we don't have whole population)
This too creates some error

i.e. $f() - f'()$

↓
reducible error.

True parameters $\rightarrow \beta_0 \beta_1 \beta_2 \dots$
Computed parameters $\rightarrow b_0 b_1 b_2 \dots$
based on sample

So, as we try to bring the values of b_i close to β_0 the error $f() - f'()$ reduces.

$$\Rightarrow Y = f'(X_1, X_2, \dots) + \text{reducible error} + \text{irreducible error}$$

↑
estimated f^n .
using model

xxxxx
 \rightarrow Now, we will use hypothesis testing to know whether the estimated parameters are useful or not

\rightarrow Then we will calculate confidence intervals to get range of values of parameters.

xxxxx

xxxxx

→ To study how much each feature (X_i) contribute to Y .

- prediction ^{→ what} ⇒ To give output of particular individual based on trained data.
- Inference ⇒ To give relation b/w features due to which we get the predicted value.
_{why.}

prediction $\propto \frac{1}{\text{inference}}$

Jiska prediction achi unka inference sahi nahi and vice versa.

- Statsmodel LR:

$$X \rightarrow Y$$

- ① Is there a relationship. (f-statistic)
- ② If yes, is it linear
- ③ If yes, how strong. (R^2 score)
- ④ If it is strong, then how is the relationship of every individual component of X with Y . (individual t-test)

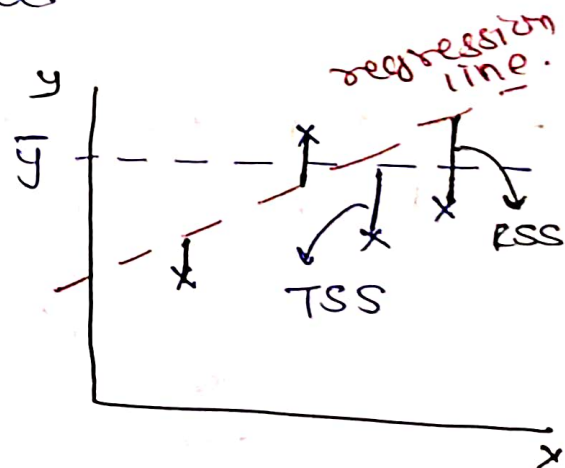
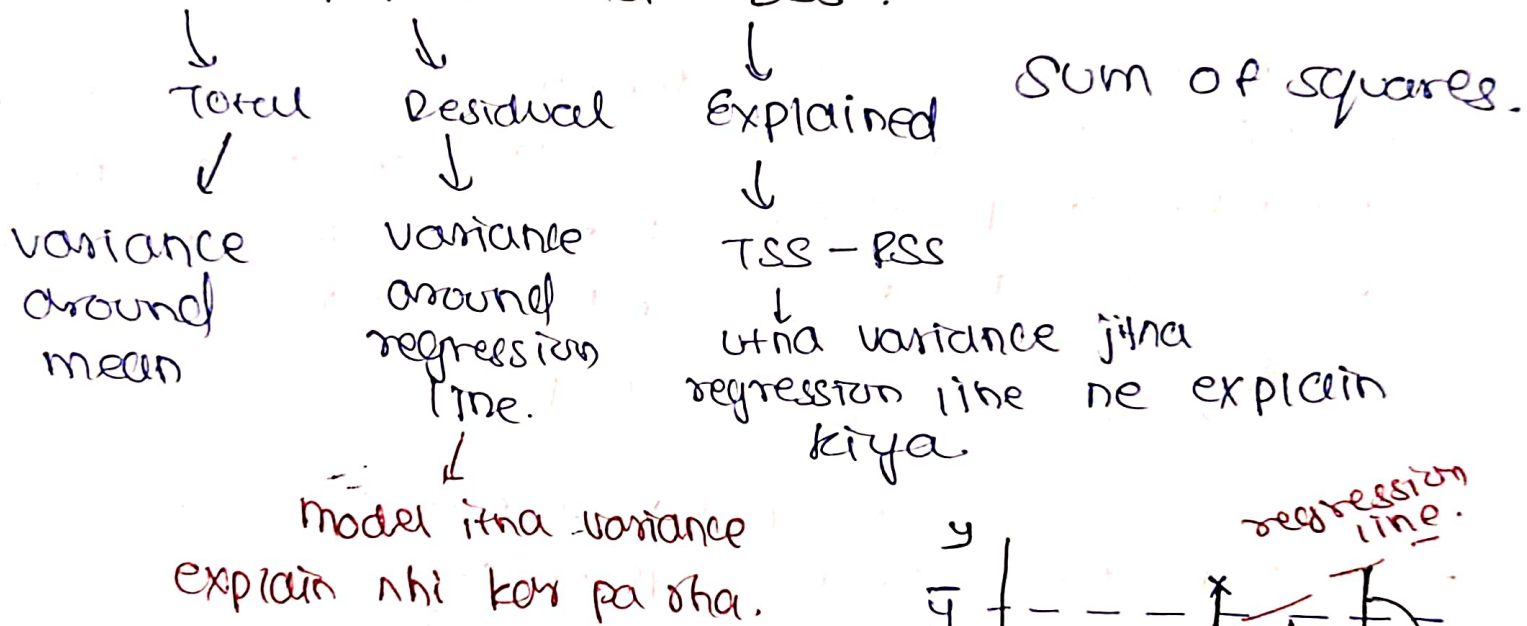
```
import statsmodels.api as sm
X = sm.add_constant(X) → for intercep.
model = sm.OLS(Y, X).fit()
print(model.summary())
```


