# "Navigating Shipments: A Predictive Analysis of Delivery Success"

## Abstract

In response to the growing demand for advanced customer insights in international e-commerce, this project focuses on leveraging machine learning techniques to analyze product shipment tracking data for an electronic products company. The primary objective is to develop a classification model capable of predicting whether a product will reach its destination on time or encounter delays. By harnessing the power of predictive analytics, the study aims to provide valuable insights into the shipment process, enabling the company to optimize logistics operations and enhance customer satisfaction. Through the exploration of shipment data patterns and the implementation of cutting-edge classification algorithms, this project endeavors to equip the company with actionable intelligence for more efficient and reliable delivery management in the competitive electronic products market.

## Table of Contents

# 1. Introduction

a. **Background and Motivation:**
   The global e-commerce industry has witnessed exponential growth, fueled by technological advancements and changing consumer behaviors. This project delves into the realm of international e-commerce, where companies strive to leverage data-driven insights to enhance customer experiences and operational efficiency. The increasing complexity of supply chain logistics and customer expectations necessitates the adoption of advanced analytics techniques to remain competitive in the market.

b. **Objectives of the Study:**
   The primary aim of this study is to harness the power of machine learning and data analytics to extract valuable insights from the customer database of an international e-commerce company. By analyzing customer behavior patterns, purchase trends, and shipment tracking data, the study seeks to uncover actionable intelligence that can drive strategic decision-making and improve business outcomes.

c. **Importance of Shipment Tracking in E-commerce:**
   Shipment tracking plays an important role in the success of e-commerce operations, serving as a key determinant of customer satisfaction and loyalty. Timely and reliable delivery of products is crucial for maintaining trust and credibility with customers, particularly in the highly competitive online retail space. By implementing robust tracking mechanisms and predictive analytics models, e-commerce companies can optimize supply chain operations, reduce delivery risks, and enhance overall service quality.

d. **Challenges and Opportunities in Logistics Management:**
   Logistics management poses unique challenges and opportunities in the e-commerce landscape. Efficient inventory management, route optimization, and last-mile delivery logistics are critical components of successful e-commerce operations. Previous studies have highlighted challenges such as supply chain disruptions, delivery bottlenecks, and rising customer expectations. Understanding these challenges and exploring opportunities for improvement through advanced analytics and predictive modeling are essential objectives of this project.

# 2. Data Collection and Preprocessing

a. **Description of the Product Shipment Tracking Dataset:** This section provides an overview of the dataset used for the analysis, including the number of observations, variables, and their respective meanings. It outlines the structure of the dataset and explains the relevance of each variable in the context of shipment tracking and delivery prediction.

| | ID | Warehouse_block | Mode_of_Shipment | Customer_care_calls | Customer_rating | Cost_of_the_Product | Prior_purchases | Product_importance | Gender | Discount_offered | Weight_in_gms | Reached.on.Time_Y.N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | D | Flight | 4 | 2 | 177 | 3 | low | F | 44 | 1233 | 1 |
| 1 | 2 | F | Flight | 4 | 5 | 216 | 2 | low | M | 59 | 3088 | 1 |
| 2 | 3 | A | Flight | 2 | 2 | 183 | 4 | low | M | 48 | 3374 | 1 |
| 3 | 4 | B | Flight | 3 | 3 | 176 | 4 | medium | M | 10 | 1177 | 1 |
| 4 | 5 | C | Flight | 2 | 2 | 184 | 3 | medium | F | 46 | 2484 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 10994 | 10995 | A | Ship | 4 | 1 | 252 | 5 | medium | F | 1 | 1538 | 1 |
| 10995 | 10996 | B | Ship | 4 | 1 | 232 | 5 | medium | F | 6 | 1247 | 0 |
| 10996 | 10997 | C | Ship | 5 | 4 | 242 | 5 | low | F | 4 | 1155 | 0 |
| 10997 | 10998 | F | Ship | 5 | 2 | 223 | 6 | medium | M | 2 | 1210 | 0 |
| 10998 | 10999 | D | Ship | 2 | 5 | 155 | 5 | low | F | 6 | 1639 | 0 |

