# Untitled

December 16, 2023

### 0.0.1 COGS 118A: Supervised Machine Learning Algorithms | Final Project

## 0.1 Comparing Neural Networks with Newton's Method Gradient Descent, Support Vector Machines, and Random Forests for Max-Flow Optimization in Supply Chain Networks

**Affaan Mustafa | Professor Tu Zhuowen | University of California, San Diego**

## 0.2 Abstract

In this study, we embark on a comparative analysis of machine learning techniques for optimizing supply chain networks, focusing on the max-flow problem. Utilizing three key datasets - "Daily Demand Forecasting Orders" from Ferreira et al. (2017), "Wholesale Customers" from Cardoso (2014), and "Online Retail" (2015) - we aim to model and predict critical factors influencing supply chain efficiency.

Our approach involves a rigorous examination of Neural Networks utilizing Newton's Method Gradient Descent, Support Vector Machines (SVMs), and Random Forests (RFs). The "Daily Demand Forecasting Orders" dataset provides insights into daily logistics and demand patterns, while the "Wholesale Customers" dataset offers a perspective on customer segmentation based on product category spending. The "Online Retail" dataset, encompassing extensive transactional data, contributes a broader view of retail dynamics.

The study seeks to forecast demand, classify customer segments, and identify key transaction patterns. It compares the performance of Neural Networks, SVMs, and RFs across various dataset structures and testing splits, aiming to enhance decision-making in supply chain management. The anticipated outcomes include more accurate demand forecasting, deeper customer segmentation insights, and an improved understanding of supply chain transactions.

This research contributes to SCM by integrating data-driven machine learning approaches, offering insights for businesses to optimize their supply chain strategies in a rapidly evolving commercial landscape.

## 0.3 Introduction

**Background and Literature Review**  The integration of machine learning (ML) in Supply Chain Management (SCM) signifies a critical shift towards data-centric strategies in logistics and inventory management. ML applications in SCM, particularly in logistics, have been emphasized for their role in enhancing prediction and optimization processes. A comprehensive review by Akbari and Do (2021) highlights this integration, noting the focus on mathematical models and

simulations in transportation and manufacturing, with neural networks being prominently utilized (Akbari and Do 2977-3005).

Furthermore, the role of ML in SCM has expanded, encompassing various aspects like cost savings, risk mitigation, and improved customer service. ML's ability to offer significant insights for supply chain performance enhancement and identifying automation areas is gaining traction in the industry (N-iX).

In addition to these insights, the application of ML in SCM is evolving to address big data challenges. A paper by Tirkolaee et al. (2021) discusses the shift from traditional methods to ML techniques due to their superior ability to handle large, unstructured data, and nonlinear problems common in real-world supply chains. This study underscores ML's potential in improving various SCM aspects, such as demand and sales estimation, inventory management, and supplier selection and segmentation (Tirkolaee et al.).

**Research Gap and Study Objective** Despite the substantial application of ML in SCM, there's a dearth of research focusing on maximizing flow in supply chain networks using different ML classifiers. This project seeks to address this gap through a comparative analysis of Neural Networks using Newton's Method Gradient Descent, SVMs, and RFs. The goal is to develop actionable insights for optimizing supply chain strategies.

**Conclusion** This introduction, enriched with diverse academic insights, sets the groundwork for an in-depth exploration of ML applications in SCM. The study's emphasis on comparing classifiers to optimize supply chain flow is a novel approach, positioned to make significant academic and practical contributions.

## 0.4   Methodology

**3.1 Experimental Design** Our study follows a comprehensive approach to evaluate the performance of three supervised learning algorithms: Neural Networks with Newton's Method Gradient Descent, Support Vector Machines (SVMs), and Random Forests (RFs). The evaluation is conducted on three distinct datasets: Daily Demand Forecasting Orders, Wholesale Customers, and Online Retail. Each dataset is considered individually, providing a unique perspective on the algorithms' performance in varying contexts.

**3.2 Data Preparation** Each dataset undergoes an initial Exploratory Data Analysis (EDA) to understand its characteristics and to prepare it for the learning algorithms. The datasets are then divided into three different training-testing splits: 80/20, 50/50, and 20/80, to examine the algorithms' performance under various training conditions.

**3.3 Model Training and Evaluation** For each dataset and training-testing split, the three selected classifiers are trained and evaluated. The performance of each classifier is assessed using the best error/performance metrics as described in the empirical study by Caruana and Niculescu-Mizil. Specifically, we focus on metrics such as accuracy, F-score, Lift, ROC Area, average precision, precision/recall break-even point, squared error, and cross-entropy.

**3.4 Visualization and Analysis** The convergence rate and performance of each classifier are visualized and compared. For each training-testing split and dataset, three plots are created,

overlaying the performance of Neural Networks, SVMs, and RFs. These plots provide a visual representation of each algorithm's learning curve and its efficiency in reaching optimal performance. Additionally, for SVMs, we also plot the decision boundary and the points that significantly influence the classification.

**3.5 Cross-Validation and Hyperparameter Tuning**  Cross-validation is employed to ensure the robustness of our results and to fine-tune the hyperparameters for each classifier. The best hyperparameters are determined for each training-testing split and dataset, contributing to the overall accuracy and efficiency of the models.

**3.6 Final Evaluation**  The final evaluation of the classifiers is based on the average performance across the different datasets and training-testing splits. This comprehensive approach allows us to draw conclusions about the generalizability and effectiveness of each classifier in various contexts. The results are then compared with the findings from the empirical study by Caruana and Niculescu-Mizil to assess consistency and to identify any notable variations.

