

Linear Regression

A Look at supervised learning

Training a ml task for every input with a corresponding target, is called Supervised learning.

Supervised learning for price prediction of used cars

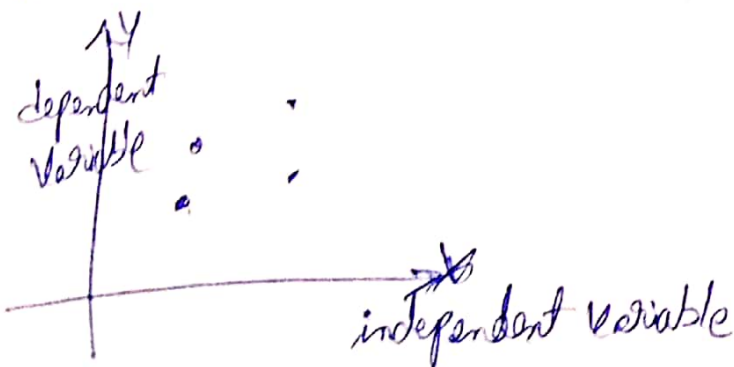
↳ predict the price of an used car.

↳ has predictor (independent) variables and an outcome (dependent) variable.

What is Regression

Predicts values of a Continuous dependent variable

Using independent explanatory variable(s)

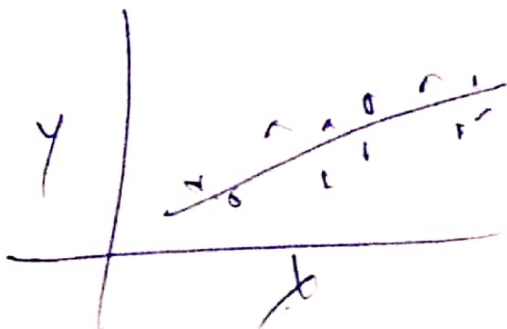


Regression Use Cases

- ① Real estate: model house prices as a function of home's living area, no of bedrooms, bathrooms & lot size.
- ② Marketing: model relationship b/w online advertising costs and monthly e-commerce sales.
- ③ Medicine: forecast medical expenses for insured population based on attributes like age, gender, smoker, BMI,
- ④ Weather: Model relationship b/w crop yield based on amount of rainfall.

What is linear regression?

A form of regression that models linear relationship b/w dependent & independent variables.



x & y are linearly related.

What makes Linear Regression Popular?

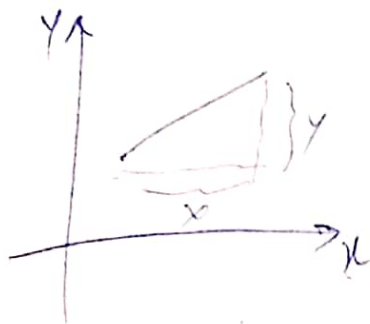
→ Ease of implementation =

→ Powerful =

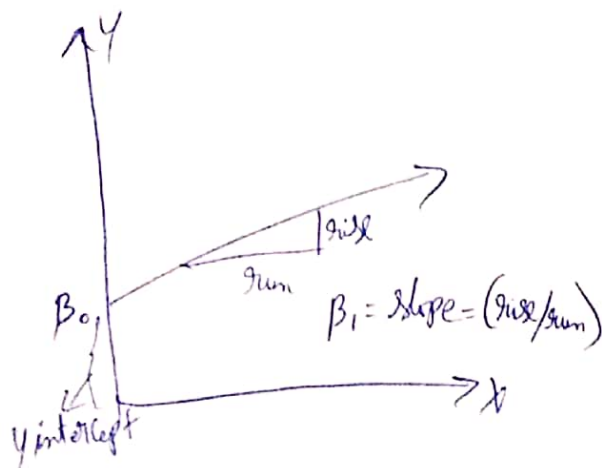
→ Easily Solvable =

ML - ~~not~~ needing more data = DL - require tons of Data.