10999 rows × 12 columns

- **ID:** ID Number of Customers.
- **Warehouse block:** The Company has a big Warehouse which is divided into blocks such as A,B,C,D,E.
- **Mode of shipment:** The Company Ships the products in multiple ways such as Ship, Flight and Road.
- **Customer care calls:** The number of calls made from enquiry for enquiry of the shipment.
- **Customer rating:** The company has a rating from every customer. 1 is the lowest (Worst), 5 is the highest (Best).
- **Cost of the product:** Cost of the Product in US Dollars.
- **Prior purchases:** The Number of Prior Purchases.
- **Product importance:** The company has categorized the product in various parameters such as low, medium, high.
- **Gender:** Male and Female.
- **Discount offered:** Discount offered on that specific product.
- **Weight in gms:** It is the weight in grams.
- **Reached on time:** It is the target variable, where 1 Indicates that the product has NOT reached on time and 0 indicates it has reached on time.

```
[ ]  df["Reached.on.Time_Y.N"] = np.where(df["Reached.on.Time_Y.N"] == 1, 0, 1)
     df.drop(["ID"], axis=1, inplace=True)
```

Now, 1 Means Reached on time. 0 Means not reached on time

b. **Data Cleaning:** Discusses the steps taken to clean the dataset before analysis. This includes handling missing values, outlier handling.
- Null Treatment

## Checking Null Values

```
] df.isnull().sum()

Warehouse_block        0
Mode_of_Shipment       0
Customer_care_calls    0
Customer_rating        0
Cost_of_the_Product    0
Prior_purchases        0
Product_importance     0
Gender                 0
Discount_offered       0
Weight_in_gms          0
Reached.on.Time_Y.N    0
dtype: int64
```

There are no missing values in the dataset, which is excellent for analysis as it means no imputation is necessary.

- Outlier Treatment
  There can be potential outliers in the Numerical columns, We'll try outlier detection and removal techniques as part of next data cleaning step

- There are two columns Prior_purchases and Discount_offered has possible outliers.
- We'll use the IQR method to detect the outliers.
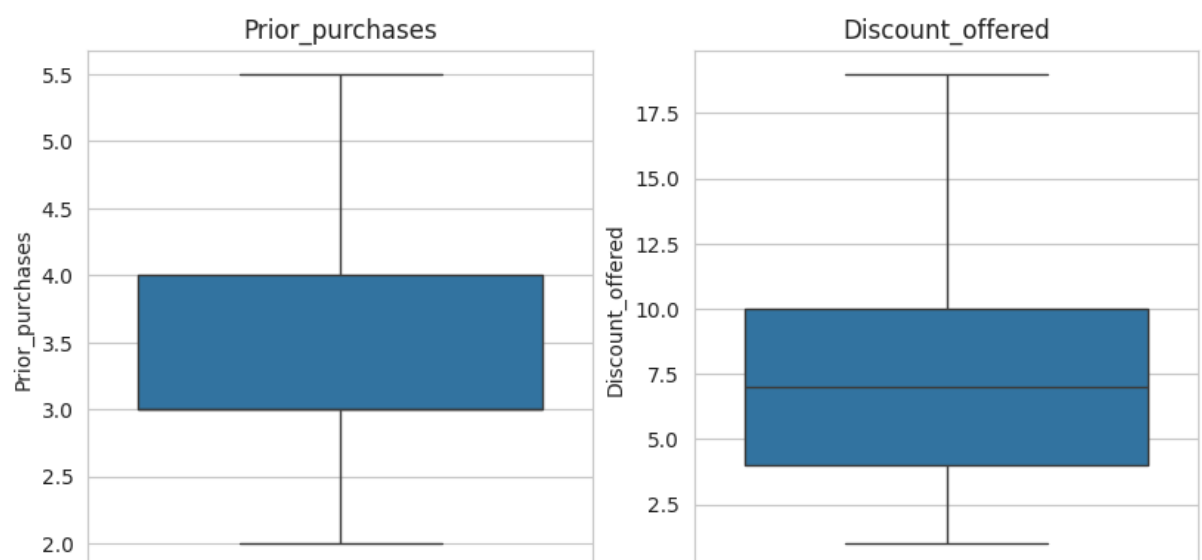- As we don't want to remove outliers Instead We'll perform capping.