• Summary :

- upper right corner \Rightarrow relationship b/w X & Y
- middle part \Rightarrow Individual relation b/w feature & y colm.
- bottom part \Rightarrow assumptions.

① relationship b/w X & Y

• TSS , RSS and ESS :



$$TSS = \sum (y_i - \bar{y})^2$$

$$RSS = \sum (y_i - \hat{y}_i)^2$$

$$ESS = TSS - RSS$$

due to both reducible & irreducible error.

• Degree of freedom.. $(n-1)$

$$df_{\text{total}} = df_{\text{model}} + df_{\text{residual}}$$

$$\begin{array}{ccc} & \checkmark & \downarrow \\ & \underbrace{k} & \underbrace{n-k-1} \\ \text{no. of input cols} & & \text{no. of rows} \\ \text{or} & & \\ \text{no. of independent} & & \end{array}$$

* f-statistic :-

= f-test used to determine whether a linear regression model is statistically significant, meaning it provides a better fit to data than just using mean of dependent variable.

↓
varient of ANOVA test.

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots$$

→ f-test for overall significance :-

① $H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$

$H_a : \beta_1 \neq \beta_2 \neq \dots = \beta_k \neq 0$

↓
at least one independent variable contributes.

② fit LR model to get estimated regression coefficients.

③ calculate TSS, ESS, RSS.

④ Compute \rightarrow MSR (Mean square regression)

$$MSR = \frac{ESS}{k}$$

\rightarrow MSE (Mean square Error)

$$MSE = \frac{RSS}{n-k-1}$$

⑤ calculate f-statistic : $\frac{MSR}{MSE}$

⑥ Determine p-value

⑦ Derive results by comparing p-value & α .

* R^2 score :- (Coefficient of determination)

\rightarrow quantifies the proportion of variance in dependent variable that can be explained by independent variables.

$R^2 \rightarrow 1$ (better fit of a model)

$$R^2 = \frac{ESS}{TSS}$$

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

$$\text{Adjusted } R^2 = 1 - \left[\frac{(1-R^2)(n-1)}{(n-k-1)} \right]$$

- R^2 always increases or stays same with addition of new predictor variables, regardless of those variables which do not contribute.
- Adjusted R^2 score penalizes the model for adding unnecessary complexity.

* T-test :-

- ① $H_0: \beta_k = 0$
 $H_a: \beta_k \neq 0$ Same for β_0
- ② Estimate β_0, β_k using LR model $\rightarrow b_0, b_1$
- ③ calculate std. errors for slope and intercept coefficients $[SE(\beta_0) \& SE(\beta_k)]$

$$SE(b_0) = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{(n-2)} \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right]}$$

$$SE(b_1) = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{(n-2) \sum (x_i - \bar{x})^2}}$$

- ④ Compute t-test statistic for b_0 & b_1

$$(t\text{-value})_{b_0} = \frac{b_0 - 0}{SE(b_0)}$$

$$(t\text{-value})_{b_1} = \frac{b_1 - 0}{SE(b_1)}$$

As taken
in
Null
hypothesis.

* confidence interval:

$$\left. \begin{aligned} CI_{b_0} &= b_0 \pm t\text{-value} * SE(b_0) \\ CI_{b_1} &= b_1 \pm t\text{-value} * SE(b_1) \end{aligned} \right\}$$

calculate using

$$(df = n - 2)$$

$$(\alpha)^{\frac{1}{2}}$$