

Submitted by:

Ramandeep Mehra

Program name:

Executive PG Programme in Machine Learning & AI - Dec 2023

Assignment-based Subjective Questions

1. From your analysis of the categorical variables from the dataset, what could you infer about their effect on the dependent variable? (3 marks)

Answer:

We converted these 4 categorical variable into encoded dummy variable for the model prediction: 'season', 'weathersit', 'mnth', 'weekday'

```

OLS Regression Results
=====
Dep. Variable: cnt R-squared: 0.840
Model: OLS Adj. R-squared: 0.836
Method: Least Squares F-statistic: 218.0
Date: Thu, 29 Feb 2024 Prob (F-statistic): 3.50e-189
Time: 17:57:04 Log-Likelihood: 506.36
No. Observations: 510 AIC: -986.7
Df Residuals: 497 BIC: -931.7
Df Model: 12
Covariance Type: nonrobust
=====
            coef  std err      t    P>|t|   [0.025  0.975]
-----
const      0.2739  0.024  11.321  0.000    0.226  0.321
yr         0.2348  0.008  28.880  0.000    0.219  0.251
holiday   -0.1054  0.026  -4.087  0.000   -0.156 -0.055
temp        0.4318  0.031  13.986  0.000    0.371  0.492
windspeed  -0.1482  0.025  -5.956  0.000   -0.197 -0.099
spring     -0.1006  0.015  -6.513  0.000   -0.131 -0.070
winter     0.0453  0.012   3.680  0.000    0.021  0.069
mist       -0.0821  0.009  -9.464  0.000   -0.099 -0.065
snow       -0.2930  0.024  -11.991 0.000   -0.341 -0.245
jan        -0.0436  0.018  -2.473  0.014   -0.078 -0.009
jul        -0.0668  0.017  -3.888  0.000   -0.101 -0.033
sep        0.0574  0.016   3.690  0.000    0.027  0.088
sun        -0.0482  0.012  -4.170  0.000   -0.071 -0.026
-----
Omnibus: 61.960 Durbin-Watson: 2.030
Prob(Omnibus): 0.000 Jarque-Bera (JB): 162.608
Skew: -0.609 Prob(JB): 4.90e-36
Kurtosis: 5.484 Cond. No. 14.5

```

To infer about their effect on the target variable, we can conclude from the final Multiple Linear Regression Summary Table:

- From the Categorical Dummy variable **season** we have considered spring and winter seasons in our model, where spring shows a negative coefficient and winter shows a positive coefficient. As we can see sales go up in winters and drop in spring.
- From the Categorical Dummy variable **month**, in the month of jan, jul the sales show a negative coefficient difference wherein in the month of sept it is showing a positive relationship with demand of the bikes.
- From the Categorical Dummy variable **week day**, as we can see here on sun the demand is dropping, as coefficient value is negative.

- From the Categorical Dummy variable **weathersit**, we can conclude that neither snow weather nor mist weather is showing a positive coefficient.

2. Why is it important to use `drop_first=True` during dummy variable creation? (2 mark)

Answer:

We use `drop_first=True` during dummy variable creation because it helps avoid multicollinearity and reduces the risk of introducing redundant information in the model.

a). Multicollinearity avoidance:

- Multicollinearity refers to a situation where two or more predictor variables in a regression model are highly correlated with each other. This can lead to issues in interpreting the individual coefficients of these variables and can make the model unstable.

- When you create dummy variables without dropping the first category, you essentially introduce perfect multicollinearity. This happens because if you have n categories for a categorical variable, you only need n-1 dummy variables to represent all the information.

- Dropping the first category ensures that one category serves as the reference category (the baseline) against which all other categories are compared. This eliminates multicollinearity as there is no perfect correlation between the dummy variables.

b). Reducing dimensionality:

- Including all dummy variables without dropping the first one increases the dimensionality of your dataset, which can be problematic in some cases. High-dimensional data can lead to increased computation time, overfitting, and a need for more data to train the model effectively.

- By dropping the first category, you reduce the dimensionality by one, making your model more manageable and potentially improving its performance, especially if you have a limited amount of data.

c). Interpretability:

- Including all dummy variables can make the interpretation of model coefficients more complex. When you drop the first category, the coefficients for the remaining categories represent the difference between each category and the reference category.

- This makes it easier to understand the effect of each category on the target variable in comparison to the reference category.

Here's an example to illustrate the concept:

Suppose you have a categorical variable "Color" with three categories: Red, Blue, and Green. If you create dummy variables without dropping the first one, you'll have three dummy variables (e.g., Color_Red, Color_Blue, Color_Green). However, using `'drop_first=True'`, you'll only create two dummy variables (e.g., Color_Blue, Color_Green), and Red becomes the reference

category. Now, the coefficients for Color_Blue and Color_Green will tell you how they differ from Red in terms of their effect on the target variable.

In summary, using `drop_first=True` when creating dummy variables is important to prevent multicollinearity, reduce dimensionality, and enhance the interpretability of your machine learning models.