Capping

```python
columns = ["Discount_offered", "Prior_purchases"]

for i in columns:
  q3 = df[i].quantile(0.75)
  q1 = df[i].quantile(0.25)
  IQR = q3 - q1

  upper_limit = q3 + 1.5 * IQR
  lower_limit = q1 - 1.5 * IQR

  df[i] = np.where(df[i] > upper_limit, upper_limit, np.where(df[i] < lower_limit, lower_limit, df[i]))
```

Now our data has no outliers

c. **Exploratory Data Analysis (EDA):** Explores the dataset through descriptive statistics, data visualization, and preliminary analysis. It aims to gain insights into the distribution of variables, identify patterns or trends, and detect any anomalies or outliers that may impact the analysis. EDA also helps in understanding the relationships between different variables and informs subsequent modeling decisions.
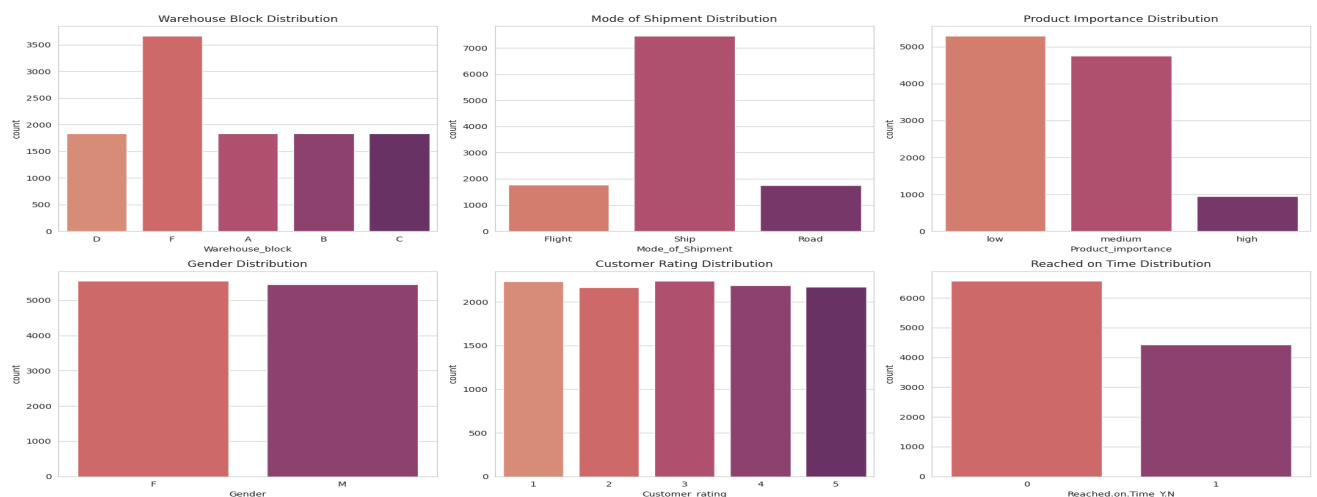
Statistics of Numeric Columns

```
df.describe()
```

|  | Customer_care_calls | Customer_rating | Cost_of_the_Product | Prior_purchases | Discount_offered | Weight_in_gms | Reached.on.Time_Y.N |
|---|---|---|---|---|---|---|---|
| count | 10999.000000 | 10999.000000 | 10999.000000 | 10999.000000 | 10999.000000 | 10999.000000 | 10999.000000 |
| mean | 4.054459 | 2.990545 | 210.196836 | 3.421629 | 8.590963 | 3634.016729 | 0.403309 |
| std | 1.141490 | 1.413603 | 48.063272 | 1.136903 | 6.095461 | 1635.377251 | 0.490584 |
| min | 2.000000 | 1.000000 | 96.000000 | 2.000000 | 1.000000 | 1001.000000 | 0.000000 |
| 25% | 3.000000 | 2.000000 | 169.000000 | 3.000000 | 4.000000 | 1839.500000 | 0.000000 |
| 50% | 4.000000 | 3.000000 | 214.000000 | 3.000000 | 7.000000 | 4149.000000 | 0.000000 |
| 75% | 5.000000 | 4.000000 | 251.000000 | 4.000000 | 10.000000 | 5050.000000 | 1.000000 |
| max | 7.000000 | 5.000000 | 310.000000 | 5.500000 | 19.000000 | 7846.000000 | 1.000000 |

