

## Linear Regression Subjective Questions

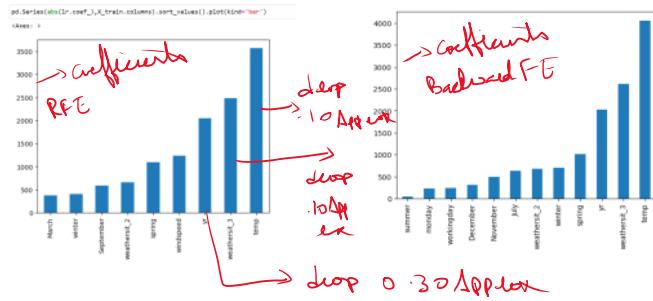
Tuesday, October 10, 2023 11:49 PM

### Assignment Based

1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable?
- (A) Yr has the highest positive correlation (.57) with the target variable among categorical variables. It is also the 2<sup>nd</sup> highest coefficient (slope) among the categorical variables. While weathersit\_3 (a component generated from the weather sit column) has the highest coefficient weight among categorical variables. Spring a component of dummy encoding the seasons column is important too, also spring got a correlation score ( $r$ -Pearson's Correlation coefficient) of -0.56, a negative correlation really close to Yr, but in the opposite direction. It was observed that the yr column had considerable impact on the r2-score. Removing yr column led to drop of approx .30 in the r2-score, this was observed while experimenting with different feature sets. Another notable but less important feature would be components of mnth(month).

```
print(stats_model1_training(x_train_copy(),y_train_copy()))
OLS Regression Results
Dep. Variable:    sales
R-squared:       0.823
Model: OLS
Method: Least Squares
Date: Tue, 10 Oct 2023
Time: 23:44:33
F-statistic: 1.41e+10, 129 (df=1, 379)
Log-Likelihood: -4245.4
AIC: 8490.8
BIC: 8497.0
DF Residuals: 379
DF Model: 2
Covariance Type: nonrobust
```

```
print(stats_model1_training(x_train_copy(),y_train_copy()))
OLS Regression Results
Dep. Variable:    sales
R-squared:       0.823
Model: OLS
Method: Least Squares
Date: Tue, 10 Oct 2023
Time: 23:44:33
F-statistic: 1.41e+10, 129 (df=1, 379)
Log-Likelihood: -4245.4
AIC: 8490.8
BIC: 8497.0
DF Residuals: 379
DF Model: 2
Covariance Type: nonrobust
```



2. Why is it important to use drop\_first=True during dummy variable creation?

(A) It is important to use drop\_first=True, because not doing so can lead to multicollinearity. The variance of the nth dummy variable created explained by the n-1 dummy variables that come are extracted. Multicollinearity cause problems in model interpretation and it also makes the model more complex by adding a additional features. The interpretation that the coefficient of x1 , m1 represent the change in y wrt x1, when x1 increases by 1 unit and all other independent variables are constant doesn't hold.

Eg: Let's say you have a categorical variable that represents the final result of a student exam

Result

Pass	
Fail	
Absent	

→ dummy encoding

	Pass	Fail	Absent
Pass	1	0	0
Fail	0	1	0
Absent	0	0	1

Pass & Fail can completely explain Absent. This also gives a high VIF (Variance Inflation Factor) Score

∴ It can be concluded that drop\_first=True, is very important

→ As year increasing year led to such a large drop in r2-score

Although going by the coeffs temp is most imp, going by drop in r2-score based on feature absence year is most important.

Eg from assignment

```
# Dummy encoding categorical variables
df = pd.concat([df,pd.get_dummies(df['season'],drop_first=True)],axis=1)
df = pd.concat([df,pd.get_dummies(df['weathersit'],drop_first=True)],axis=1)
df = pd.concat([df,pd.get_dummies(df['weekday'],drop_first=True)],axis=1)
df = pd.concat([df,pd.get_dummies(df['mnth'],drop_first=True)],axis=1)
```

Eg if what happens when there is multicollinearity

↳ it can be inferred from here the variance/length of the Absent column can be inferred from Pass & Fail

Features VIF  
1 holiday inf  
2 workingday inf  
7 monday inf  
8 tuesday inf  
9 wednesday inf  
10 thursday inf  
11 friday inf  
14 temp 4.71  
8 windspeed 4.82  
0 yr 2.06  
3 spring 1.84  
5 weathersit\_2 1.53  
4 winter 1.41  
12 March 1.26  
13 September 1.16  
6 weathersit\_3 1.09

