

# Procurement Purchase Prediction

|                   |   |
|-------------------|---|
| Business Use case | To identify if a purchase order would be cancelled or not |
| Presented by      | Madhu   |

**Meta data**

RangeIndex: 85759 entries, 0 to 85758

Data columns (total 19 columns):

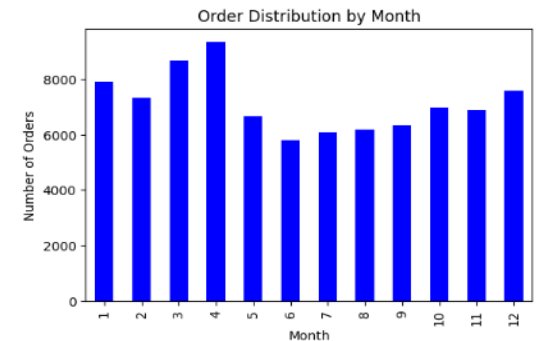
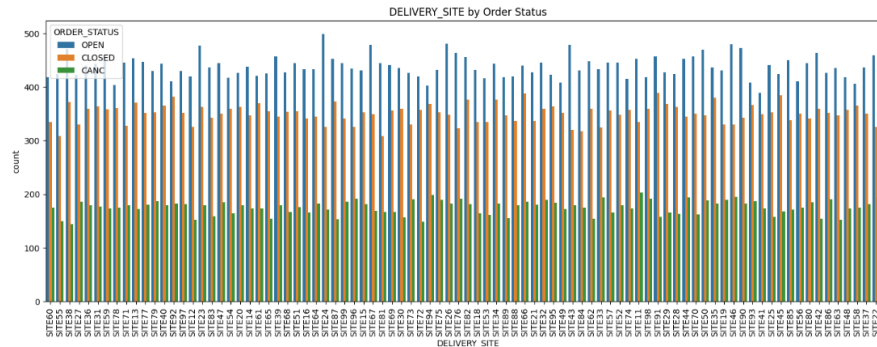
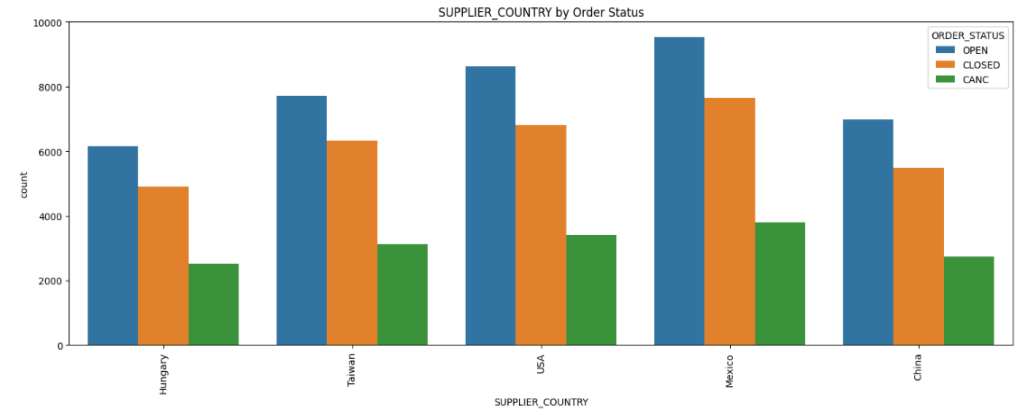
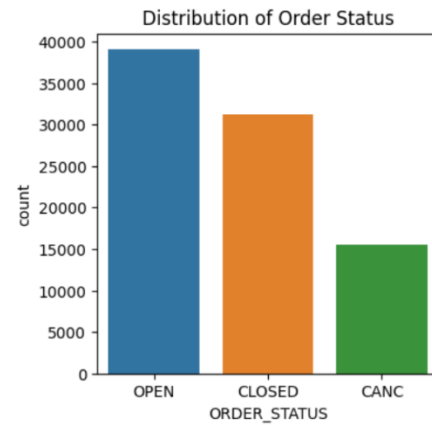
ORDER\_ID, ORDER\_DATE, NEED\_DATE, ORDER\_STATUS, ORDER\_CHANNEL, DELY\_QTY, DELY\_DATE, DELIVERY\_SITE, BUYER\_ID, BUYER\_NAME, SUPPLIER\_ID, SUPPLIER\_NAME, SUPPLIER\_COUNTRY, ITEM\_ID, ORDER\_QUANTITY, ORDER\_UNIT\_PRICE, ORDER\_COST, ORDER\_CURRENCY\_CODE