## Summary

- The dataset contains 10,999 entries.
- Customer care calls range from 2 to 7, with an average of approximately 4 calls per customer.
- Customer ratings are evenly distributed between 1 (lowest) and 5 (highest).
- The cost of the product varies significantly, with a minimum of $96 and a maximum of $310.
- Prior purchases range from 2 to 5, indicating the number of times customers have purchased before.
- Discount offered varies, with a maximum discount of 19%.
- The weight of the products ranges from 1001 grams to 7846 grams.
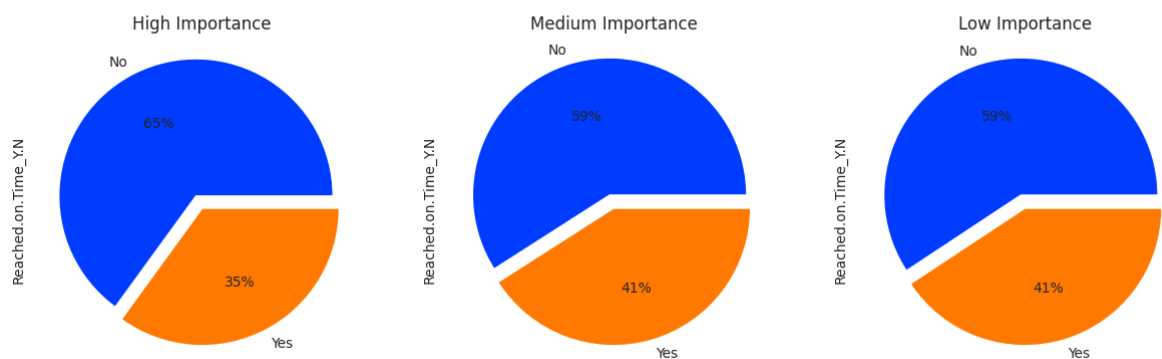- Approximately 59.7% of the shipments did not reach on time.

# Distribution of Categorical Columns

- **Warehouse Block Distribution:** The products are evenly distributed across the warehouse blocks A to D. Block F has highest frequency.
- **Mode of Shipment Distribution:** Ship is the most common mode of shipment compared to the Flight and Road.
- **Product Importance Distribution:** Most products are categorized as low importance, followed by medium and high.
- **Gender Distribution:** The dataset contains a nearly equal distribution of male and female customers.
- **Customer Rating Distribution:** Customer ratings are evenly distributed from 1 to 5.
- **Reached on Time Distribution:** There's a higher number of shipments that did not reach on time compared to those that did.

## Impact of Product Importance on Delivery



Impact of Product Importance on Delivery

This visualization shows that if the product importance is high there is a slight chance that it will not reach on time.

## Customer Rating and On-Time Delivery.

This chart illustrates that the distribution of customer ratings (from 1 to 5) is fairly consistent across both on-time and not on-time deliveries. This suggests that the customer rating may not be directly influenced by whether the product was delivered on time or not. Each rating level shows a similar pattern of distribution between products that were delivered on time and those that were not.

## Impact of Shipment Mode on Delivery



Impact of Shipment Mode on Delivery

All shipment modes show a similar ratio of Reach on time and Not Reached on time shipments.

# Distribution of Cost.

Customer Satisfaction:

| | Customer_care_calls | Not Reached on Time | Reached on Time |
|---|---|---|---|
| 0 | 2 | 0.652038 | 0.347962 |
| 1 | 3 | 0.625117 | 0.374883 |
| 2 | 4 | 0.597695 | 0.402305 |
| 3 | 5 | 0.584192 | 0.415808 |
| 4 | 6 | 0.516288 | 0.483712 |
| 5 | 7 | 0.516260 | 0.483740 |

This table shows that when Customer raises more queries to the customer care, There slight chance that Customer Support will address the issue and delivery will make it on time.

Scatter plot of Cost and Discount



As per the visualization It is clear that products with heavy discounts (10%++) are more likely to Not reach on time.

## Cost vs Customer Care Calls



We can see that the number of customer calls are increasing if the Cost of the Product is high.