→ The result of VIF = Inf is b'cause holiday, working day & weekdays can completely explain variance of the

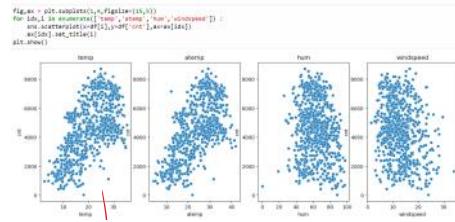
```
print(vif_method(x_train_copy()))
```

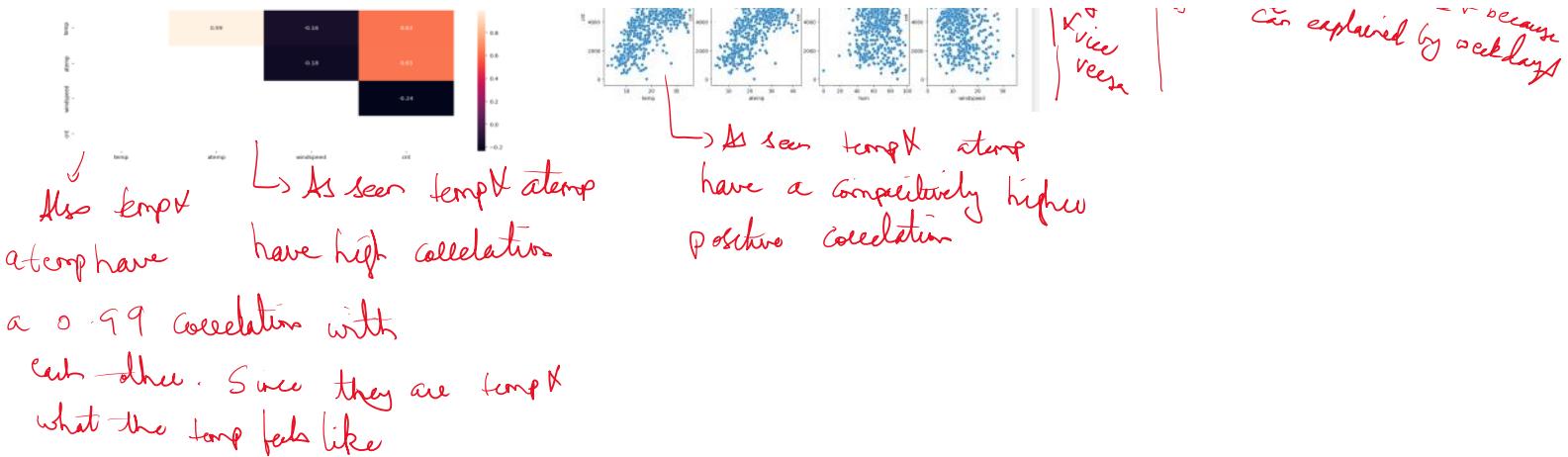
Features VIF  
1 workingday 50.31  
8 wednesday 8.58  
7 tuesday 7.27  
9 thursday 6.98  
6 monday 6.57  
13 temp 4.71  
14 windspeed 4.71  
0 yr 2.86  
2 spring 1.84  
4 weathersit\_2 1.53  
5 winter 1.41  
11 March 1.26  
13 September 1.16  
5 weathersit\_3 1.09

→ effect of removing holiday \* working day still has high vif because its variance can be explained by weekdays