Describe equation of a line

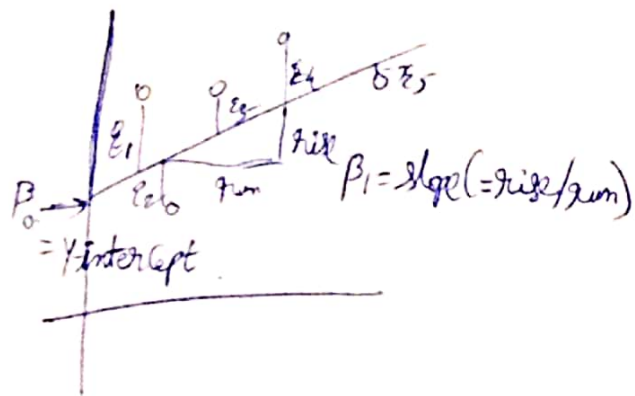


Dependent Variable \rightarrow Independent Variable
 $Y = mx + b$
Coefficient, rate & slope of line \rightarrow Y-intercept
Where line crosses the Y-axis



Predicted value \rightarrow Input variable
 $Y = \beta_1 X + \beta_0 + \epsilon$
Intercept \rightarrow Error term
Slope of regression line

Slope has great impact on your model (β_1)

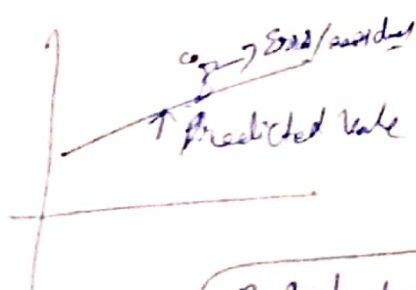


Errors are not directly added

$$\hat{y} = (\beta_1)x + \beta_0 + \epsilon$$

Slope of regression ~~line~~ line

Errors or residuals



>> Residual is the difference b/w observed value of dependent variable and predicted value.

Residual ~~var~~ = observed values - predicted values

Metrics Used to Evaluate L.R. model

Accuracy = Closeness of actual predicted \rightarrow as high as possible.

Loss function = Actual - Predicted (Error & Residual)
 as low as possible.

What is the best fitting line - Line of Best fit.

Goal: Find the best slope and intercept that fits data points
 Soln: Minimize residual errors.

The following metrics are used to measure model performance

MAE = Mean absolute Error.

MSE = Mean Squared Error.

RMSE = Root mean Squared Error.

MAE is the average absolute difference between actual and predicted values.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Characteristics

→ Neutral to outliers.

→ Error is in same units as that of the data points.

MSE is the average squared difference between actual and predicted values.

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

Characteristics :

→ Accounts impact of outliers.

→ Error is not in same unit as that of the data points

RMSE is the square root of the average squared difference between actual and predicted values.

$$RMSE = \sqrt{\frac{1}{n} \sum_{j=1}^n (y_j - \hat{y}_j)^2}$$

Characteristics :-

- Accounts impact of outliers
- Error is in same unit as that of the data points.

RMSE & MAE = which metric is better without outliers.

↳ Both RMSE & MAE are used to measure the model's accuracy.

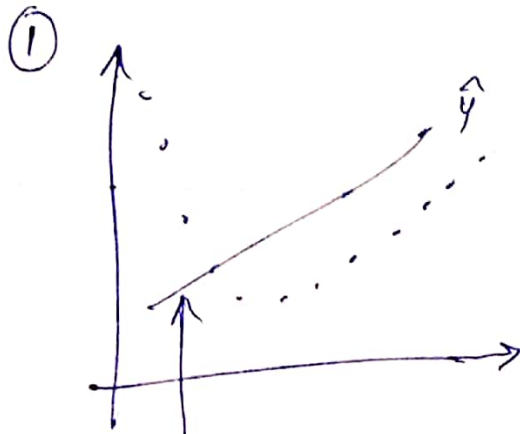
RMSE & MAE = which metric is better with outliers

↳ In RMSE, since errors are squared before they are averaged, RMSE gives relatively high weight to large error.

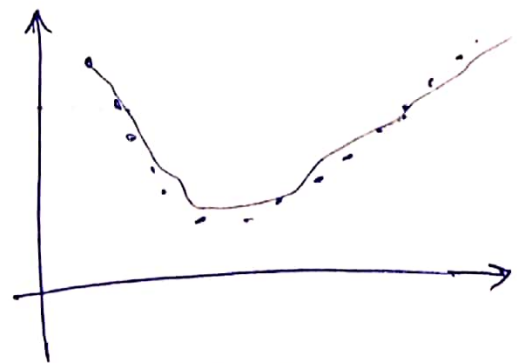
↳ Thus ~~RMSE~~ RMSE is more useful when large errors occur.