## Correlation Matrix



Conclusions from Correlation matrix :-

- Discount Offered has a high negative correlation with Reached on Time 40%.
- Weights in gram have positive correlation with Reached on Time 27%.
- Discount Offered and weights in grams have negative correlation 38%.
- Customer care calls and weights in grams have a negative correlation 28%.
- Customer care calls and cost of the product have a positive correlation of 32%.
- Prior Purchases and Customer care calls have slightly positive correlation.

# 3. Feature Engineering

a. **Encoding Categorical Variables:** Explains the process of encoding categorical variables into numerical format suitable for machine learning algorithms. Technique used is one-hot encoding to be utilized to transform categorical variables into a format that the model can understand and interpret effectively. The choice of encoding method depends on the nature of the categorical variables and the requirements of the modeling task.

b. **Feature Construction:**

```python
# Overall Inteaction
df_["Over All Interaction"] = df_["Customer_care_calls"] * df_["Customer_rating"]

# Cost per Gram
df_["Cost Per Unit"] = df_["Cost_of_the_Product"] / df["Weight_in_gms"]
```

c. **Feature Scaling and Transformation Methods:** We employed robust scaling to handle outliers in key predictive features, using the median and interquartile range (IQR) for scaling instead of mean and standard deviation. Additionally, to address skewed distributions and enhance model robustness, we applied the Yeo-Johnson power transformation. This method adjusts data distributions to resemble a Gaussian distribution, contributing to improved model performance by mitigating the impact of outliers and non-normal data distributions.

## Scaling

```python
from sklearn.preprocessing import RobustScaler


rbs = RobustScaler()

X_train = rbs.fit_transform(X_train)
X_test = rbs.transform(X_test)
```

## Power Transformation

```python
from sklearn.preprocessing import PowerTransformer

pt = PowerTransformer()

X_train = pt.fit_transform(X_train)
X_test = pt.transform(X_test)
```

# Model Development

a. **Model Selection and Evaluation Criteria:** Discusses the criteria used for selecting the most appropriate model for the given prediction task. This may include considerations such as model performance metrics (accuracy, precision, recall, F1-score), computational efficiency, interpretability, and scalability to larger datasets.

b. **Training and Testing Data Split:** Explains the process of splitting the dataset into training and testing subsets. Typically, a portion of the data is reserved for training the model, while the remaining data is used to evaluate the model's performance and generalization capabilities.

```python
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 123, test_size = 0.2)
```

c. **Implementation of Machine Learning Models:**
   1. *Logistic Regression:* Describes the logistic regression algorithm and its application in binary classification tasks. It explains the logistic function, parameter estimation, and interpretation of coefficients.
      Training:

      ```python
      from sklearn.linear_model import LogisticRegression

      lr = LogisticRegression()

      lr.fit(X_train, y_train)
      ```

   2. *Random Forest:* Introduces the random forest algorithm, which is an ensemble learning method based on decision trees. It covers the concept of bagging, tree construction, and feature importance in random forest models.
      Training:

      ```python
      from sklearn.ensemble import RandomForestClassifier

      rf = RandomForestClassifier()

      rf.fit(X_train, y_train)
      ```

   3. *Support Vector Machines (SVM):* Discusses the SVM algorithm and its ability to perform linear and nonlinear classification. It explains the concept of hyperplanes, kernel methods, and optimization techniques used in SVM.
      Training:

      ```python
      from sklearn.svm import SVC

      svc = SVC()

      svc.fit(X_train, y_train)
      ```

4. *XGBoost:* Introduces the XGBoost algorithm, which is an efficient implementation of gradient boosting. It covers the principles of boosting, tree ensembles, and regularization techniques used in XGBoost models.
Training:

```python
import xgboost as xgb

xgb = xgb.XGBClassifier()

xgb.fit(X_train, y_train)
```