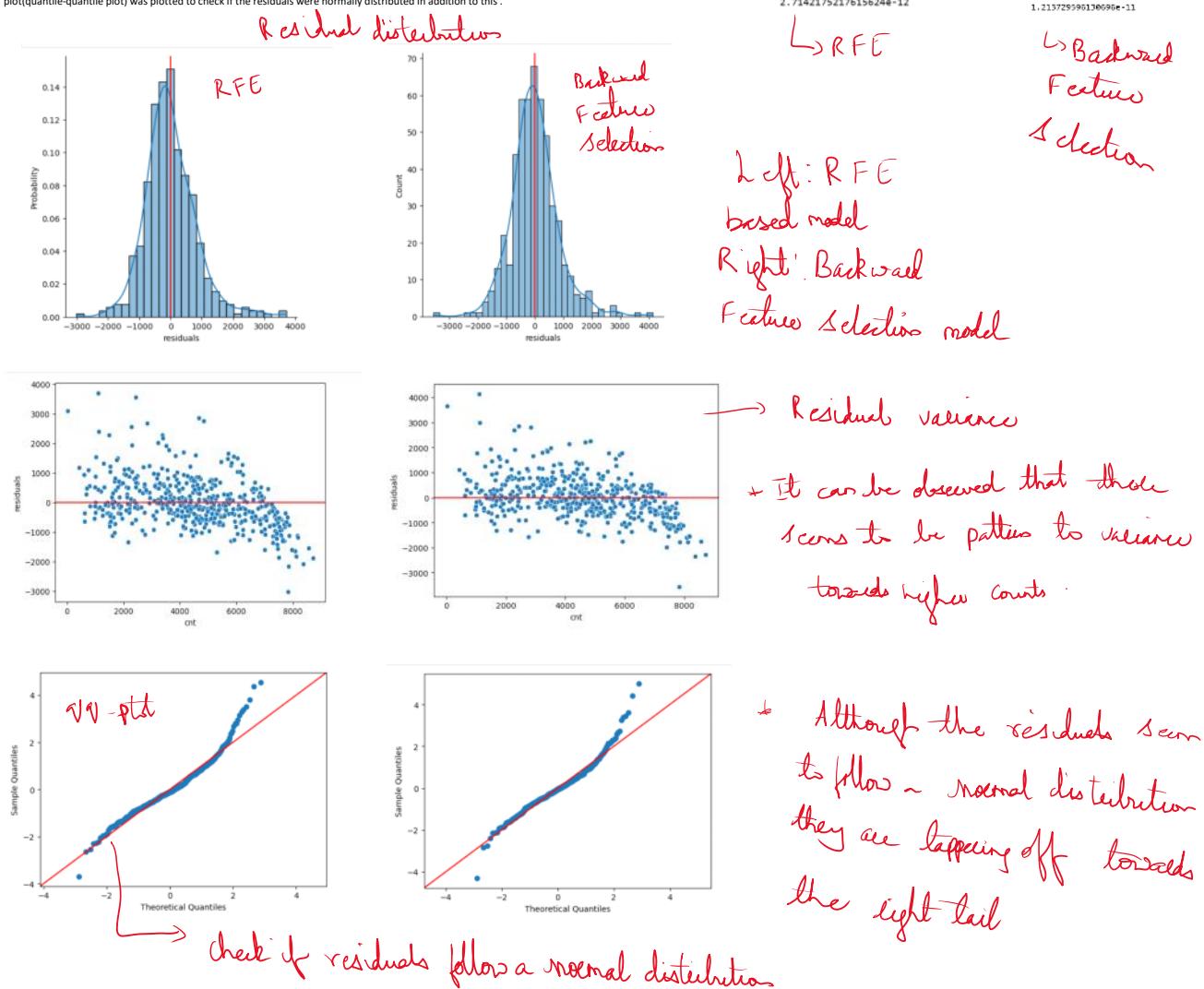
3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable?

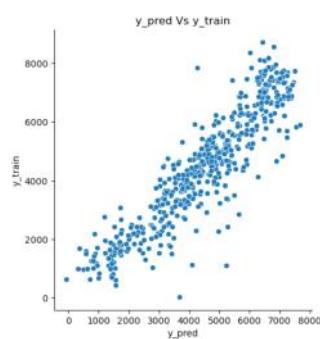
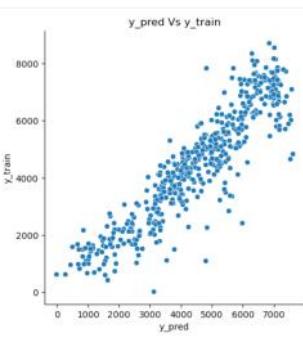
(A) Disregarding registered and casual , which are a part of cnt (cnt = registered + casual) . The features with highest correlation to the target variable are temp and atemp





4. How did you validate the assumptions of Linear Regression after building the model on the training set ?  
 (A) The model was validated by checking the mean of the residuals . The residuals followed a normal distribution and mean of the residuals was infinitesimally close to zero . Also there weren't any large notable patterns overall . There was slight decrease in residuals as cnt(target variable) increased. Also a qq-plot (quantile-quantile plot) was plotted to check if the residuals were normally distributed in addition to this .