3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (1 mark)

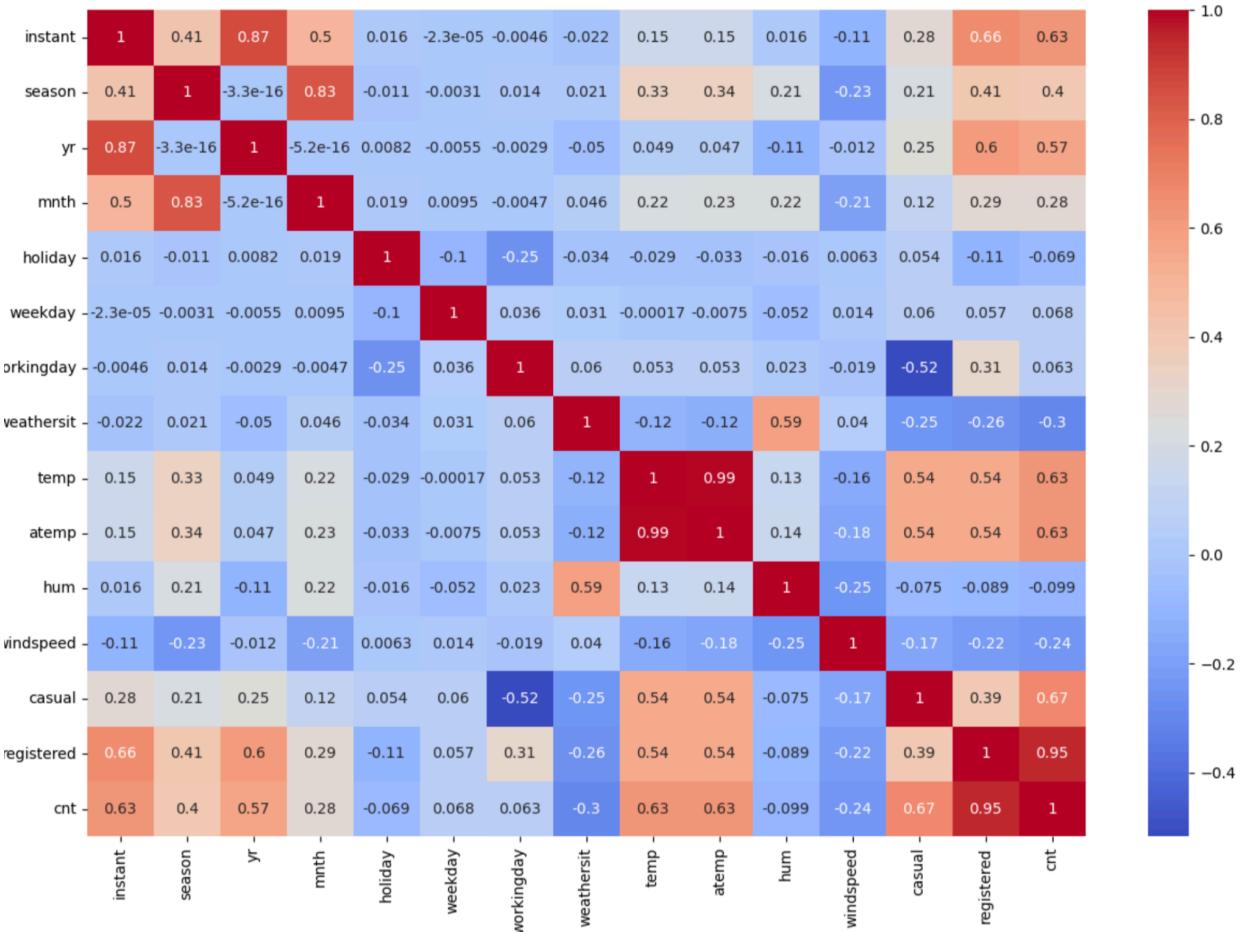
Answer:

Before model building and splitting the dataset into training and testing, the pair plot shows highest correlation for **registered and cnt**, **casual and cnt** variables having correlation 0.95 and 0.67. Hence we are not using casual and registered in our pre-processed training data for model building. As **casual + registered = cnt**. cnt Variable is the sum of casual and registered variables and the model might get overfit during the model building and have the multicollinearity issue.

We also excluded this variable atemp. As it has the highest correlation with temp and one variable is explaining each other, to avoid the multicollinearity and overfitting we kept only one variable.

As per the correlation heatmap, correlation coefficient between atemp and cnt is 0.63. And correlation coefficient between temp and cnt is 0.63.

Instant variable also has the highest correlation with the target variable cnt which is 0.63.