# Model Evaluation

a. **Performance Metrics for Classification Models:** Discusses the evaluation metrics used to assess the performance of classification models. Common metrics include accuracy, precision, recall, F1-score, ROC-AUC (Receiver Operating Characteristic - Area Under Curve), and confusion matrix. Each metric provides insights into different aspects of model performance, such as prediction accuracy, ability to correctly identify positive cases, and trade-offs between true positive and false positive rates.

```
Model : Logistic Regression

Accuracy : 0.6445454545454545
ROC-AUC  : 0.6218024047226023
F1 Score : 0.53341288878281623

Classification Report:
              precision    recall  f1-score   support

           0       0.69      0.73      0.71      1322
           1       0.56      0.51      0.53       878

    accuracy                           0.64      2200
   macro avg       0.63      0.62      0.62      2200
weighted avg       0.64      0.64      0.64      2200

Confusion Matrix:
[[971 351]
 [431 447]]
```

```
Model : SVM

Accuracy : 0.6568181818181819
ROC-AUC  : 0.6758130326453671
F1 Score : 0.6416706217370669

Classification Report:
              precision    recall  f1-score   support

           0       0.79      0.58      0.67      1322
           1       0.55      0.77      0.64       878

    accuracy                           0.66      2200
   macro avg       0.67      0.68      0.66      2200
weighted avg       0.70      0.66      0.66      2200

Confusion Matrix:
[[769 553]
 [202 676]]


Model : Random Forest

Accuracy : 0.6468181818181818
ROC-AUC  : 0.6458797845467799
F1 Score : 0.5916973200210194

Classification Report:
              precision    recall  f1-score   support

           0       0.73      0.65      0.69      1322
           1       0.55      0.64      0.59       878

    accuracy                           0.65      2200
   macro avg       0.64      0.65      0.64      2200
weighted avg       0.66      0.65      0.65      2200

Confusion Matrix:
[[860 462]
 [315 563]]
```

```
Model : XGBoost

Accuracy : 0.6695454545454546
ROC-AUC  : 0.6672769221756226
F1 Score : 0.613092070250133

Classification Report:
            precision    recall   f1-score   support

        0       0.75       0.68      0.71       1322
        1       0.58       0.66      0.61        878

  accuracy                           0.67       2200
 macro avg       0.66       0.67      0.66       2200
weighted avg     0.68       0.67      0.67       2200

Confusion Matrix:
[[897 425]
 [302 576]]
```

**Accuracy:** The accuracy metric measures the overall correctness of the model's predictions. Among the models listed, XGBoost has the highest accuracy of approximately 66.95%, followed closely by SVM with 65.68%.
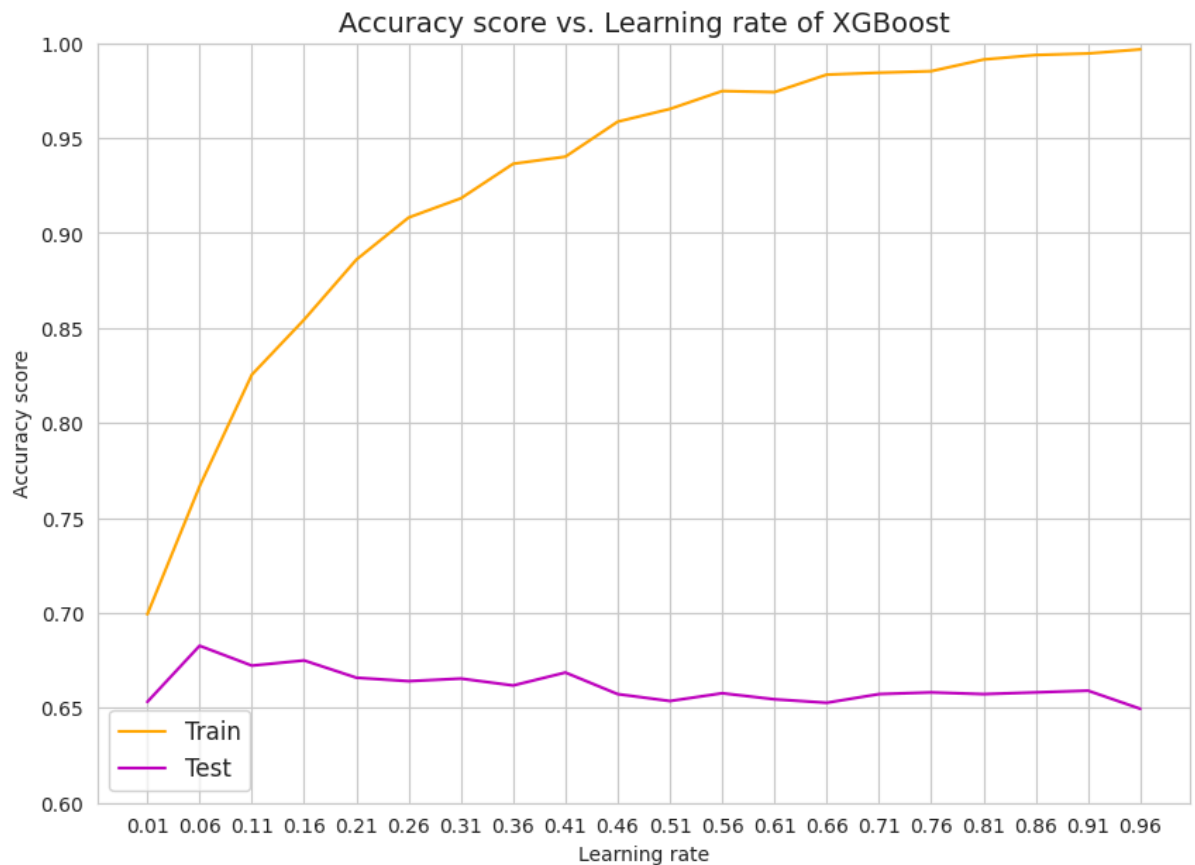
**ROC-AUC:** The ROC-AUC metric evaluates the model's ability to distinguish between positive and negative classes. A higher ROC-AUC score indicates better discrimination. XGBoost has the highest ROC-AUC score of approximately 0.67, followed by SVM with 0.676 and Random Forest with 0.645.

