

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) .

The categorical variables are dependent and it depends on the number of variables that we are using if the dependent variable is normally distributed. Any change in one variable will impact the other one.

Season: Season_3 shows the maximum bookings with the median of greater than 5000, followed by season_2, season_4 and season_1 (where season_1:spring, season_2:summer, season_3:fall, season_4:winter). It is a good predictor variable to be considered in the model.

Mnth: As per the box plot the bookings are above 4000 for the months 4, 5, 6, 7, 8, 9,10.

Weathersit: Maximum bookings are done in weather1 with a median of around 5000, followed by weather2, weather3. It is a good predictor variable to be considered in the model.
weathersit : - 1: Clear, Few clouds, Partly cloudy, Partly cloudy - 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist - 3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds - 4: Heavy Rain + Ice Pellets + Thunderstorm + Mist, Snow + Fog)

Holiday: Majority of the booking about 97.2% is happening on non-holiday. The data doesn't seem to be distributed so it can't be considered as a good predictor. • Weekday: There is no clear difference in trend on all the days of the week. • Workingday: Bookings done on a working day have a median of greater than 5000. It can be a good predictor variable to be considered in the model.

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

For that column only the other dummy columns are created, so once the dummies are created that column doesn't have any meaning of its own .Ifcategorical variables are N then we use n-1 to represent the dummy variable

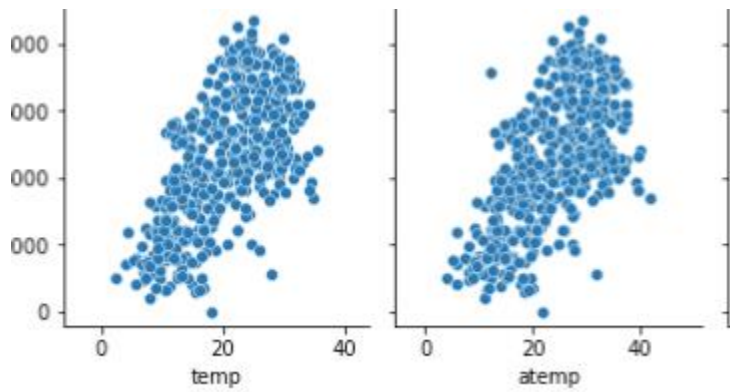
3. 3. Looking at the pair-plot among the numerical variables, which one has the highest correlation with the target variable? (1 mark)

Temperature and atemp has the highest correlation

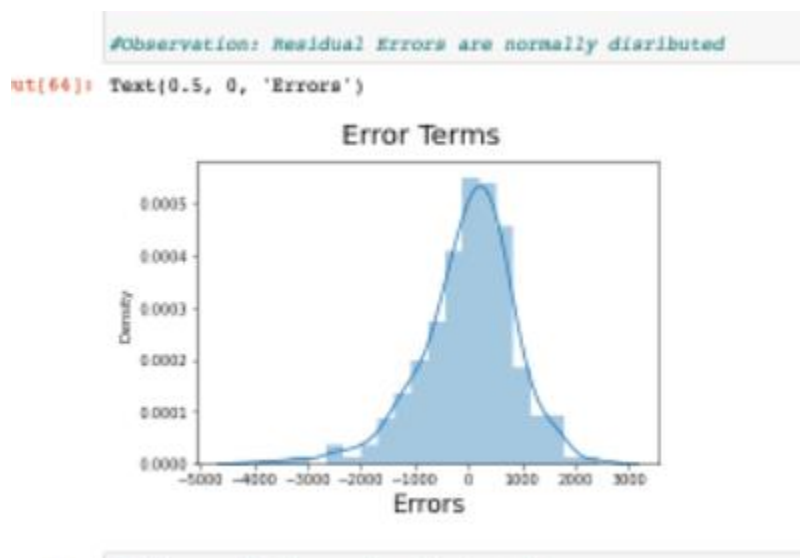
4. 4. How did you validate the assumptions of Linear Regression after building the model on the training set? (3 marks)

Analysis results from the Data set :