Dtypes: datetime64[ns](3), float64(3), int64(1), object(12)

|                |           |       |
|----------------|-----------|-------|
| Missing Values | DELY_QTY  | 54590 |
|                | DELY_DATE | 54590 |

- General Observation**
1. Target Variable: ORDER\_STATUS (Canc/open/closed)
  2. Whenever the order is closed delivery quantity and delivery date are available.
  3. ORDER\_ID is a unique identifier

## Exploratory Data Analysis



# Exploratory Data Analysis- Buyer and Supplier

Top 5 Buyers based on total amount spent:

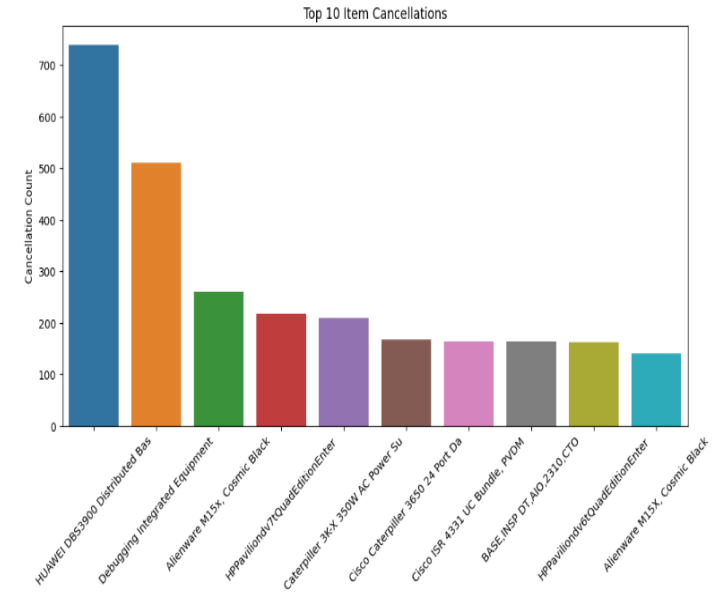
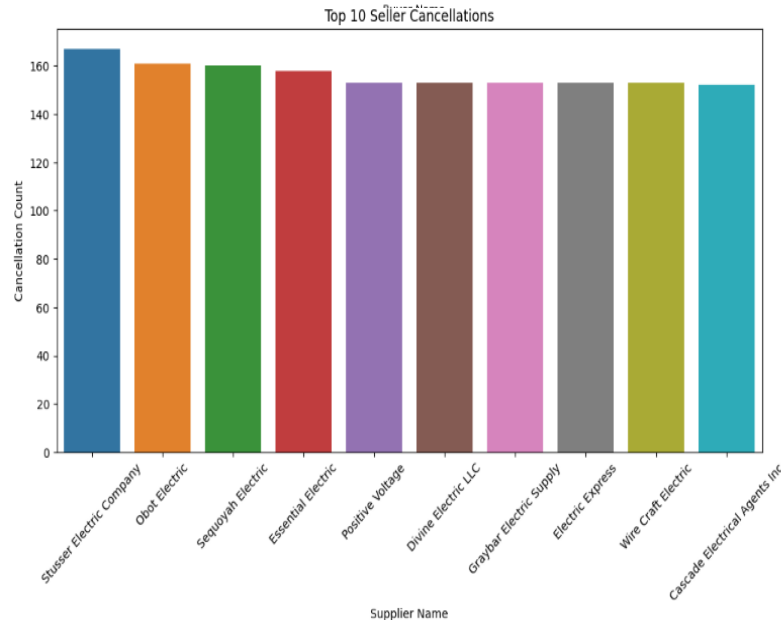
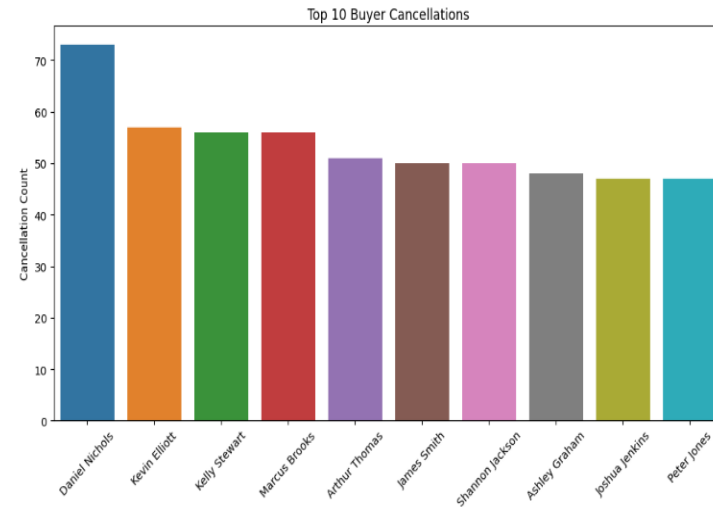
| BUYER_NAME       | Amount Spent |
|------------------|--------------|
| Shannon Jackson  | 2.730083e+08 |
| Daniel Nichols   | 2.566291e+08 |
| Michael Callahan | 1.964493e+08 |
| Brian Green      | 1.919278e+08 |
| Scott Bennett    | 1.911048e+08 |

