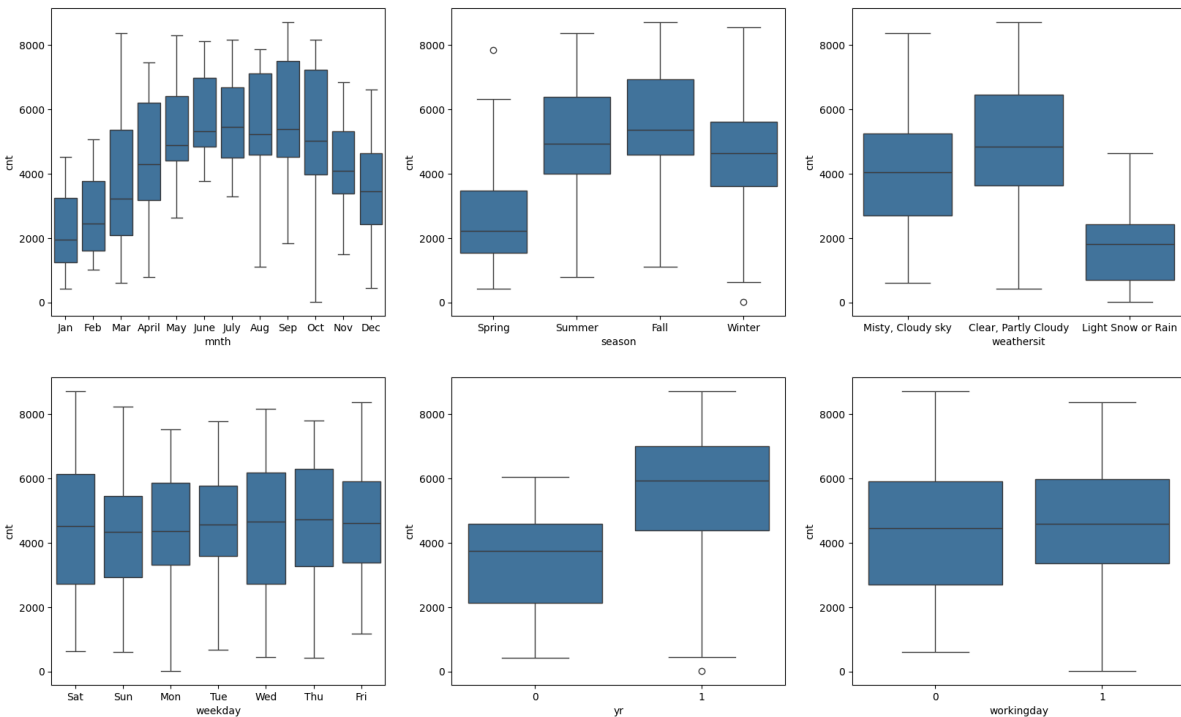


Assignment-based Subjective Questions

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

- There were more users when the weather was clear, average under cloudy conditions and least when it rains.
- The users signed up / used has increased from 2018 - 2019
- The demand was least in spring, this can be seen from the month wise data as well.



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

- This is to not add the columns that do not offer additional information. For example, if there are 4 seasons, Spring, Summer, Fall and Winter, if it is not Summer / Fall or Spring, it must be Winter so we do not need to waste more space conveying the same message.