- Linearity : Dependent variable always shows linearity



- Homoscedasticity: The variance of the error term has to be same across all values of the independent variable. c).
- Normal Error: The error term should be normally distributed



- Multi-Collinearity is a phenomenon when two or more independent variables are highly correlated. Variance Inflation Factor (VIF) can be used to identify multi collinearity. A high value of VIF indicates high collinearity. VIF values of the Bike sharing case study

Out[60]:

	Features	VIF
2	temp	3.63
3	windspeed	2.97
0	yr	2.00
4	season_2	1.55
5	season_4	1.35
6	mnth_9	1.20
7	weathersit_3	1.06
1	holiday	1.03

- Multi-Collinearity is a phenomenon when two or more independent variables are highly correlated. Variance Inflation Factor (VIF) can be used to identify multi collinearity. A high value of VIF indicates high collinearity. VIF values of the Bike sharing case study

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

Temperature (temp) - A coefficient value of '4919.47' indicated that a unit increase in temp variable increases the bike hire numbers by 4919.47 units.

Weather Situation 3 (weathersit_3) - A coefficient value of '-1665.04' indicated that, with respect to Weathersit, a unit increase in Weathersit3 variable decreases the bike hire numbers by 1665.04 units, where weathersit_3 refers to Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds

yr - A coefficient value of '-2475.44' indicated that a unit increase in yr variable increases the bike hire numbers by 2475.44 units.

yr - A coefficient value of '1968.21' indicated that a unit increase in yr variable increases the bike hire numbers by 1968.21 units.

General Subjective Questions

1. Explain the linear regression algorithm in detail. (4 marks) 2.

Linear regression is a machine learning algorithm based on supervised learning, it shows a linear relationship between a dependent (y) and one or more independent (x) variables. Simple Linear Regression is represented by the equation $y = mx + c$ where m is the slope and c is the intercept with the axis. It helps in finding how the value of the dependent variable is changing according to the value of the independent variable.

Two types of Linear regressions:

1. Simple Linear Regression
2. Multiple Linear Regression

The strength of the linear regression model can be assessed using 2 metrics:

R^2 or Coefficient of Determination and Residual Standard Error (RSE) R^2 is a number which explains what portion of the given data variation is explained by the developed model. It always takes a value between 0 & 1. $R^2 = 1 - (RSS / TSS)$ Where, RSS: Residual Sum of squares,

TSS: Sum of errors of data from mean Assumptions of

simple linear regression are:

1. Linear relationship between X and Y
2. Error terms will be normally distributed (not X, Y)
3. Error terms have constant variance (homoscedasticity) Hypothesis testing of the co-efficients: The p-value is used to determine the significance of the variable.
4. Error terms are independent of each other

Parameters to assess a model are:

1. t-statistic: Used to determine the p-value and hence, helps in determining whether the coefficient is significant or not.
2. F statistic: Used to assess whether the overall model fit is significant or not. The higher the value of F statistic, the more significant a model turns out to be.
3. R-squared: After it has been concluded that the model fit is significant, the R-squared value tells the extent of the fit.

Multiple Linear Regression is used to understand the relationship between one dependent variable and several independent variables.

Multicollinearity: Multicollinearity is the effect of having related predictors in the multiple linear regression model that is one variable related to more than one other independent variable. It can be determined by the correlation matrix, VIF (Variance Inflation Factor). Higher the VIF value greater is the correlation. The variables with high VIF can be dropped one by one to reduce the multicollinearity

Feature Scaling: The independent variables in a model, might be on different scales which will lead a model to result in unexpected coefficients that might be difficult to interpret ex area for multiple variables. So we use the Standardizing or the Min Max Scaling approach to scale the variables.

Creating dummy variable for the categorical variables is essential to build a good model. For categorical variables with 'n' levels, you create 'n-1' new columns each indicating whether that level exists or not using a zero or one.

Adjusted R-sq penalizes the models based on the no. of variables present in it. If a variable is added and the adjusted R-sq drops then that variable is insignificant for the model and it has to be dropped.