Assumptions in Linear Regression

- A Linear Relation between dependent & independent Variables.
- No Multicollinearity between independent Variables.
- Residuals (errors) are homoscedastic.
- Residuals are normally distributed.



Regression line doesn't capture much information due to distribution of data points

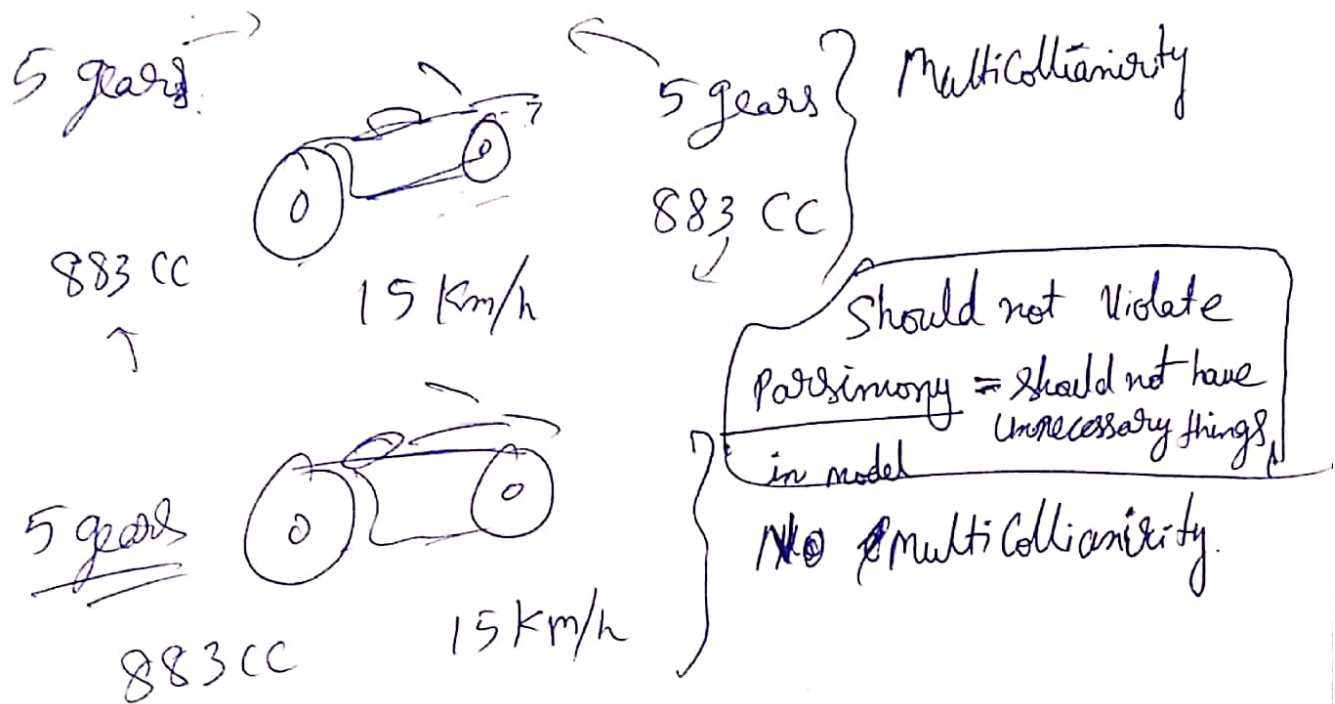


A non-linear approach works better in this scenario

② A statistical phenomenon where,

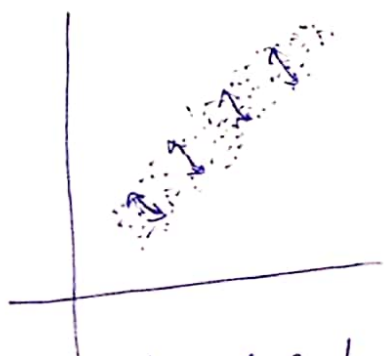
→ Independent explanatory variables are highly positive correlated with one another.

→ Highly correlated independent variables explain the same variance about dependent variable. Thus introducing Redundancy.