Name: ORDER\_COST, dtype: float64

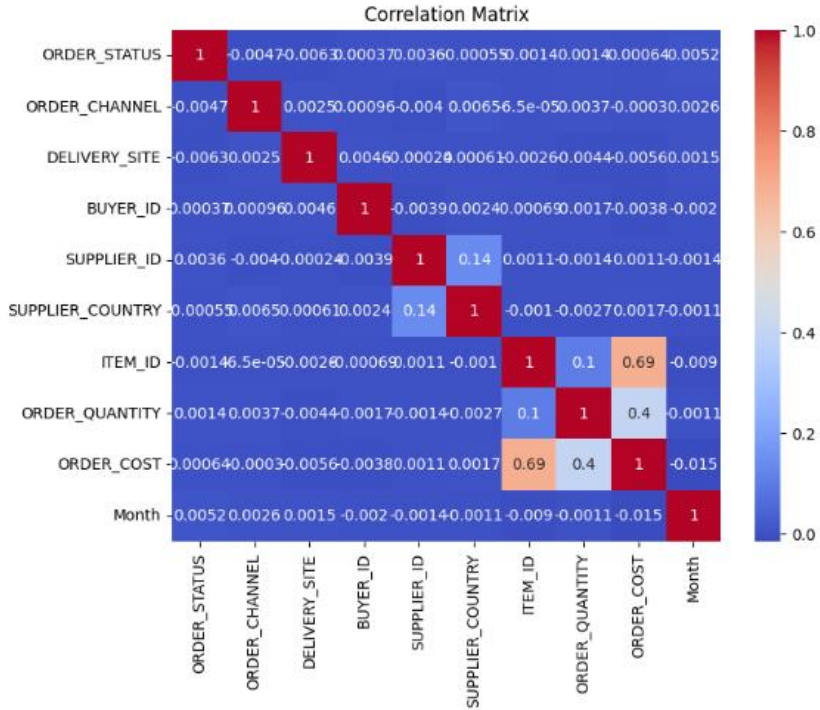
Top 5 Suppliers based on total amount received:

| SUPPLIER_NAME                | Amount Received |
|------------------------------|-----------------|
| EC Electric                  | 6.470283e+08    |
| Divine Electric LLC          | 6.173887e+08    |
| Electric Express             | 6.170383e+08    |
| NW Wind & Solar              | 6.106914e+08    |
| VECA Electric & Technologies | 5.986096e+08    |

Name: ORDER\_COST, dtype: float64



# EDA - Feature Engineering



A correlation coefficient of 0.69 indicates a strong positive relationship between ITEM\_ID and ORDER\_COST. This means that as one variable increases, the other variable tends to increase as well.