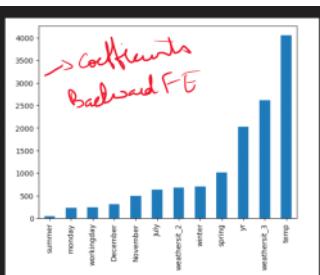
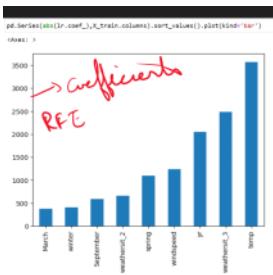


→  $y_{pred}$  vs  $y_{train}$

There is a +ve linear relationship here which is a good thing

5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes?

(a) It can be observed from the below plots that represent coefficients that temp , weather\_sit3 and yr are the most important features to the model. This was verified across 2 models that used RFE and the other backward feature elimination.



### Base Model (FF)

```
DLS Regression Results
=====
Dep. Variable:          y_train
   OLS
Model: Least Squares
Date: Mon, 11 Oct 2023
Time: 20:48:45
Jmp. No. of Observations: 510
Df Residuals: 491
Df Model: 18
Covariance Type: nonrobust
...
coef std err t P>|t| [0.025 0.975]
const 1833.3562 278.226 6.598 0.000 1286.713 2379.999
yr 1.41e-173 0.000 4.427 0.000 0.411 2.420
workingday -284.7758 81.019 3.514 0.000 97.712 411.448
weather_sit1 -239.0899 101.098 2.306 0.021 -94.052 394.984
weather_sit2 -236.0899 101.098 2.306 0.021 -94.052 394.984
weather_sit3 -235.6663 107.987 2.160 0.031 -83.263 520.266
temp -476.7044 79.651 6.000 0.000 -645.832 314.458
...
Prob(Omnibus): 0.888 Jarque-Bera (2R): 209.263
Skew: 0.442 Kurtosis: 5.824 Cond. No.: 28.4
```



```
X_train = X_train.drop(columns = ['yr'])
print(stats_model_training(X_train.copy(),y_train.copy()))
DLS Regression Results
=====
Dep. Variable:          y_train
   OLS
Model: Least Squares
Date: Mon, 11 Oct 2023
Time: 20:48:45
Jmp. No. of Observations: 510
Df Residuals: 491
Df Model: 18
Covariance Type: nonrobust
...
coef std err t P>|t| [0.025 0.975]
const 3880.0002 143.878 26.916 0.000 3680.034 4018.982
workingday -285.5397 83.812 3.382 0.000 119.798 2519.101
weather_sit1 -238.9892 103.282 2.318 0.021 -97.322 374.976
weather_sit2 -236.9892 103.282 2.318 0.021 -97.322 374.976
weather_sit3 -235.5663 107.887 2.160 0.031 -83.052 520.266
temp -476.7044 79.651 6.000 0.000 -645.832 314.458
...
Prob(Omnibus): 0.885 Jarque-Bera (2R): 207.648
Skew: 0.442 Kurtosis: 5.821 Cond. No.: 19.2
```

R-squared 0.823 to 0.558

→  $y_e$  removed

\* Note this is not final model . this is the output of feature Elimination.

being displayed for example purposes

→ temp removed

```
X_train = X_train.drop(columns = ['temp'])
print(stats_model_training(X_train.copy(),y_train.copy()))
DLS Regression Results
=====
Dep. Variable:          y_train
   OLS
Model: Least Squares
Date: Mon, 11 Oct 2023
Time: 20:48:45
Jmp. No. of Observations: 510
Df Residuals: 491
Df Model: 17
Covariance Type: nonrobust
...
coef std err t P>|t| [0.025 0.975]
const 3880.0002 143.878 26.916 0.000 3680.034 4018.982
workingday -285.5397 83.812 3.382 0.000 119.798 2519.101
weather_sit1 -238.9892 103.282 2.318 0.021 -97.322 374.976
weather_sit2 -236.9892 103.282 2.318 0.021 -97.322 374.976
weather_sit3 -235.5663 107.887 2.160 0.031 -83.052 520.266
...
Prob(Omnibus): 0.885 Jarque-Bera (2R): 207.648
Skew: 0.442 Kurtosis: 5.821 Cond. No.: 19.2
```