One of the automated options of Feature selection is Recursive Feature Elimination (RFE). In this approach, the RFE is applied repeatedly on the train dataset to identify the best model.

In the process the variables with high p-value and high VIF values are dropped one by one and we arrive at a final model. Then the test dataset is used to make the final predictions.

2. Explain the Anscombe's quartet in detail. (3 marks) 3.

Anscombe's Quartet can be defined as a group of four data sets which are nearly identical in simple descriptive statistics, but there are some peculiarities in the dataset that fools the regression model if built. They have very different distributions and appear differently when plotted on scatter plots.

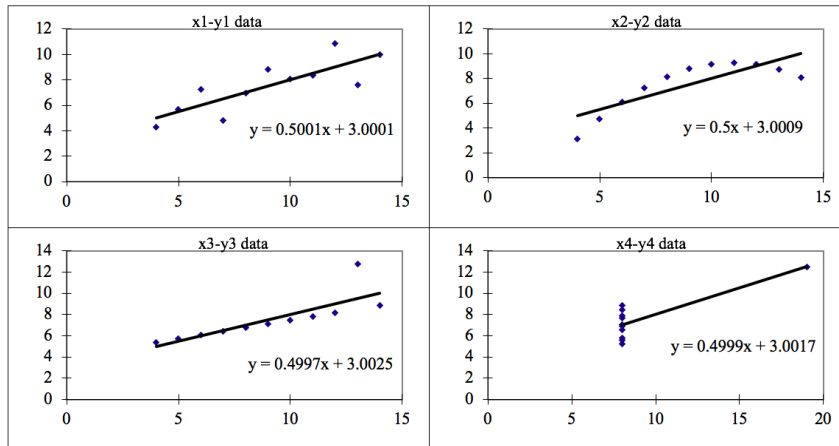
This tells us about the importance of visualising the data before applying various algorithms out there to build models out of them which suggests that the data features must be plotted in order to see the distribution of the samples that can help you identify the various anomalies present in the data like outliers, diversity of the data, linear separability of the data, etc. Also, the Linear Regression can be only be considered a fit for the data with linear relationships and is incapable of handling any other kind of datasets. These four plots can be defined as follows:

Anscombe's Data											
Observation	x1	y1		x2	y2		x3	y3		x4	y4
1	10	8.04		10	9.14		10	7.46		8	6.58
2	8	6.95		8	8.14		8	6.77		8	5.76
3	13	7.58		13	8.74		13	12.74		8	7.71
4	9	8.81		9	8.77		9	7.11		8	8.84
5	11	8.33		11	9.26		11	7.81		8	8.47
6	14	9.96		14	8.1		14	8.84		8	7.04
7	6	7.24		6	6.13		6	6.08		8	5.25
8	4	4.26		4	3.1		4	5.39		19	12.5
9	12	10.84		12	9.13		12	8.15		8	5.56
10	7	4.82		7	7.26		7	6.42		8	7.91
11	5	5.68		5	4.74		5	5.73		8	6.89

The statistical information for all these four datasets are approximately similar and can be computed as follows:

Anscombe's Data											
Observation	x1	y1		x2	y2		x3	y3		x4	y4
1	10	8.04		10	9.14		10	7.46		8	6.58
2	8	6.95		8	8.14		8	6.77		8	5.76
3	13	7.58		13	8.74		13	12.74		8	7.71
4	9	8.81		9	8.77		9	7.11		8	8.84
5	11	8.33		11	9.26		11	7.81		8	8.47
6	14	9.96		14	8.1		14	8.84		8	7.04
7	6	7.24		6	6.13		6	6.08		8	5.25
8	4	4.26		4	3.1		4	5.39		19	12.5
9	12	10.84		12	9.13		12	8.15		8	5.56
10	7	4.82		7	7.26		7	6.42		8	7.91
11	5	5.68		5	4.74		5	5.73		8	6.89
				Summary Statistics							
N	11	11		11	11		11	11		11	11
mean	9.00	7.50		9.00	7.500909		9.00	7.50		9.00	7.50
SD	3.16	1.94		3.16	1.94		3.16	1.94		3.16	1.94
r	0.82			0.82			0.82			0.82	

When these models are plotted on a scatter plot, all datasets generates a different kind of plot that is not interpretable by any regression algorithm which is fooled by these peculiarities and can be seen as follows:



The four datasets can be described as:

Dataset 1: this fits the linear regression model pretty well.