Chi-square test of independence for DELIVERY\_SITE:  
Chi2 value: 173.65447662614497  
P-value: 0.5358539431319053

Chi-square test of independence for SUPPLIER\_COUNTRY:  
Chi2 value: 6.677618438268236  
P-value: 0.5717804797175376

Chi-square test of independence for ITEM\_ID:  
Chi2 value: 524.7806142430564  
P-value: 0.7553159447592004

Chi-square test of independence for ORDER\_STATUS:  
Chi2 value: 171518.0  
P-value: 0.0

Chi-square test of independence for ORDER\_CHANNEL:  
Chi2 value: 4.4797719858775755  
P-value: 0.6120390995202432

Chi-square test of independence for BUYER\_ID:  
Chi2 value: 837.4504216664104  
P-value: 0.9259634592102286

Chi-square test of independence for SUPPLIER\_ID:  
Chi2 value: 206.37796492029145  
P-value: 0.8210821825390068

## The interpretation of the chi-square test results : (Categorical Variables)

DELIVERY\_SITE: Chi2 value: 173.65 P-value: 0.54 The chi-square test suggests that there is no significant association between the delivery site and the order status. The p-value is greater than the significance level of 0.05, indicating that the variables are independent.

SUPPLIER\_COUNTRY: Chi2 value: 6.68 P-value: 0.57 The chi-square test suggests that there is no significant association between the supplier country and the order status. The p-value is greater than 0.05, indicating independence between the variables.

ITEM\_ID: Chi2 value: 524.78 P-value: 0.76 The chi-square test suggests that there is no significant association between the item ID and the order status. The p-value is greater than 0.05, indicating independence.

ORDER\_CHANNEL: Chi2 value: 4.48 P-value: 0.61 The chi-square test suggests that there is no significant association between the order channel and the order status. The p-value is greater than 0.05, indicating independence.

BUYER\_ID: Chi2 value: 837.45 P-value: 0.93 The chi-square test suggests that there is no significant association between the buyer ID and the order status. The p-value is greater than 0.05, indicating independence.

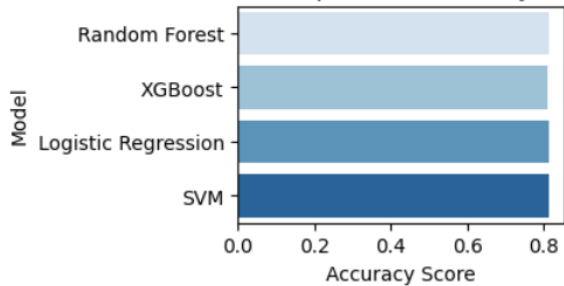
SUPPLIER\_ID: Chi2 value: 206.38 P-value: 0.82 The chi-square test suggests that there is no significant association between the supplier ID and the order status. The p-value is greater than 0.05, indicating independence.

# Model Insights and Analysis

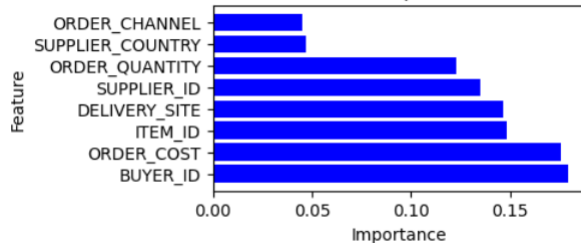
## Insignificant Features

'ORDER\_ID', 'NEED\_DATE', 'DELY\_DATE', 'ORDER\_CURRENCY\_CODE', 'ORDER\_DATE', 'ORDER\_UNIT\_PRICE', 'ITEM\_DESC', 'SUPPLIER\_NAME', 'BUYER\_NAME', 'DELY\_QTY'