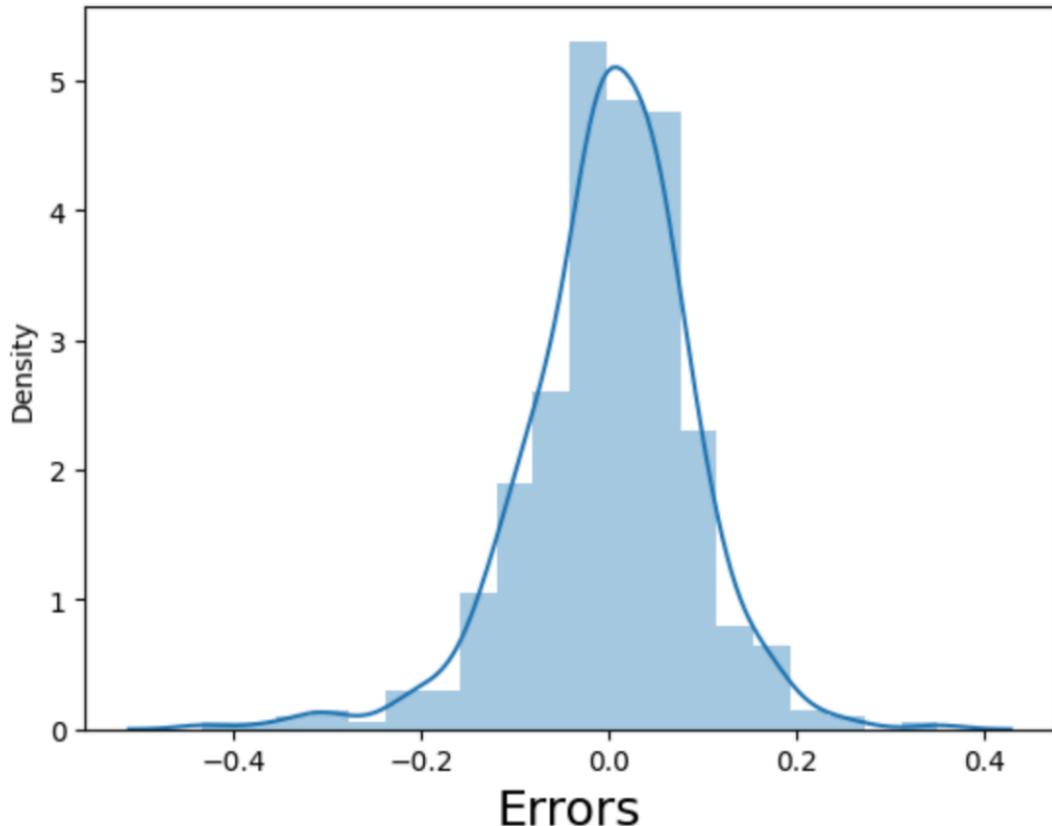
4. How did you validate the assumptions of Linear Regression after building the model on the training set? (3 marks)

Answer:

To validate assumptions of the model, and hence the reliability for inference, we go with the following procedures:

- Residual Analysis: We need to check if the error terms are also normally distributed (which is in fact, one of the major assumptions of linear regression). I have plotted the histogram of the error terms and this is what it looks like:

Error Terms



The residuals are following the normally distributed Curve with a mean 0 and maximum data points are distributed around the 0.

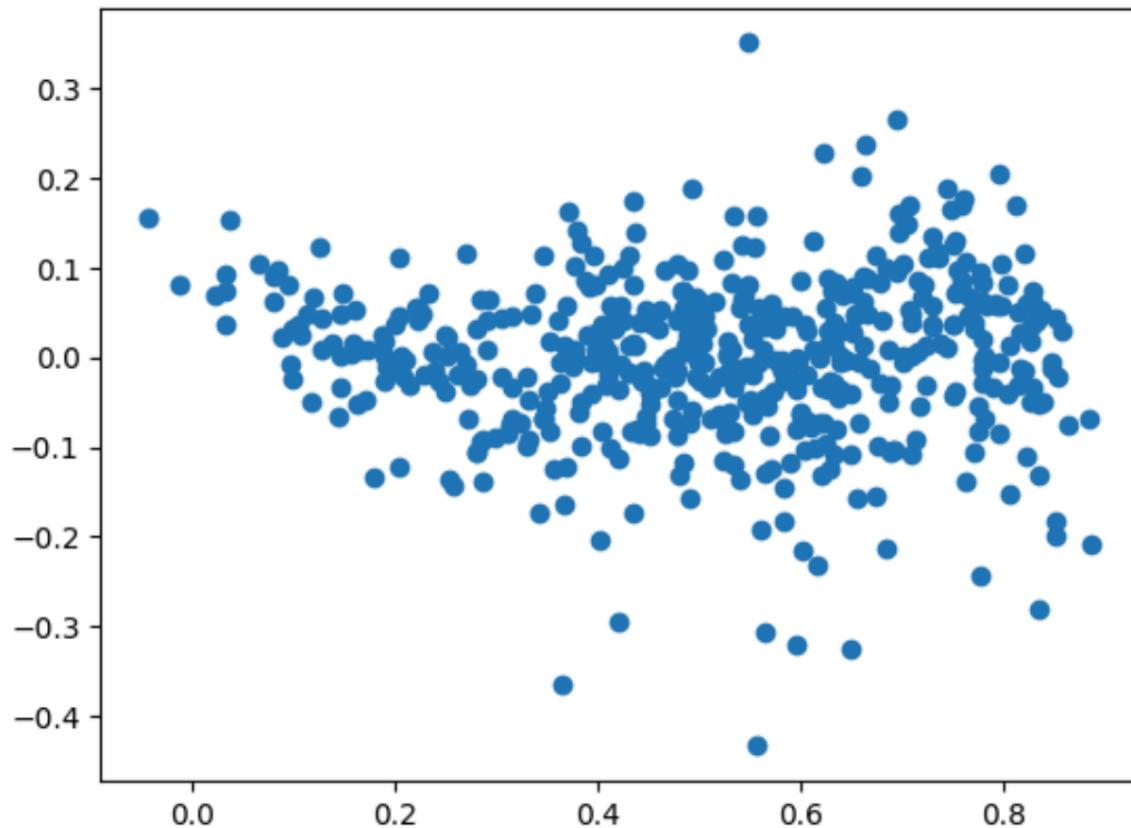
- Linear relationship between predictor variables and target variable:

This is happening because all the predictor variables are statistically significant (p-values are less than 0.05). Also, R-Squared value on training set is 0.84 and the adjusted R-Squared value on the training set is 0.83600. This means that variance in data is being explained by all these predictor variables and the target variable is explained by the independent variable well.

- Error terms are independent of each other:

Handled properly in the model. The predictor variables are independent of each other. Multicollinearity issue is not there because VIF (Variance Inflation Factor) for all predictor variables are below 5.

From the below plot we can see that residuals (also known as error terms) are independent. This means that there is no correlation between the residuals and the predicted values, or among the residuals themselves.



We are confident that the model fit isn't by chance, and has decent predictive power. The normality of residual terms allows some inference on the coefficients.

5. Based on the final model, which are the top 3 features contributing significantly towards explaining the demand of the shared bikes? (2 marks)

Answer:

Top 3 features which are significantly contributing towards demand of shared bikes are mentioned below:

- 1) **temp** which has highest positive coefficient as: 0.4318
- 2) **yr** which has second highest positive coefficient as: 0.2348
- 3) **sep** which has third highest positive coefficient as: 0.0574

General Subjective Questions

1. Explain the linear regression algorithm in detail.

(4 marks)

Answer:

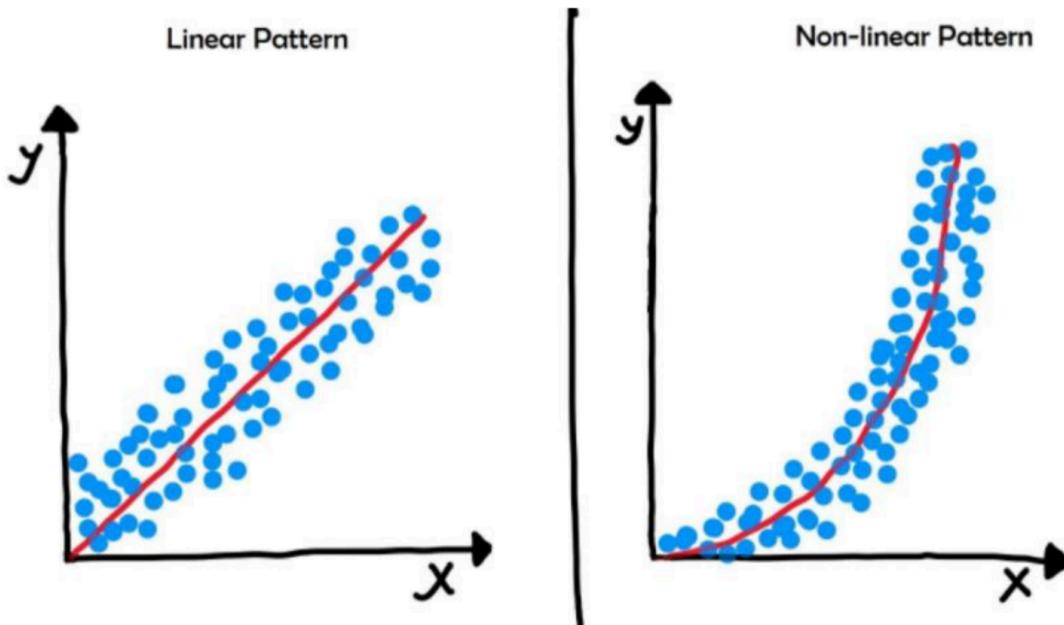
Linear Regression finds the best linear relationship between the independent and dependent variables.

It is a method of finding the best straight-line fitting to the given data.

In technical terms, linear regression is a machine learning algorithm that finds the best linear-fit relationship on any given data, between independent and dependent variables. It is mostly done by the Sum of Squared Residuals Method.

The assumptions of linear regression are:

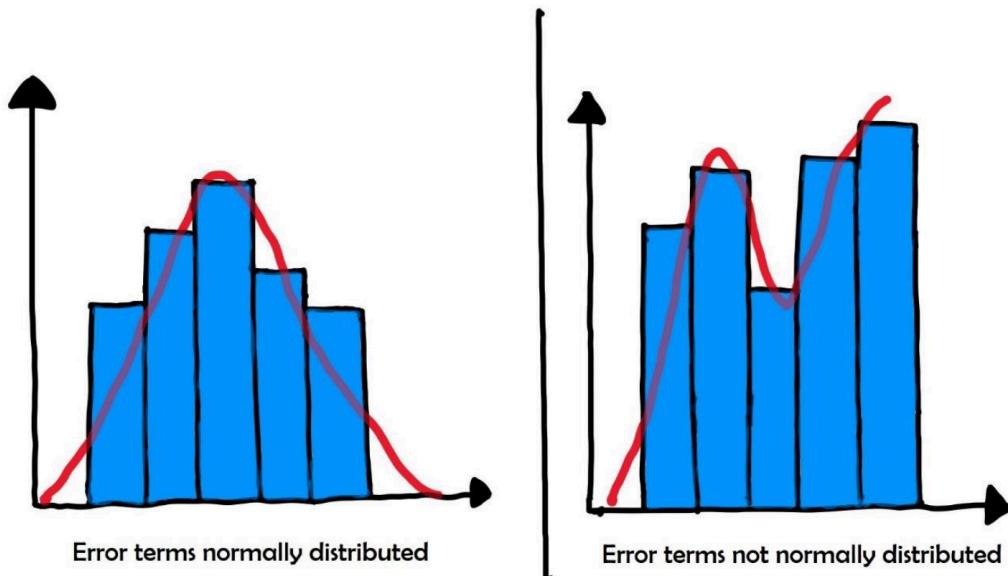
- The assumption about the form of the model: It is assumed that there is a linear relationship between the dependent and independent variables.



- Assumptions about the residuals:

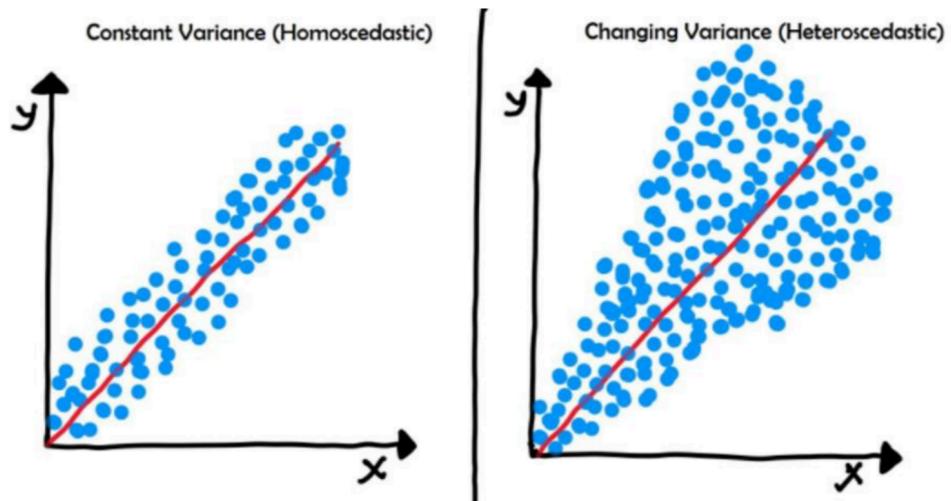
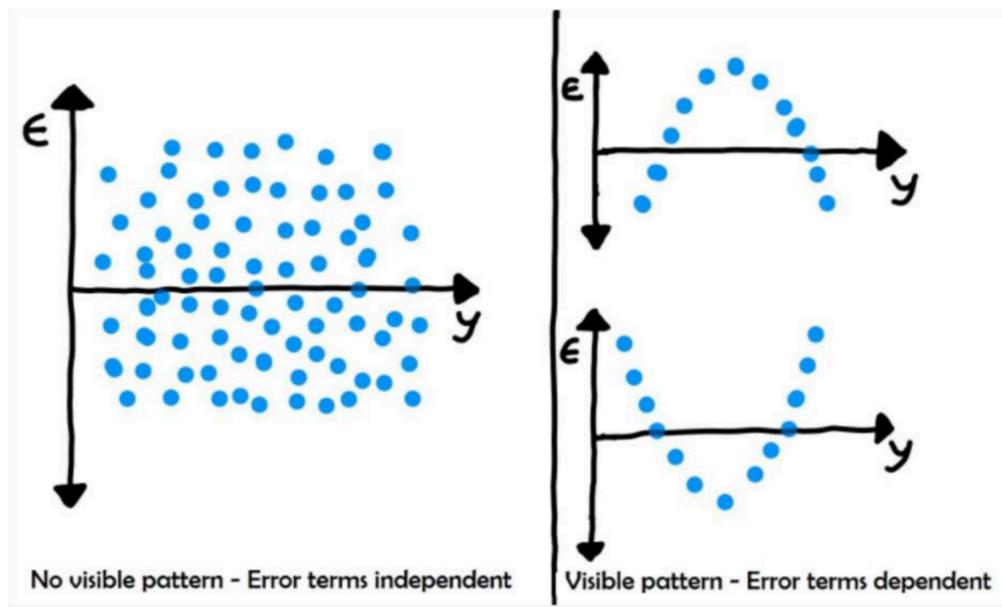
- 1) Normality assumption: It is assumed that the error terms, $\varepsilon(i)$, are normally distributed.
- 2) Zero mean assumption: It is assumed that the residuals have a mean value of zero, i.e., the error terms are normally distributed around zero.
- 3) Constant variance assumption: It is assumed that the residual terms have the same (but unknown) variance, sigma square. This assumption is also known as the assumption of homogeneity or homoscedasticity.