**F1 Score:** The F1 score balances precision and recall and is particularly useful when dealing with imbalanced datasets. XGBoost also has the highest F1 score of approximately 0.61, followed by SVM with 0.64.

Given these metrics, XGBoost appears to be the most promising model for further optimization through grid search. It demonstrates competitive performance across accuracy, ROC-AUC, and F1 score metrics, making it a strong candidate for fine-tuning hyperparameters.

Performance of XGB

## Accuracy score vs. Learning rate of XGBoost



**Hyper Parameter Tuning of XGB:**

```python
from sklearn.metrics import make_scorer

# ROC AUC Scorer
roc_auc_scorer = make_scorer(roc_auc_score, greater_is_better=True,
                             needs_threshold=True)
```

```python
# Define hyperparameters to tune
params = {
    'max_depth': [3, 4, 5],
    'learning_rate': [0.1, 0.01, 0.05],
    'n_estimators': [100, 200, 300],
    'min_child_weight': [1, 3, 5],
    'gamma': [0, 0.1, 0.2],
    'subsample': [0.8, 0.9, 1.0]
}
```

```python
from sklearn.model_selection import GridSearchCV

# Perform GridSearchCV for hyperparameter tuning
grid_search = GridSearchCV(estimator=xgb, param_grid=params,
                           scoring=roc_auc_scorer, cv=5,
                           n_jobs=-1)
grid_search.fit(X_train, y_train)
```

```
# Print best parameters and best score
print("Best parameters found: ", grid_search.best_params_)
print("Best ROC_AUC found: ", grid_search.best_score_)

# Evaluate the model on test set
best_model = grid_search.best_estimator_
test_accuracy = best_model.score(X_test, y_test)
print("Accuracy on test set: ", test_accuracy)
```

```
Best parameters found:  {'gamma': 0.2, 'learning_rate': 0.01, 'max_depth': 3, 'min_child_weight': 3, 'n_estimators': 300, 'subsample': 1.0}
Best ROC_AUC found:  0.7539247770661708
Accuracy on test set:  0.6868181818181818
```

In the grid search conducted on the XGBoost classifier, the best combination of hyperparameters was found to be {'gamma': 0.2, 'learning_rate': 0.01, 'max_depth': 3, 'min_child_weight': 3, 'n_estimators': 300, 'subsample': 1.0}. This optimized configuration yielded a notable improvement in the model's performance, with the best ROC-AUC score reaching 0.7539. Additionally, the accuracy on the test set improved to 0.6868. This indicates that the tuned XGBoost classifier achieved better discrimination between positive and negative classes, resulting in enhanced overall predictive performance compared to the default parameter settings. Overall, the grid search optimization process succeeded in fine-tuning the XGBoost classifier, leading to improved model accuracy and discriminative power.

