# FIT5196 Reflective Report in Assessment 2

Student Name: Anushka Jemima

Student ID: 33617457

Student Name: Shruthi Shashidhara Shastry

Student ID: 33684669

Date: 18-10-2024

1. Dirty Data

**Feedback provided :**
   a. Order of finding and fixing anomalies is important.
   b. IS_EXPIDITED column is related to delivery charge.
   c. There can be more than one anomaly present in a column.

**Improvements after Feedback :**

After the feedback session we began examining the order in which we had discovered and fixed anomalies, we noticed for season and nearest warehouses the only anomaly was not the capitalization issue but there was mismatch between the season and months , and the nearest warehouse based on minimum distance.

We also ensured that we are first fixing the customer latitude and longitude column anomaly before the nearest warehouse as this caused 2 anomalies in the same row if we did the nearest warehouse first and then the customer latitude and longitude , to track this we added an index during exploration and maintained an excel file to ensure we have only one anomaly per row.



For is_expidited we decided to leverage the linear dependence of the is_expidited column on delivery charge and built a model to decide at which threshold(for which rows) the is_expidited column should be updated, and used R squared to evaluate the accuracy.

2. Outlier Data

**Feedback provided:**

    a.  The model needs to be checked for accuracy after dropping the outliers.

**Improvements after Feedback :**

We dropped the records that were detected as outliers, and checked the accuracy of our model, through this we were able to alter our outlier detection metric , i.e we had use the 3 sigma rule and in order to decide a threshold for outliers we had attempted different values for the confidence interval , since when 3 sigma was giving us a very low accuracy , not all outliers had been dropped we attempted a 2 standard deviation and this gave us a better R squared value.

3. Task 2:

    **Feedback provided:**

    a.  Research if the target variable needs to be transformed.

    **Improvements after Feedback:**

On performing research through various online posts and papers we were able to draw a conclusion that target variables and not usually transformed. This is because we want to maintain the distribution and we want to avoid interpretability issues once predictions are made.