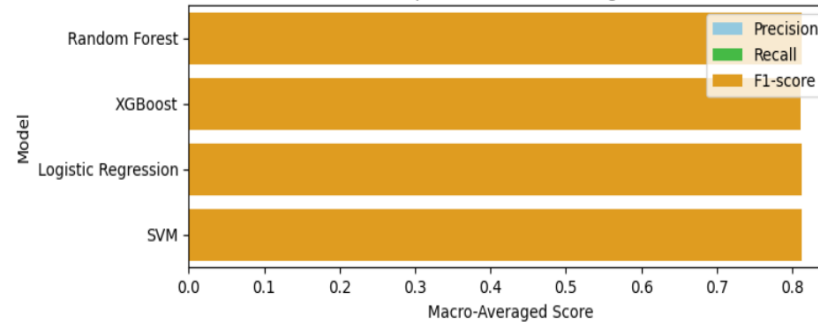
### Model Comparison - Accuracy Scores



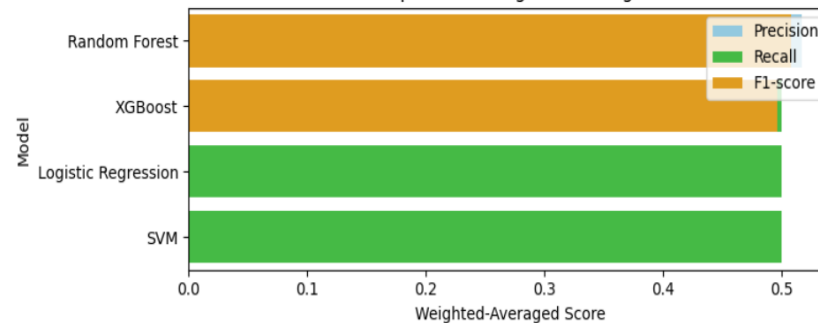
### Feature Importances



### Model Comparison - Macro-Averaged Scores



### Model Comparison - Weighted-Averaged Scores



- ✓ Experimented Random Forest, XGBoost, Logreg and SVM models.
- ✓ All the models achieves an overall accuracy of 0.81, indicating decent overall performance.
- ✓ However, when considering the performance across individual classes, the model's performance is relatively poor in terms of precision, recall, and F1-score (as indicated by macro-average).
- ✓ The weighted-average scores provide a more representative evaluation, accounting for class imbalance, and show better performance compared to the macro-average scores.
- ✓ Of these Random Forest has a better Macro avg.

# Summary

|                            |   |
|----------------------------|---|
| <b>EDA</b>                 | <ul style="list-style-type: none"><li>✓ The p-value is greater than the significance level of 0.05, indicating that the variables DELIVERY_SITE, SUPPLIER_COUNTRY, ITEM_ID, BUYER_ID, SUPPLIER_ID are independent.</li><li>✓ There is positive relationship between ITEM_ID and ORDER_COST</li></ul>  |
| <b>Feature Engineering</b> | <ul style="list-style-type: none"><li>✓ Label Encoding is performed for Categorical variables 'ORDER_CHANNEL', 'DELIVERY_SITE', 'BUYER_ID', 'SUPPLIER_ID', 'SUPPLIER_COUNTRY', 'ITEM_ID'</li><li>✓ Selection Bias - Insufficient data for the class : Cancel; To overcome the same, Oversampling and downsampling is performed.</li></ul>   |
| <b>Models</b>              | <ul style="list-style-type: none"><li>✓ Removed Buyer name, Seller name and Item description as they can be interpreted using their IDs</li><li>✓ Since it is a classification problem with known target variable Supervised learning is performed (insights are explained in the previous slide)</li><li>✓ Stratified sampling is performed.</li><li>✓ As all the models have given an accuracy of 81%, it is understood that they have learnt similar features.</li><li>✓ F1 score of 90 signifies a strong classification performance and indicates that the models are effectively capturing the underlying patterns and relationships in the data.</li><li>✓ Have built a binary classification model combining Closed and Open as one class and retained Cancel as the other class</li><li>✓ XGBoost is a powerful algorithm for handling imbalanced datasets, but the algorithm also has given an accuracy of 81%.</li></ul> |
| <b>Problem</b>             | <ul style="list-style-type: none"><li>✓ Classification with class imbalance</li><li>✓ The class imbalance issue is challenging because the model is struggling to learn patterns from the minority class due to the dominance of the majority class.</li></ul>  |
| <b>Next Steps</b>          | <ul style="list-style-type: none"><li>✓ Explore other algorithms: Consider trying other machine learning algorithms that are suitable for imbalanced datasets, such as neural networks.</li><li>✓ Other feature engineering techniques specifically designed for addressing class imbalance, such as SMOTE (Synthetic Minority Over-sampling Technique) or ADASYN (Adaptive Synthetic Sampling) can be studied</li><li>✓ More feature Engineering to create new features using domain knowledge.</li><li>✓ Hyperparameter Tuning like tuning max_depth, learning_rate, n_estimators etc</li><li>✓ Experiment ensemble such as bagging or boosting techniques like AdaBoost, CatBoost</li></ul>  |