4) Independent error assumption: It is assumed that the residual terms are independent of each other, i.e., their pairwise covariance is zero.



c. Assumptions about the estimators:

- 1) The independent variables are measured without error.
- 2) The independent variables are linearly independent of each other, i.e., there is no multicollinearity in the data.

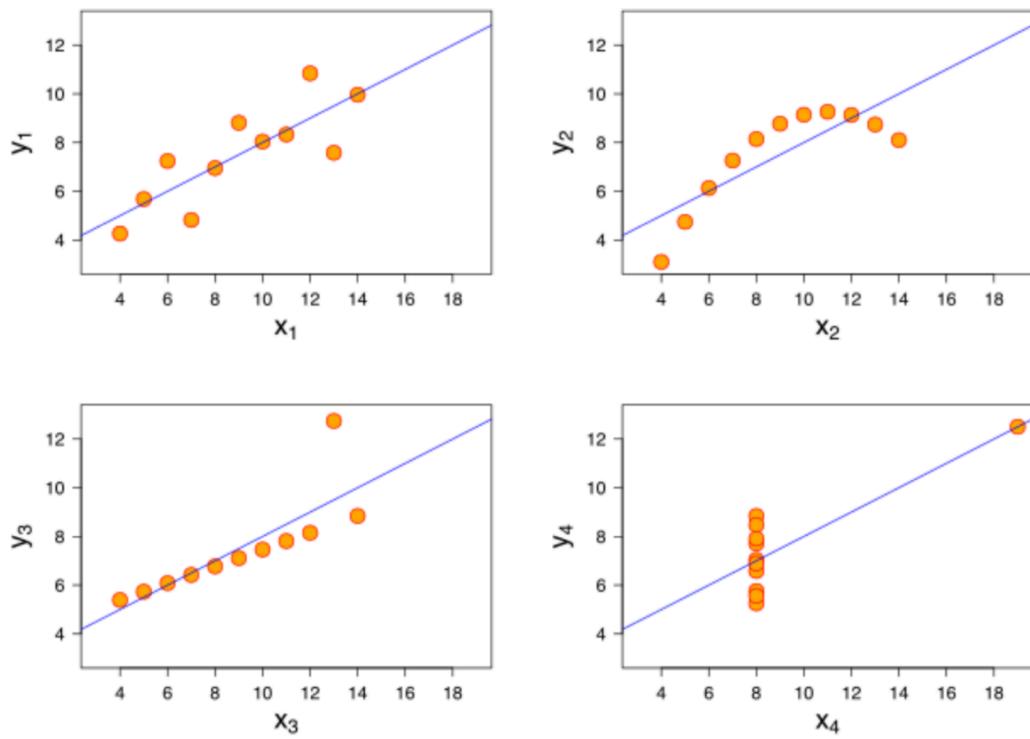


2. Explain the Anscombe's quartet in detail.

(3 marks)

Answer:

Anscombe's quartet comprises four data sets that have nearly identical simple descriptive statistics, yet have very different distributions and appear very different when graphed.



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

- 1) The first scatter plot (top left) appears to be a simple linear relationship, corresponding to two variables correlated where y could be modeled as gaussian with mean linearly dependent on x .
- 2) The second graph (top right) is not distributed normally; while a relationship between the two variables is obvious, it is not linear, and the Pearson correlation coefficient is not relevant. A more general regression and the corresponding coefficient of determination would be more appropriate.
- 3) In the third graph (bottom left), the distribution is linear, but should have a different regression line (a robust regression would have been called for). The calculated regression is offset by the one outlier which exerts enough influence to lower the correlation coefficient from 1 to 0.816.
- 4) the fourth graph (bottom right) shows an example when one high-leverage point is enough to produce a high correlation coefficient, even though the Other data points do not indicate any relationship between the variables.

Property	Value	Accuracy
Mean of x	9	exact
Sample variance of $x : s_x^2$	11	exact
Mean of y	7.50	to 2 decimal places
Sample variance of $y : s_y^2$	4.125	± 0.003
Correlation between x and y	0.816	to 3 decimal places
Linear regression line	$y = 3.00 + 0.500x$	to 2 and 3 decimal places, respectively
Coefficient of determination of the linear regression : R^2	0.67	to 2 decimal places

Anscombe's quartet

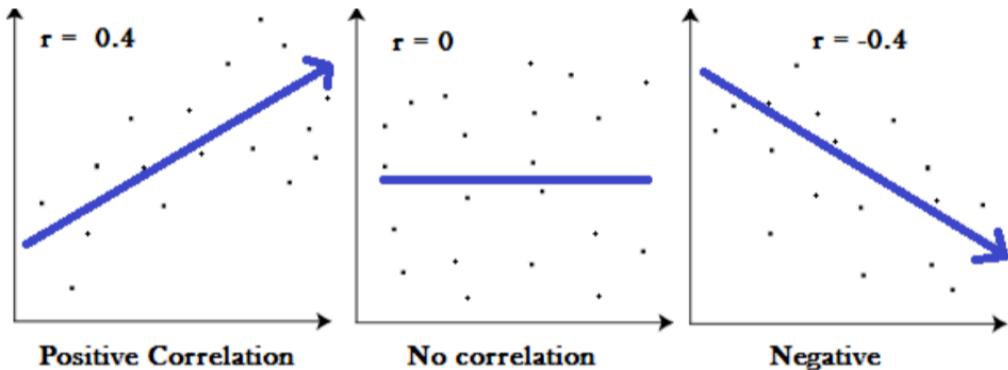
I		II		III		IV	
x	y	x	y	x	y	x	y
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	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

3. What is Pearson's R?

(3 marks)

Answer:

Pearson's R or correlation coefficient is a measure of linear correlation between two sets of data. It is the ratio between the covariance of two variables and the product of their standard deviations; thus, it is essentially a normalized measurement of the covariance, such that the result always has a value between -1 and 1. As with covariance itself, the measure can only reflect a linear correlation of variables, and ignores many other types of relationship or correlation.