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

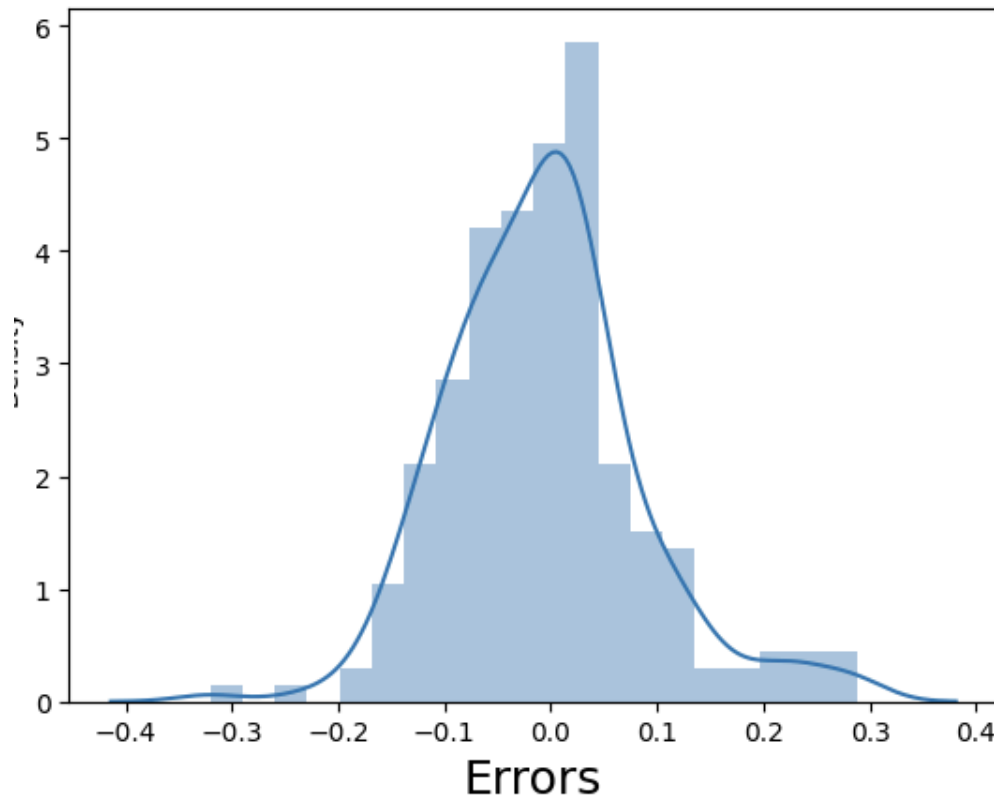
Tmp and Atemp has the highest corelation, if we remove the obvious casual and registered since cnt is sum of them.