R-Squared 0.823 to 0.767

→ WeatherSit\_3 removed

R-Squared 0.823 to 0.773

```
DLS Regression Results
=====
Dep. Variable:          y_train
   OLS
Model: Least Squares
Date: Mon, 11 Oct 2023
Time: 20:48:45
Jmp. No. of Observations: 510
Df Residuals: 491
Df Model: 16
Covariance Type: nonrobust
...
coef std err t P>|t| [0.025 0.975]
const 1313.9797 313.448 4.227 0.000 897.968 2129.999
yr 1.41e-173 0.000 4.427 0.000 0.411 2.420
workingday -233.7758 86.163 2.223 0.027 24.998 382.082
spring -148.3875 101.354 1.464 0.148 -97.733 235.468
winter -426.8289 283.216 2.140 0.036 77.553 826.183
weather_sit1 -235.3314 122.228 -1.760 0.079 -455.283 35.019
Monday -797.8803 238.478 -3.399 0.000 -1173.128 -344.042
March 184.0623 154.762 1.189 0.233 -12.071 348.976
April -248.9464 182.148 -1.333 0.182 -52.081 442.811
May -248.9464 182.481 -1.333 0.182 -52.081 442.811
...
Prob(Omnibus): 0.886 Jarque-Bera (2R): 717.815
Skew: 0.442 Kurtosis: 5.819 Cond. No.: 28.3
```

Generic Questions

1. Explain the linear regression algorithm in detail ?  
 (A) Is a supervised machine learning algorithm , meaning it required target labels for the model to be trained . It is part of the regression family of ml algorithms . Linear regression aims to fit the target variable using a straight line as its name suggests. It is used to predict the behavior of a continuous variable .

## Equation of Linear Regression

$$y = \beta_0 + \sum_{i=1}^n \beta_i x_i$$

where  $\beta_0$  is your intercept &  $\beta_i$  is the coefficient / slope of the  $i^{th}$  feature ( $x_i$ )

$n \rightarrow$  no. of features

## Assumptions of Linear Regression :

- \* There should be linear relationship between  $X$  (Independent/Predictor)
- \*  $y$  (Dependent/Target) variables -
- \* The error terms that is the residuals should be normally distributed , with their mean centered on 0.
- \* The variance of residuals should be constant (homoscedasticity)
- \* Error terms should be independent of one another. If the error at  $y=10$  is high it is not that the error at  $y=20$ , should be higher or lower. The errors shouldnt depend on the previous error

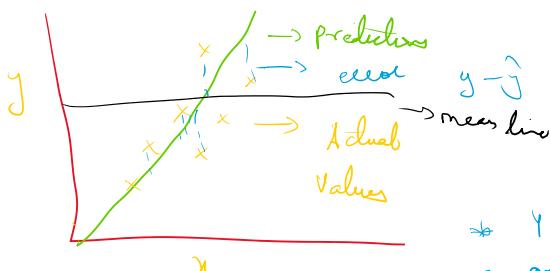
The equations of linear regression can be solved in 2 ways

OLS & Gradient Descent .

↳ closure function      ↳ Iterative approach

Error terms :

$$\text{Mean Squared Error} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|^2 = \frac{1}{n} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$$



The errors here are called residuals .

\* You want the errors to be as low as possible



\* You want the error to be as low as possible

+ Although Mean Squared Error gives you error of the model, it cannot be determined if the model you have is good or bad, as there is no values that defines it.

This is where R-Score / R-Squared comes in

$$R\text{-Score} = 1 - \frac{\text{MSE of Regression Model}}{\text{MSE of Mean Line}}$$

↓

Values lie b/w chosen for  $x$ ,  
 you are to 't' letters  
 closer you are to predicted  
 closer you are to  
 value your model predicts

MSE of Mean Line

↳ Mean of all  $y$  values is taken

↳ a line is drawn

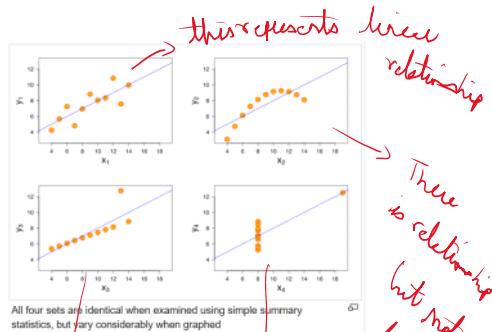
↳ calculate MSE for this

2. Explain the Anscombe's quartet in detail?

(A) Anscombe's quartet contains 4 datasets which extremely similar descriptive stats i.e. the same mean, variance and correlations. However as seen in the same they have very different distribution. This emphasizes the need to always plot the data, as these stats can be misleading. Anscombe's quartet provides important of plotting graphs especially in linear regression. These 4 datasets give extremely similar coefficients, intercepts and  $R^2$ -scores for. If not plotted it could lead to the misconception that the datasets are identical.