b. **Cross-validation Techniques:** Explains the use of cross-validation methods to assess the generalization performance of the models. Technique used is k-fold cross-validation. Cross-validation helps estimate the model's performance on unseen data and reduces the risk of overfitting.

```
from sklearn.model_selection import cross_val_score

np.mean(cross_val_score(best_model, X_train, y_train, cv=5, scoring=roc_auc_scorer))
```

```
0.7539247770661708
```

## Deployment and Integration

a. **Model Deployment Strategies:** Various strategies for deploying the trained predictive model into production environments are possible. This may include deploying the model as a standalone application, integrating it into existing software systems, or using cloud-based platforms for scalability and accessibility.

b. **Integration of Predictive Model into E-commerce Systems:** The process of integrating the predictive model seamlessly into the company's e-commerce systems. This involves considerations such as data input and output formats, API (Application Programming Interface) design, and compatibility with existing software architecture.

c. **Considerations for Real-time Shipment Tracking:** The challenges and considerations for implementing real-time shipment tracking capabilities using the predictive model. This includes ensuring timely data updates, handling streaming data, and maintaining model performance in dynamic environments.

By considering these deployment and integration strategies, the company can effectively leverage the predictive model to optimize logistics operations, improve delivery efficiency, and enhance customer satisfaction in the competitive e-commerce industry.

# Conclusion

**Recap of Key Findings and Implications**
Based on the analysis and findings from the capstone project, several conclusions can be drawn:

1. **Customer Interaction and Satisfaction:** The number of customer care calls seems to increase with the cost of the product, indicating that customers may have higher expectations and demand better service for higher-priced items. This suggests that the company should focus on improving customer service channels and responsiveness to inquiries, especially for more expensive products.

2. **Impact of Discounts on Timeliness:** Products with heavy discounts are more likely to not reach on time. This could imply that logistics and fulfillment operations may need to be optimized to handle high-demand periods, especially during discount seasons or sales events. Offering discounts while maintaining delivery efficiency could enhance customer satisfaction and loyalty.

3. **Product Importance and Timeliness:** There's a slight indication that high-importance products may have a lower likelihood of reaching on time. This warrants further investigation into the logistics and handling processes for high-priority items to ensure timely delivery. Strengthening supply chain management practices for such products could help mitigate delays and improve overall service reliability.

4. **Correlation Insights:** The correlation matrix reveals valuable insights into the relationships between various features and the target variable (reached on time). For instance, the positive correlation between customer care calls and the cost of the product suggests that higher-priced items may require more support or attention from customer service teams.

5. **Machine Learning Model Performance:** The machine learning model achieved a respectable ROC AUC score of 75%, indicating its ability to effectively classify whether a product will reach on time or not. While the model performs reasonably

well, continuous monitoring and refinement are necessary to adapt to changing business dynamics and evolving customer needs.

---

**Practical Recommendations for E-commerce Companies**
Based on these conclusions, here are some suggestions for the company:

1. **Optimize Logistics and Fulfillment Processes:** Invest in improving warehouse management systems, shipment tracking technologies, and delivery networks to enhance efficiency and accuracy in product delivery.

2. **Enhance Customer Communication Channels:** Strengthen customer care infrastructure and empower support teams with tools and resources to address inquiries and resolve issues promptly. Implement proactive communication strategies to keep customers informed about shipment status and potential delays.

3. **Segmentation and Prioritization:** Segment products based on importance, demand, and delivery requirements to prioritize resources and attention accordingly. Allocate resources strategically to ensure timely delivery of high-priority items while maintaining service standards for all products.

4. **Data-Driven Decision Making:** Continue leveraging data analytics and machine learning techniques to gain actionable insights into customer behavior, market trends, and operational performance. Foster a culture of data-driven decision-making to drive continuous improvement and innovation across the organization.

By implementing these recommendations and leveraging the insights gained from the capstone project, the company can strengthen its competitive position, enhance customer satisfaction, and drive sustainable growth in the dynamic e-commerce landscape.