yr	1	0.0082	-0.0029	0.049	0.047	-0.11	-0.012	0.25	0.6	0.57	4.1e-17	3.7e-17	2.8e-16	5.5e-17	1.5e-16	2.1e-16	1.6e-16	7.7e-17	3e-17	8.6e-17	2.1e-17	3.2e-16	1.8e-17	1.1e-16	-0.074	-0.0058	0.0039	-0.0039	0.0039	-1.1e-16	-2e-16	-0.0039
holiday	-0.0082	1	-0.25	-0.029	-0.033	-0.016	0.0063	0.054	-0.11	-0.069	-0.052	0.0064	0.012	0.036	0.0064	-0.052	-0.052	0.0064	0.068	0.0064	0.0082	0.035	-0.024	0.017	-0.03	-0.019	0.28	-0.071	-0.071	-0.023	-0.047	-0.046
workingday	-0.0029	-0.25	1	0.053	0.053	0.023	-0.019	-0.52	0.31	0.063	0.038	-0.015	-0.0031	0.025	-0.015	0.021	0.028	0.0065	-0.011	-0.004	-0.011	-0.028	0.015	-0.0046	0.029	0.049	0.15	-0.6	-0.6	0.26	0.27	0.27
temp	0.049	-0.029	0.053	1	0.99	0.13	-0.16	0.54	0.54	0.63	0.35	-0.29	-0.31	-0.43	0.43	0.31	-0.17	0.17	-0.21	-0.018	0.2	-0.62	0.15	-0.23	-0.058	-0.096	-0.0048	-0.03	-0.027	0.019	0.019	0.023
atemp	0.047	-0.033	0.053	0.99	1	0.14	-0.18	0.54	0.54	0.63	0.33	-0.28	-0.31	-0.44	0.43	0.3	-0.17	0.17	-0.2	-0.0048	0.19	-0.62	0.16	-0.21	-0.065	-0.093	9.1e-05	-0.031	-0.023	0.02	0.022	0.021
hum	-0.11	-0.016	0.023	0.13	0.14	1	-0.25	-0.075	-0.089	-0.099	0.022	0.082	-0.13	-0.09	-0.064	-0.11	-0.084	0.13	-0.0058	0.14	0.18	-0.18	-0.0029	0.16	0.27	0.49	0.029	-0.023	2e-05	-0.052	0.041	0.046
windspeed	-0.012	0.0063	-0.019	-0.16	-0.18	-0.25	1	-0.17	-0.22	-0.24	-0.069	-0.055	0.096	0.062	-0.096	-0.02	0.13	-0.03	-0.026	-0.06	-0.095	0.18	0.097	-0.14	0.12	-0.037	0.001	0.032	-0.01	0.0058	0.007	-0.014
casual	0.25	0.054	-0.52	0.54	0.54	-0.075	-0.17	1	0.39	0.67	0.14	-0.22	-0.25	-0.29	0.18	0.17	-0.059	0.16	-0.1	0.051	0.14	-0.43	0.22	-0.099	-0.17	-0.17	-0.1	0.37	0.29	-0.15	-0.17	-0.17
registered	0.6	-0.11	0.31	0.54	0.54	-0.089	-0.22	0.39	1	0.95	0.16	-0.12	-0.23	-0.33	0.13	0.17	-0.13	0.093	-0.0042	0.11	0.18	-0.51	0.085	0.12	-0.23	-0.14	0.0014	-0.15	-0.2	0.11	0.077	0.094
cnt	0.57	-0.069	0.063	0.63	0.63	-0.099	-0.24	0.67	0.95	1	0.18	-0.17	-0.27	-0.37	0.17	0.2	-0.13	0.13	-0.04	0.11	0.19	-0.56	0.15	0.065	-0.24	-0.17	-0.036	0.009	-0.059	0.034	0.00056	0.014
Aug	-4.1e-17	-0.052	0.038	-0.35	0.33	0.022	-0.069	0.14	0.16	0.18	1	0.093	-0.088	-0.093	-0.093	-0.091	-0.093	-0.093	-0.091	-0.093	-0.091	-0.17	-0.18	-0.17	-0.052	-0.02	0.0012	-0.013	-0.013	0.0023	0.0023	0.018