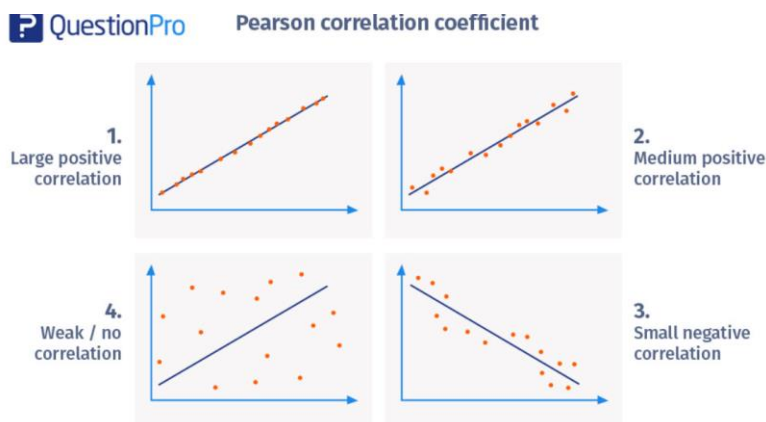
Dataset 2: this could not fit linear regression model on the data quite well as the data is non-linear.

Dataset 3: shows the outliers involved in the dataset which cannot be handled by linear regression model

Dataset 4: shows the outliers involved in the dataset which cannot be handled by linear regression model

3. What is Pearson's R? (3 marks)

The Pearson product-moment correlation coefficient (or Pearson correlation coefficient, for short) is a measure of the strength of a linear association between two variables and is denoted by r . The Pearson correlation coefficient, r , can take a range of values from +1 to -1. A value of 0 indicates that there is no association between the two variables. A value greater than 0 indicates a positive correlation and a value lesser than 0 indicates negative correlation. The below screenshot shows the Pearson correlation coefficient value of two variables and the strength of the associated variables



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

Feature scaling refers to the process of converting the variable into the same range of values. The independent variables in a model, might be on different scales which will lead a model to result in unexpected coefficients that might be difficult to interpret.

There are two types:

Normalized scaling: The feature variables are mapped to a minimum value of 0 and a maximum value of 1.

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

Standardized scaling: The feature variables are not set to 0 or 1 but are calculated based on the formula with the mean of 0 and standard deviation of 1.

$$Z = \frac{x - \mu}{\sigma}$$

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

The VIF value is infinite for the variables that have a perfect correlation. For such variables the R^2 value will be 1 which will lead to $1/(1-R^2)$.

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

Q-Q plots are also known as Quantile-Quantile plots, they plot the quantiles of a sample distribution against quantiles of a theoretical distribution. It helps us determine if a dataset follows any particular type of probability distribution like normal, uniform, exponential. In probability distributions, the data is represented in charts where the x-axis represents the possible values of the sample and the y-axis represents the probability of occurrence. The machine learning models are built based on the distributions which helps us achieve the best models.

The Q-Q plots are used to help us understand the data visually and is useful to determine the following:

- If two populations are of the same distribution.
- If residuals follow a normal distribution. Having a normal error term is an assumption in regression and we can verify if it's met using this.
- Skewness of distribution Q-Q plot for a normal distribution is as below

In statistics, a Q-Q (quantile-quantile) plot is a probability plot, which is a graphical method for comparing two probability distributions by plotting their quantiles against each other.[1] First,

the set of intervals for the quantiles is chosen. A point (x, y) on the plot corresponds to one of the quantiles of the second distribution (y -coordinate) plotted against the same quantile of the first distribution (x -coordinate). Thus the line is a parametric curve with the parameter which is the number of the interval for the quantile.