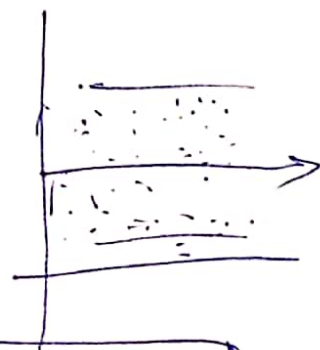


③ Homoscedasticity?

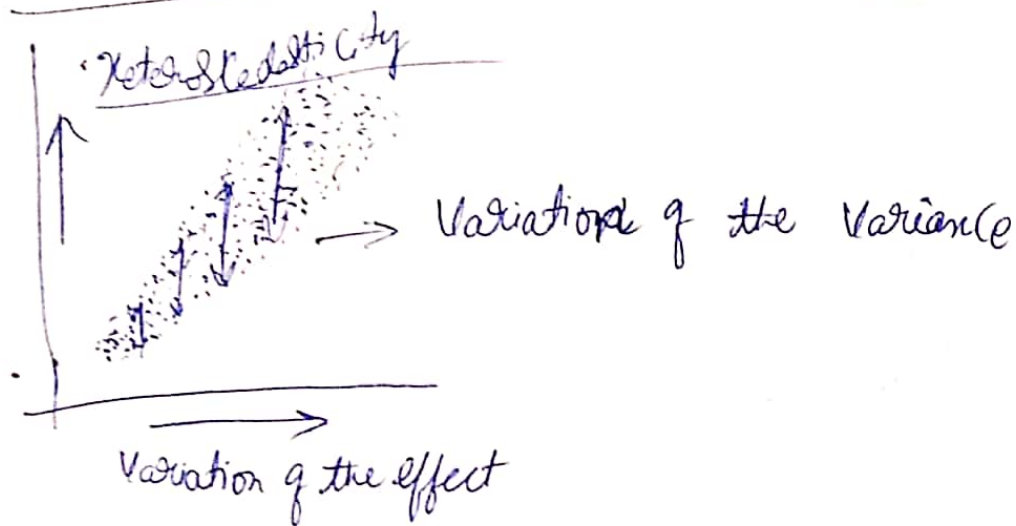
Homo = equal + Scedastic = dispersion = Homoscedastic implies Constant Variance.



Spread of residual values is Constant as we move along x-axis.



Heteroscedastic looks like this

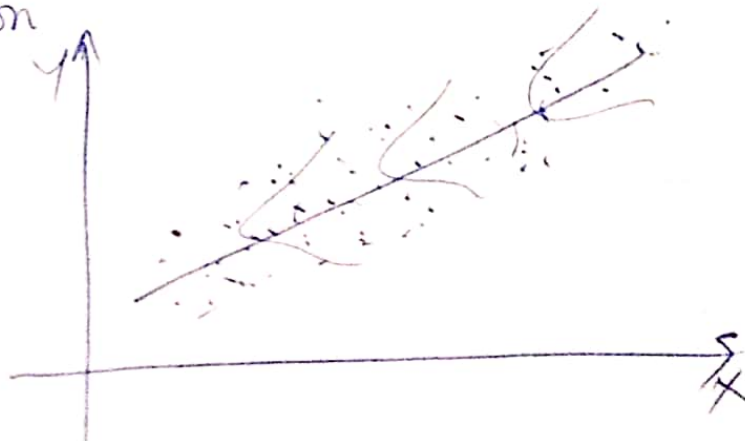


Spread of residual values increases as we move along ~~the~~ x-axis.

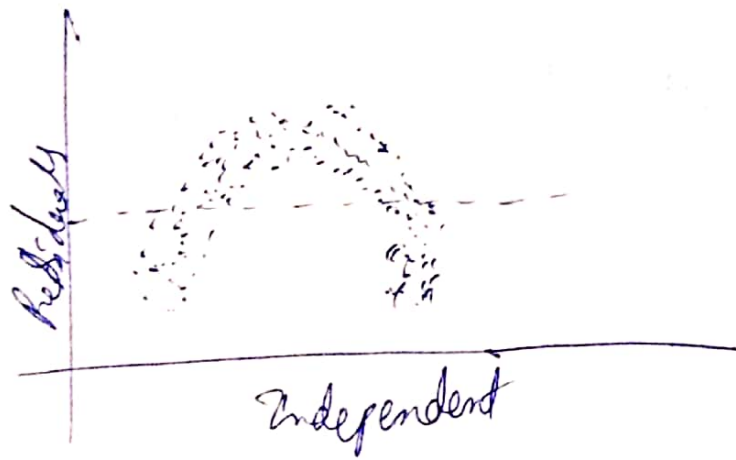
Bigger residuals imply greater error in Prediction

(i) Normal Distribution of residuals.