Anscombe's quartet							
I	II	III	IV				
$x$	$y$	$x$	$y$	$x$	$y$	$x$	$y$
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89

→ this is the dataset



All four sets are identical when examined using simple summary statistics, but very considerably when graphed

this represents linear relationship  
 True relationship but not linear  
 No linear relationship with outliers  
 relationship at all. But there is a high correlation

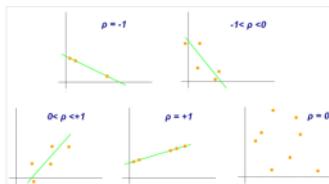
This proves that plotting data is important.  
As it is not a good idea to linear model in II & IV

III shows that a single outlier can lower the scale of a model

I seems to be simple linear relationship model

3. What is Pearson's R ?

- (A) Is pearson's correlation coefficient . It is used to measure the correlation between linear data . This is the most widely used correlation coefficient , though it doesn't well on non-linear data . Usually spearman's correlation is used to measure correlations in case data is monotonic . The value of r ranges from -1 to 1 . Closer to 1 indicates a strong positive correlation . Closer to zero indicates a very weak correlation



→ Covariance btw X & Y

$$r = P(x,y) = \frac{\text{Cov}(x,y)}{\sqrt{\text{Var}(x)\text{Var}(y)}} \rightarrow \text{Covariance of } y \\ \hookrightarrow \text{Variance of } x$$

pd corr () → gives you pearson correlation coefficient

(r) by default for all combinations in the data frame

4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling ?

- (A) Scaling is the process of bringing continuous variables in a dataset to a similar range . Scaling is performed to prevent a single feature from overly influencing the outcome , especially in models like liner regression , logistic regression and knn , as values with a higher scale usually tend to dominate the predictions .

Standardized Scaling or Standardization : is a scaling technique which scales features into a standard normal distribution with mean = 0 and std = 1 . Also 99.7% of the data is present within 3 std from the mean .

Normalization or Normalized Scaling or MinMax Scaling : is a scaling technique which brings all the features into a range between 0 and 1 . The max value would be 1 and the minimum value would be 0 . It also handles outliers . It compresses the outliers , so that need they are no longer considered outliers after scaling

Standardization

$$z_{\text{Scale}} = \frac{x - \mu}{\sigma}, \quad \mu = 0, \sigma = 1$$

$$\text{Min Max Scaling} = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \in [0, 1] \\ \hookrightarrow \text{handles outliers}$$

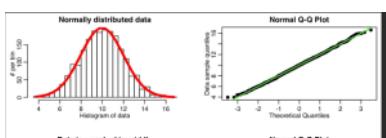
5. You might have observed that sometimes the value of VIF is infinite . Why does this happen ?

- (A) VIF also known as Variance Inflation Factor estimates how much variance of the a dependent variables is being explained by the other dependent variables . The VIF would be infinite if the variance of the dependent variable is being completely explained by the other dependent variables in the dataset . It is sign of very high multi collinearity .

$$VIF = \frac{1}{1 - R^2} \quad , \quad VIF = \infty , \text{ when } R^2 = 1 \rightarrow \underbrace{\text{Values completely explained}}_{\hookrightarrow \text{perfect linear relations}}$$

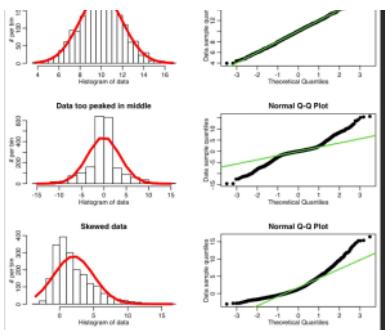
6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression.

- (A) Q-Q plots that expand to Quantile-Quantile Plots are used to measure if 2 datasets come from the same distribution . In linear regression q-q plots can be used to determine if the residuals follow normal distribution . q-q plots have a theoretical quantile and a sample quantile . The theoretical quantile in this case of linear regression would be a normal distribution and the sample quantiles would be the percentile values of the residuals .



Interpretation of Q-Q Plot

If the sample distribution is same as



If the sample distribution is same as theoretical distribution the pts would lie on the straight line -

∴ The more the pts lie on the straight line - more similar the distributions

Inversely, as more & more pts taper away from the line, more distinct the distributions