- 1) A correlation coefficient of 1 means that for every positive increase in one variable, there is a positive increase of a fixed proportion in the other. For example, shoe sizes go up in (almost) perfect correlation with foot length.
- 2) A correlation coefficient of -1 means that for every positive increase in one variable, there is a negative decrease of a fixed proportion in the other. For example, the amount of gas in a tank decrease in (almost) perfect correlation with speed.
- 3) Zero means that for every increase, there isn't a positive or negative increase. The two just aren't related.

The absolute value of the correlation coefficient gives us the relationship strength. The larger the number, the stronger the relationship. For example, $|-0.95| = .95$, which has a stronger relationship than .55.

4. What is scaling? Why is scaling performed? What is the difference between normalized scaling and standardized scaling? (3 marks)

Answer:

Scaling is a method used to normalize the range of independent variables or features of data.

Since the range of values of raw data varies widely, in some machine learning algorithms, objective functions will not work properly without normalization. For example, many classifiers calculate the distance between two points by the Euclidean distance. If one of the features has a broad range of values, the distance will be governed by this particular feature. Therefore, the range of all features should be normalized so that each feature contributes approximately proportionately to the final distance.

Another reason why feature scaling is applied is that gradient descent converges much faster with feature scaling than without it.

Normalization:

Also known as min-max scaling or min-max normalization, it is the simplest method and consists of rescaling the range of features to scale the range in [0, 1]. The general formula for normalization is given as:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Here, $\max(x)$ and $\min(x)$ are the maximum and the minimum values of the feature respectively.

Standardization:

Feature standardization makes the values of each feature in the data have zero mean and unit variance. The general method of calculation is to determine the distribution mean and standard deviation for each feature and calculate the new data point by the following formula:

$$x' = \frac{x - \bar{x}}{\sigma}$$

Here, σ is the standard deviation of the feature vector, and \bar{x} is the average of the feature vector.

5. You might have observed that sometimes the value of VIF is infinite. Why does this happen? (3 marks)

Answer:

If there is perfect correlation, then $VIF = \infty$. This shows a perfect correlation between two independent variables. In the case of perfect correlation, we get $R^2 = 1$, which lead to $1/(1-R^2) = \infty$. To solve this problem, we need to drop one of the variables from the dataset which is causing this perfect multicollinearity.

An infinite VIF value indicates that the corresponding variable may be expressed exactly by a linear combination of other variables (which show an infinite VIF as well).

6. What is a Q-Q plot? Explain the use and importance of a Q-Q plot in linear regression. (3 marks)

Answer:

Q-Q Plots (Quantile-Quantile plots) are plots of two quantiles against each other. A quantile is a fraction where certain values fall below that quantile. For example, the median is a quantile where 50% of the data fall below that point and 50% lie above it. The purpose of Q-Q plots is to find out if two sets of data come from the same distribution. A 45-degree angle is plotted on the Q-Q plot; if the two data sets come from a common distribution, the points will fall on that reference line.

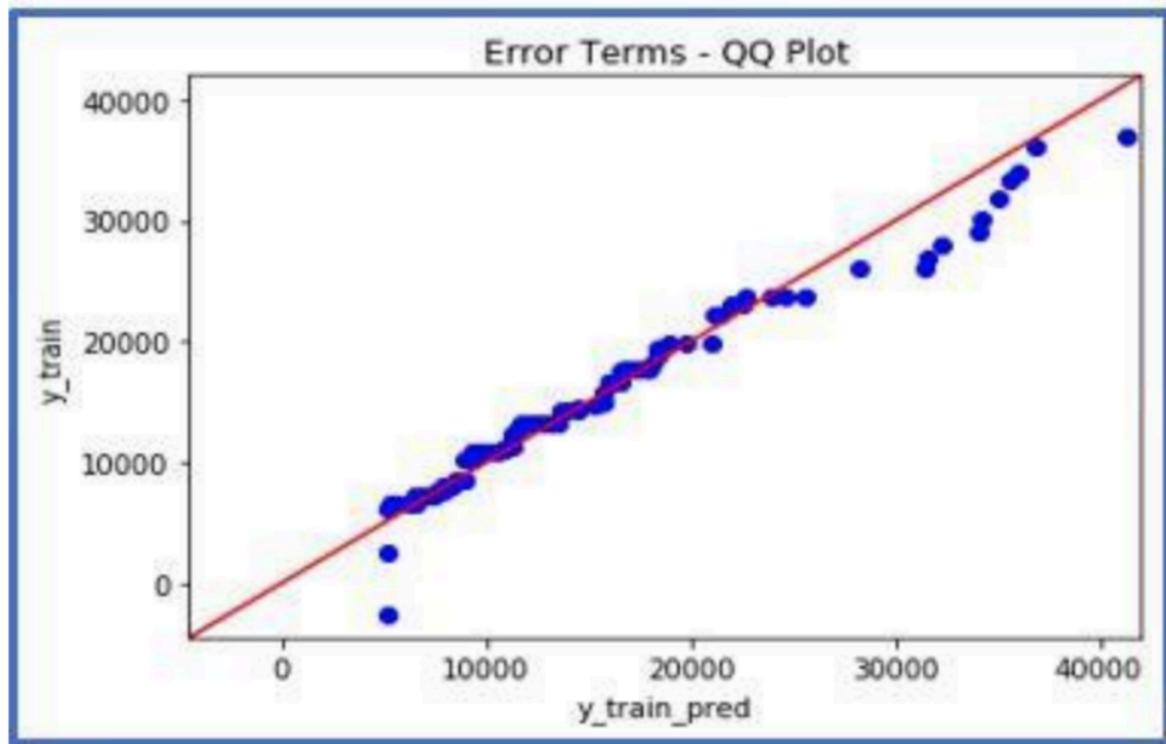
This helps in a scenario of linear regression when we have training and test data set received separately and then we can confirm using Q-Q plot that both the data sets are from populations with the same distributions.