Data points should be random and follow a normal distribution



Errors follow a curved pattern thus introducing bias



Exceptions to normal distribution assumption

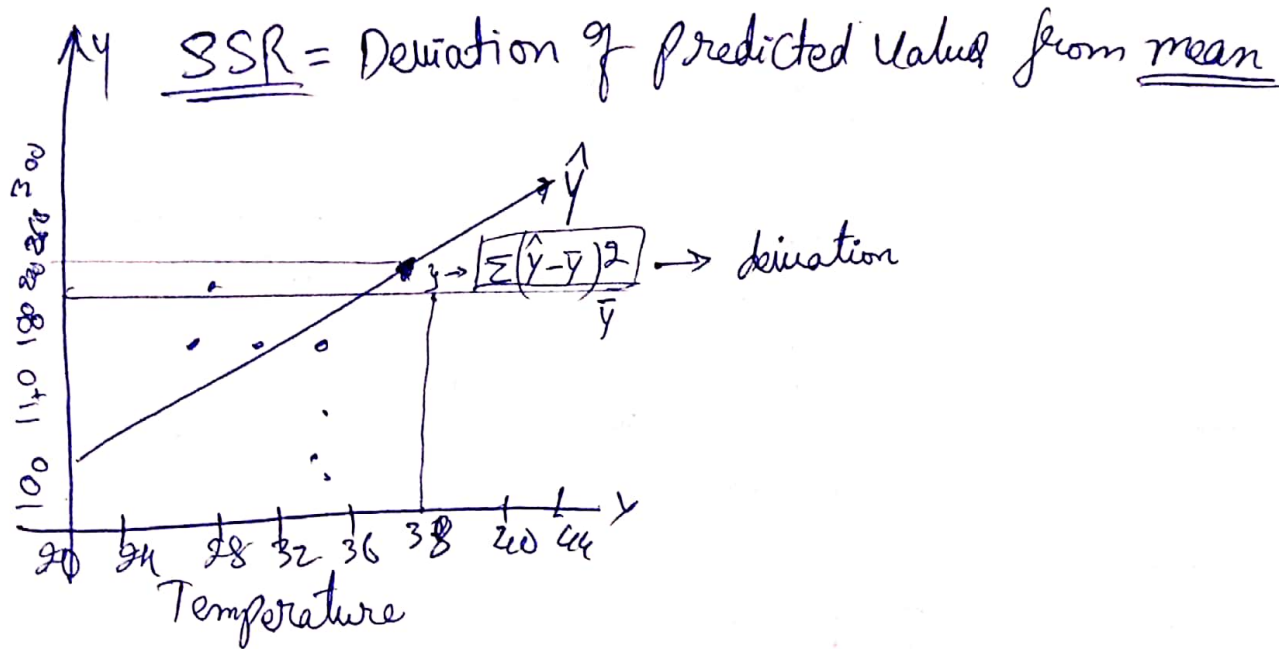
- ↳ For larger samples, residuals may follow any pattern or biasness but still linear regression
- ↳ Normality is the least important assumption of linear regression.

Terminologies associated with error

↳ SSR = ~~sum~~ Sum square due to regression.
↳ Explained deviation from mean.

↳ SSE = Sum square due to Error.
↳ Unexplained deviation from mean.