Dec	-3.7e-17	0.0064	-0.015	-0.29	-0.28	0.082	-0.055	-0.22	-0.12	-0.17	-0.093	1	-0.088	-0.093	-0.093	-0.091	-0.093	-0.093	-0.091	-0.093	-0.091	0.077	-0.18	0.28	0.036	0.063	0.0012	0.015	0.0012	0.0023	-0.012	-0.011
Feb	-2.8e-16	0.012	-0.0031	-0.31	-0.31	-0.13	0.096	-0.25	-0.23	-0.27	-0.088	-0.088	1	-0.088	-0.088	-0.086	-0.088	-0.088	-0.086	-0.088	-0.086	0.5	-0.17	-0.16	-0.019	-0.0095	-0.0008	-0.0008	-0.0008	0.00032	0.00032	0.0015
Jan	-5.5e-17	0.036	-0.025	-0.43	-0.44	-0.09	0.062	-0.29	-0.33	-0.37	-0.093	-0.093	-0.088	1	-0.093	-0.091	-0.093	-0.093	-0.091	-0.093	-0.091	0.53	-0.18	-0.17	-0.023	0.032	0.015	0.0012	0.015	-0.012	0.0023	-0.011
July	-1.5e-16	0.0064	-0.015	0.43	0.43	-0.064	-0.096	0.18	0.13	0.17	-0.093	-0.093	-0.088	-0.093	1	-0.091	-0.093	-0.093	-0.091	-0.093	-0.091	-0.17	-0.18	-0.17	-0.023	-0.11	0.0012	0.0012	0.015	-0.012	0.0023	-0.011
June	-2.1e-16	-0.052	0.021	0.31	0.3	-0.11	-0.02	0.17	0.17	0.2	-0.091	-0.091	-0.086	-0.091	-0.091	1	-0.091	-0.091	-0.09	-0.091	-0.09	-0.17	0.29	-0.17	-0.052	-0.066	-0.009	0.0053	-0.009	0.0065	-0.0078	0.0077
Mar	-1.6e-16	-0.052	0.028	-0.17	-0.17	-0.084	0.13	-0.059	-0.13	-0.13	-0.093	-0.093	-0.088	-0.093	-0.093	-0.091	1	-0.093	-0.091	-0.093	-0.091	0.28	0.072	-0.17	0.0064	0.032	-0.013	0.0012	-0.013	0.016	0.0023	0.0036
May	-7.7e-17	0.0064	0.0065	0.17	0.17	0.13	-0.03	0.16	0.093	0.13	-0.093	-0.093	-0.088	-0.093	-0.093	-0.091	-0.093	1	-0.091	-0.093	-0.091	-0.17	0.52	-0.17	-0.052	0.043	0.0012	-0.013	0.0012	0.0023	0.016	0.0036
Nov	-3e-17	0.068	-0.011	-0.21	-0.2	-0.0058	-0.026	-0.1	-0.0042	0.04	-0.091	-0.091	-0.086	-0.091	-0.091	-0.09	-0.091	-0.091	1	-0.091	-0.09	-0.17	-0.17	0.53	0.038	-0.055	-0.009	-0.009	-0.009	0.0065	0.0065	0.0077
Oct	-8.6e-17	0.0064	-0.004	0.018	0.0048	0.14	-0.06	0.051	0.11	0.11	-0.093	-0.093	-0.088	-0.093	-0.093	-0.091	-0.093	-0.093	-0.091	1	-0.091	-0.17	-0.18	0.54	0.095	0.043	0.015	0.0012	0.0012	-0.012	0.0023	0.0036
Sep	-2.1e-17	0.0082	-0.011	0.2	0.19	0.18	-0.095	0.14	0.18	0.19	-0.091	-0.091	-0.086	-0.091	-0.091	-0.09	-0.091	-0.091	-0.09	-0.091	1	-0.17	-0.17	0.016	0.038	0.04	-0.009	0.0053	0.0053	0.0065	-0.0078	-0.0067
Spring	-3.2e-16	0.035	-0.028	-0.62	-0.62	-0.18	0.18	-0.43	-0.51	-0.56	-0.17	0.077	0.5	0.53	-0.17	-0.17	0.28	-0.17	-0.17	-0.17	-0.17	1	-0.33	-0.32	-0.022	0.029	0.00099	0.01	0.01	-0.0059	-0.0059	-0.013