Interpretation:

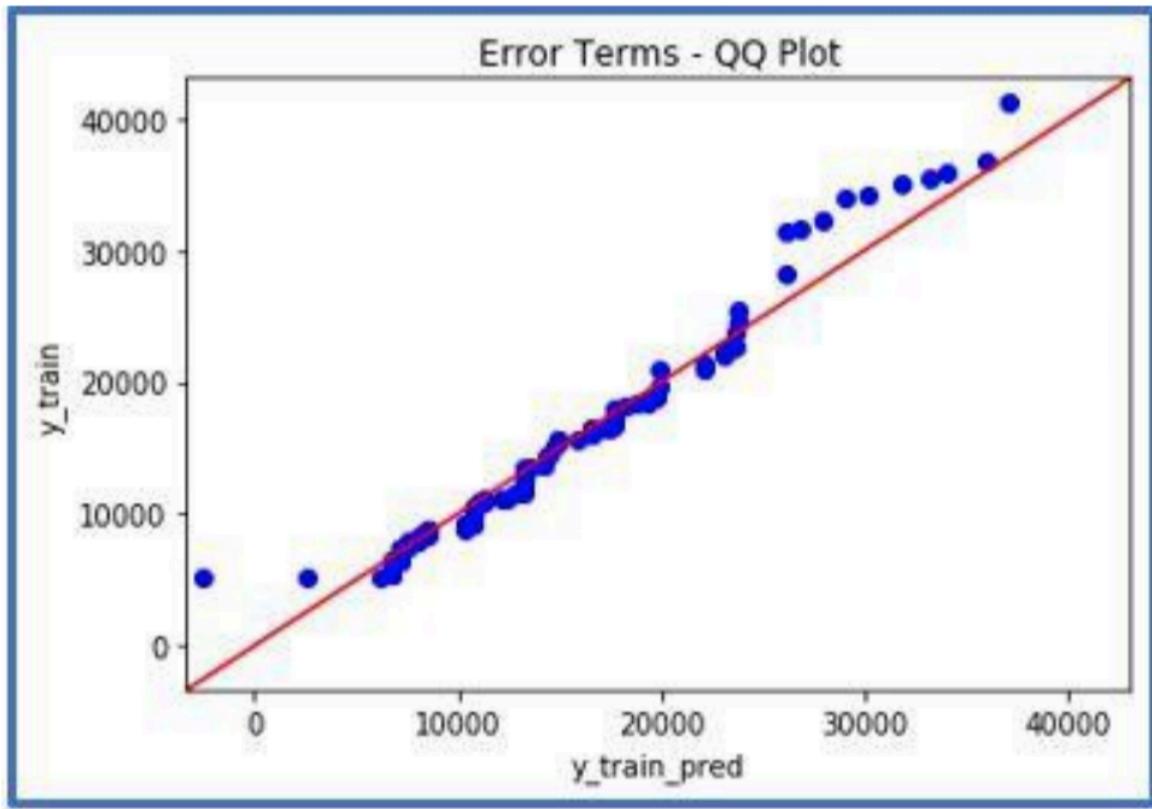
A q-q plot is a plot of the quantiles of the first data set against the quantiles of the second data set.

Below are the possible interpretations for two data sets.

- 1) Similar distribution: If all point of quantiles lies on or close to straight line at an angle of 45 degree from x -axis
- 2) Y-values < X-values: If y-quantiles are lower than the x-quantiles.



3) X-values < Y-values: If x-quantiles are lower than the y-quantiles.



- 4) Different distribution: If all point of quantiles lies away from the straight line at an angle of 45 degree from x -axis.

statsmodels.api provides **qqplot** and **qqplot_2samples** to plot Q-Q graphs for single and two different data sets respectively.