↳ SST = Sum square Total ($SSE + SSR$)

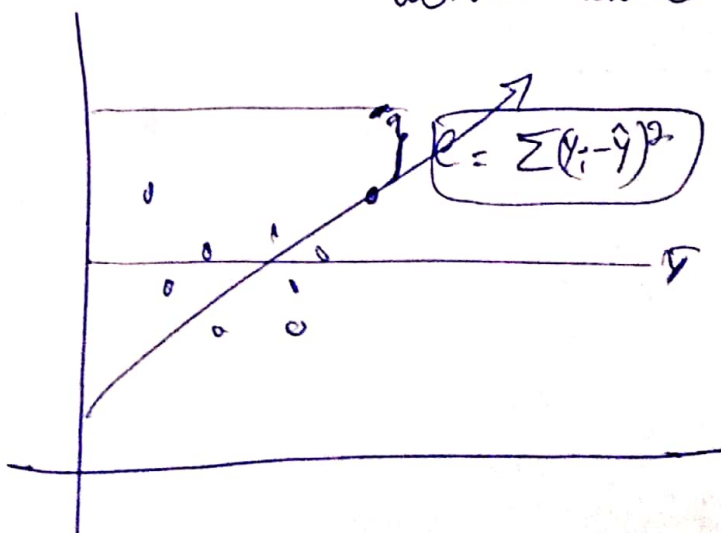


→ At 38°C on a particular day, we predicted 260 ice cream to be sold which is higher than mean \bar{Y} .

→ This deviation collectively for all predictions is termed as SSR (Sum square due to Regression) or ESS (Explained Sum of Square)

Inference : we accept this deviation & use it to calculate R^2 .

SSE = Deviation of Predicted values from actual value



→ At 38°C on a particular day, Actual no. of ice cream sold was 300 which is higher than our Prediction (260 ice cream)

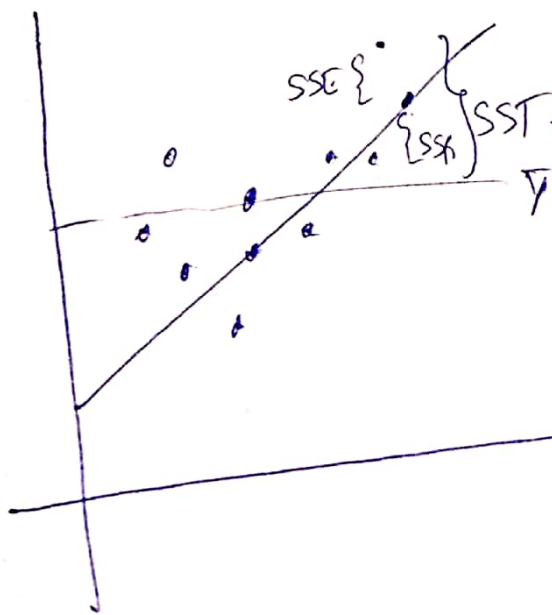
→ This deviation or difference is called SSE (Sum Squared Error). The deviation is Unexpected or is Unexplained from mean \bar{Y} .

Inference

→ This difference Can't be explained by the difference in independent variable.

→ We want this as small as possible,

SST = Deviation of actual values from mean

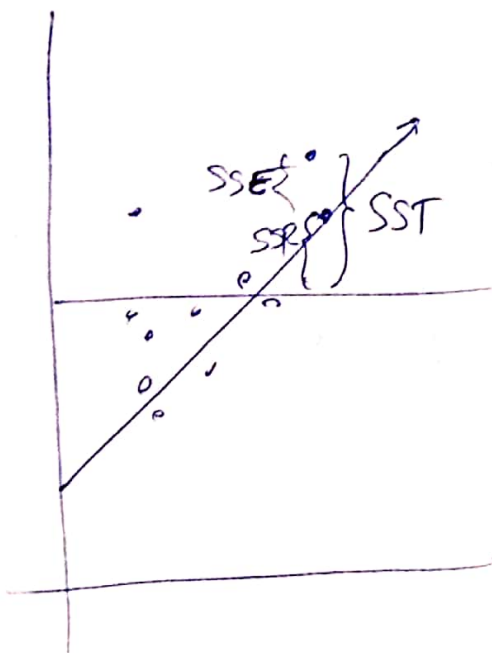


$SST = \sum (Y_i - \bar{Y})^2$ → Sum square total (SST) is the deviation of actual data point (Y_i) from mean (\bar{Y}).

→ SST is composed of two parts, SSE & SSR

$$\boxed{SST = SSE + SSR}$$

R-Squared



$$R^2 = \frac{SSR}{SST} \quad \text{or} \quad 1 - \frac{SSE}{SST}$$

↳ most popular metric to measure performance of regression model.

↳ Tells about Variability or difference in y variable explained by difference in x variable(s).

R^2 = Coefficient of determination

↳ R^2 is square of correlation coefficient (r) b/w dependent (y) & independent (x) variables.