Summer	-1.8e-17	-0.024	0.015	0.15	0.16	-0.0029	0.097	0.22	0.085	0.15	-0.18	-0.18	-0.17	-0.18	-0.18	0.29	0.072	0.52	-0.17	-0.18	-0.17	-0.33	1	-0.33	-0.043	0.04	0.0048	-0.0042	-0.0042	-0.0019	-0.0019	0.0094
Winter	-1.1e-16	0.017	-0.0046	-0.23	-0.21	0.16	-0.14	-0.099	0.12	0.065	-0.17	0.28	-0.16	-0.17	-0.17	-0.17	-0.17	-0.17	0.53	0.54	0.016	-0.32	-0.33	1	0.093	0.034	0.0036	-0.0055	0.0036	-0.0033	0.0059	-0.0011
Light Snow or Rain	-0.074	-0.03	0.029	-0.058	-0.065	0.27	0.12	-0.17	-0.23	-0.24	-0.052	0.036	-0.019	-0.023	-0.023	-0.052	0.0064	-0.052	0.038	0.095	0.038	-0.022	-0.043	0.093	1	-0.12	-0.024	0.023	-0.047	0.00019	0.024	0.095
Misty, Cloudy sky	-0.0058	-0.019	0.049	-0.096	-0.093	0.49	0.037	-0.17	-0.14	-0.17	-0.02	0.063	-0.0095	0.032	-0.11	-0.066	0.032	0.043	-0.055	0.043	0.04	0.029	0.04	0.034	-0.12	1	0.013	-0.011	-0.044	-0.0087	0.024	-0.023
Mon	-0.0039	0.28	0.15	-0.0048	9.1e-05	0.029	0.001	-0.1	0.0014	-0.036	0.0012	0.0012	-0.0008	0.015	0.0012	-0.009	-0.013	0.0012	-0.009	0.015	-0.009	0.00099	0.0048	0.0036	-0.024	0.013	1	-0.17	-0.17	-0.17	-0.17	-0.17
Sat	-0.0039	-0.071	-0.6	-0.03	-0.031	-0.023	0.032	0.37	-0.15	0.009	-0.013	0.015	-0.0008	0.0012	0.0012	0.0053	0.0012	-0.013	-0.009	0.0012	0.0053	0.01	-0.0042	-0.0055	0.023	-0.011	-0.17	1	-0.17	-0.17	-0.17	-0.17
Sun	-0.0039	-0.071	-0.6	-0.027	-0.023	2e-05	-0.01	0.29	-0.2	-0.059	-0.013	0.0012	-0.0008	0.015	0.015	-0.009	-0.013	0.0012	-0.009	0.0012	0.0053	0.01	-0.0042	0.0036	-0.047	-0.044	-0.17	-0.17	1	-0.17	-0.17	-0.17
Thu	-1.1e-16	-0.023	0.26	0.019	0.02	0.052	0.0058	-0.15	0.11	0.034	0.0023	0.0023	0.00032	-0.012	-0.012	0.0065	0.016	0.0023	0.0065	-0.012	0.0065	-0.0059	-0.0019	-0.0033	0.00019	-0.0087	-0.17	-0.17	-0.17	1	-0.17	-0.17
Tue	-2e-16	-0.047	0.27	0.019	0.022	0.041	0.007	-0.17	0.077	0.00056	0.0023	-0.012	0.00032	0.0023	0.0023	-0.0078	0.0023	0.016	0.0065	0.0023	-0.0078	-0.0059	-0.0019	0.0059	0.024	0.024	-0.17	-0.17	-0.17	-0.17	1	-0.17
Wed	-0.0039	-0.046	0.27	0.023	0.021	0.046	-0.014	-0.17	0.094	0.014	0.018	-0.011	0.0015	-0.011	-0.011	0.0077	0.0036	0.0036	0.0077	0.0036	-0.0067	-0.013	0.0094	-0.0011	0.095	-0.023	-0.17	-0.17	-0.17	-0.17	-0.17	1

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

- The error was normally distributed with mean around zero.

Error Terms



- R-squared of the predicted values was 0.82 which is decent.

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

- Temperature had positive correlation
- Weather - especially if it was Raining / Snowing had a negative correlation
- Windspeed also had a positive correlation if we do not consider 'year'.

General Subjective Questions

1. Explain the linear regression algorithm in detail.

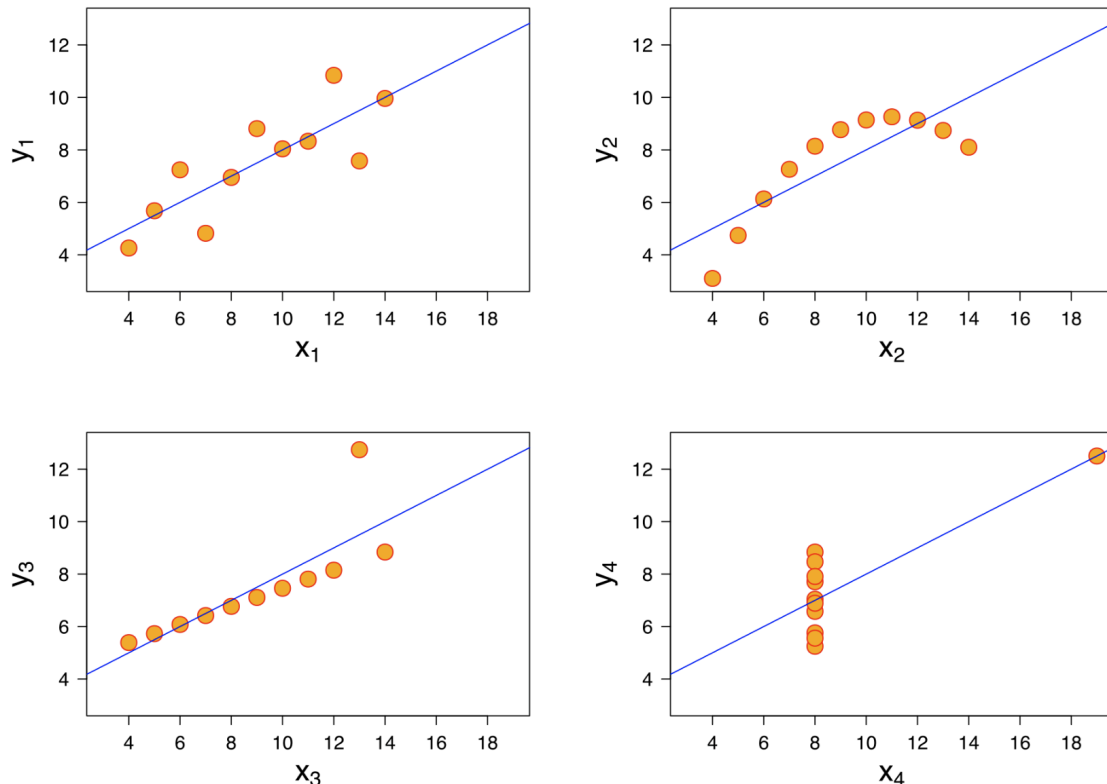
Linear regression is a statistical model which estimates the linear relationship between a scalar response and one or more explanatory variables. The algorithm tries to fit the

values represented by a straight line. There are two types Simple and Multiple Linear Regression. It involves following steps.

- Reading and understanding the data
- Visualizing the data to identify relationships
- Data cleansing and preparation
- Split to training / test set
- Train the model.
- Evaluate the model.

2. Explain the Anscombe's quartet in detail.

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. Each dataset consists of eleven (x, y) points.



3. What is Pearson's R?

Pearson's R, also called the Pearson correlation coefficient, is a numerical value that is used to assess the **linear relationship** between two continuous variables. It's a number between -1 and 1. Positive number mean a positive correlation. The closeness to 1 indicates the strength of correlation.

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

Scaling:

Scaling is a process of bring the values of all the variables in to a common range.

Why is Scaling performed:

This is to improve convergence of the data and provide equal footing to avoid larger valued data from influencing the mode.

Standardized vs Normalized:

Normalization scales the features to a range between 0 and 1 (or sometimes -1 and 1), using the min and max values.

Standardization scales the features by subtracting the mean and then dividing by the standard deviation.

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

This can happen when there is a perfect multicollinearity. Since VIF is calculated based on R-squared, if R-squared is 1, VIF can be infinity.

It is most likely because of unreliable coefficients or misleading interpretations. For example we might have assumed two variables as independent when they were not.

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

Q-Q plot(Quantile-Quantile plot) is a tool used to assess the distribution in linear regression. It helps us visually compare the distribution of the residuals (errors) in the model. 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). This defines a parametric curve where the parameter is the index of the quantile interval.