↳ r tells us the degree of association b/w x & y .

→ For Example, $R^2 = 0.81$ so $r = 0.9$, implies 90% change in y is known to x & vice versa.

↳ Remaining 10% association b/w x & y is Unknown & this is why x can't explain remaining 10% difference (SSE) in y .

Adjusted R squared

Adjusted R squared is more powerful ~~to~~ metric than R squared to measure the linear regression model.

↳ It measures the proportion of variation explained by only those independent variables that really help in explaining the dependent variable.

↳ It penalizes you for adding independent variable that do not help in predicting the independent variable.

Variables	R-squared	Adjusted R-squared
1	67.5	67.1
2	85.9	84.2
3	88.9	81.7

↳ Adding 2nd variable adj. R^2 has increased along with R^2 because the variable has real impact on the dependent variable

→ But 3rd Variable increased R^2 & Reduced Adj- R^2 implies that the Variable is not impactful at all.

R-squared	Adjusted R-squared
Every time an independent Variable added to the model, the R-square increases, even if the independent Variable is insignificant. R^2 never decrease.	Adjusted R-squared increases only when independent Variable is significant and affects dependent Variable
R^2 Can be zero	Adj R^2 Can be negative

Machine Learning Algo's

Steps involved in ML are

- ① Load Dataset
 - ② EDA \rightarrow Are exhaustive?
& Extensive one
- } PRE ML

③ Create Features & label

④ Split, Train and Test Data
and Cross Validation

⑤ Create Instant of the model

⑥ Fit the Model.

⑦ PREDICT test data using trained model

⑧ Evaluate model Performance?
Using Evaluation metrics

Sub-Steps
need to be
explored
by me



Linear Regression

All steps with intermediary steps ~~are~~ included in the algo is as below. [No template is necessary, but to ~~write~~ remember all steps are essential -]

① Import the ~~data~~ packages (load the data)

② Import Packages

③ Explore data

Bikes.info
Bikes.head

} Clean Data } \Rightarrow Pandas

} Visualize Data } \Rightarrow Seaborn & Matplotlib

Step 1 \rightarrow Create features & label

~~X = bikes.drop(['Count'], axis=1)~~

X = bikes.drop(['Count'], axis=1, inplace=False)

y = bikes['Count']

X.shape

y.shape

✍

P.T.O.

Step 2 - Split Train & Test Data only for Supervised

$X_{\text{train}}, X_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train_test_split}(X, y,$

$\text{test_size}=0.2, \text{random_state}=999)$

from sklearn.model_selection import train_test_split → any no
↓ → To get the same output
Need to verify

$X_{\text{train}}.shape$

$y_{\text{train}}.shape$

$X_{\text{test}}.shape$

$y_{\text{test}}.shape$

Step 3 - Create Instance of the model.

from sklearn.linear_model import LinearRegression.

$lm = \text{LinearRegression}()$

Step 4 - Fit model

$lm.fit(X_{\text{train}}, y_{\text{train}})$

Step 5 Predict using the trained model.

$\text{Predicted} = lm.predict(X_{\text{test}})$

$\text{Predicted}.shape$

Step 6 Evaluate model performance

metrics . mean - squared - error (y-test, predicted)
metrics . mean - absolute - error (y-test, predicted)
metrics . median - absolute - error (y-test, predicted)

→ lm.coef

→ x.column → just to see the attributes/variables

→ lm.intercept

→ Additional lines of codes to be entertained

→ In the last step we are appending the predicted house prices into the original data and computing Error in estimation for the test data.

$fdf = \text{pd.concat}([test_x, test_y], 1)$

$fdf['predicted'] = \text{np.round}(\text{predict_test}, 1)$

$fdf['prediction_error'] = fdf['House_price'] - fdf['predicted']$

→ P.T.O.

Fitting using different methods

Using Stats models to build linear Regression

```
import statsmodels.formula.api as smf
```

```
model = smf.ols('Medv ~ CRIM + ZN + ...',  
data, data=Boston_df).fit()
```

↑
all features

```
Print(model.summary())
```

Remove Insignificant Variables — Based on
P-Contributions

Print Heteroscedasticity Robust Std Errors

```
Robust_model = model.get_robustcov_results()
```

```
Print(Robust_model.summary())
```