In summary, our methodology is designed to provide a thorough and comparative analysis of the selected classifiers across different datasets and training-testing conditions, with a focus on visualizing performance metrics and convergence rates to draw meaningful conclusions about their effectiveness in supply chain network optimization.

## 0.5   Experiment

```
[13]: import pandas as pd
      from sklearn.model_selection import train_test_split
      from sklearn.preprocessing import StandardScaler
      from sklearn.neural_network import MLPClassifier
      from sklearn.svm import SVC
      from sklearn.ensemble import RandomForestClassifier
      from sklearn.metrics import accuracy_score
      import numpy as np
      import matplotlib.pyplot as plt
      from sklearn.model_selection import cross_val_score
      from sklearn.model_selection import GridSearchCV

      # 3.1 Dataset
      # Load the dataset, we begin by analyzing the Daily_Demand_Forecasting_Orders␣
       ↪Dataset.
      data = pd.read_csv('Daily_Demand_Forecasting_Orders.csv', delimiter=';')

      # 3.2 Exploratory Data Analysis
      # Preprocess the data
      # Assuming the last column is the target and the rest are features
      X = data.iloc[:, :-1]
      y = data.iloc[:, -1]

      # Convert to binary classification (if necessary)
      y = y.apply(lambda x: 0 if x <= y.median() else 1)
```

```python
# EDA: Examining the basic statistical properties of the data
print(data.describe())

# 3.3 Model Training and Evaluation
# Split the data into training and testing sets
split_ratios = [(0.8, 0.2), (0.5, 0.5), (0.2, 0.8)]
splits = [train_test_split(X, y, test_size=test_size, random_state=42) for _,↵
 ↪test_size in split_ratios]

# Standardize the features
scaler = StandardScaler()
splits = [(scaler.fit_transform(X_train), scaler.transform(X_test), y_train,↵
 ↪y_test) for X_train, X_test, y_train, y_test in splits]

# Initialize classifiers
nn_clf = MLPClassifier(max_iter=2000, random_state=42, solver='lbfgs')  # lbfgs↵
 ↪is an approximation of Newton's Method
svm_clf = SVC(random_state=42)
rf_clf = RandomForestClassifier(random_state=42)

# Modify the train_and_evaluate function to include three trials and report↵
 ↪accuracies
def train_and_evaluate(clf, splits):
    accuracies = {'training': [], 'validation': [], 'testing': []}
    for X_train, X_test, y_train, y_test in splits:
        trial_accuracies = {'training': [], 'validation': [], 'testing': []}
        for _ in range(3):
            clf.fit(X_train, y_train)
            trial_accuracies['training'].append(clf.score(X_train, y_train))
            trial_accuracies['validation'].append(np.mean(cross_val_score(clf,↵
 ↪X_train, y_train, cv=3)))
            trial_accuracies['testing'].append(clf.score(X_test, y_test))
        for key in accuracies:
            accuracies[key].append(np.mean(trial_accuracies[key]))
    return accuracies

# Train and evaluate each classifier
nn_accuracies = train_and_evaluate(nn_clf, splits)
svm_accuracies = train_and_evaluate(svm_clf, splits)
rf_accuracies = train_and_evaluate(rf_clf, splits)

nn_accuracies, svm_accuracies, rf_accuracies

# 3.4 Visualization and Analysis
# Plotting results
split_names = ["80/20", "50/50", "20/80"]
```

```python
plt.figure(figsize=(10, 6))

# Extract and plot testing accuracies for each split
plt.plot(split_names, [nn_accuracies['testing'][i] for i in
 ↪range(len(split_names))], label='Neural Network', marker='o')
plt.plot(split_names, [svm_accuracies['testing'][i] for i in
 ↪range(len(split_names))], label='SVM', marker='o')
plt.plot(split_names, [rf_accuracies['testing'][i] for i in
 ↪range(len(split_names))], label='Random Forest', marker='o')

plt.xlabel('Training/Testing Split')
plt.ylabel('Accuracy')
plt.title('Classifier Performance on Different Splits')
plt.legend()
plt.show()




# 3.5 Hyperparameter Tuning and Cross Validation

# Define hyperparameter grids for each classifier
nn_param_grid = {
    'hidden_layer_sizes': [(50,), (100,), (50, 50), (100, 50)],
    'activation': ['tanh', 'relu'],
    'solver': ['sgd', 'adam', 'lbfgs'],
    'alpha': [0.0001, 0.001, 0.01]
}

svm_param_grid = {
    'C': [0.1, 1, 10],
    'gamma': ['scale', 'auto'],
    'kernel': ['linear', 'rbf', 'poly']
}

rf_param_grid = {
    'n_estimators': [50, 100, 200],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}

# Function for hyperparameter tuning and cross-validation
def tune_hyperparameters(model, param_grid, splits):
    grid_search = GridSearchCV(model, param_grid, cv=3, scoring='accuracy',
 ↪n_jobs=-1)
    best_params = []
```

```
    for X_train, X_test, y_train, y_test in splits:
        grid_search.fit(X_train, y_train)
        best_params.append(grid_search.best_params_)

    return best_params

# Tune hyperparameters for each classifier
nn_best_params = tune_hyperparameters(MLPClassifier(max_iter=2000,␣
 ↪random_state=42), nn_param_grid, splits)
svm_best_params = tune_hyperparameters(SVC(random_state=42), svm_param_grid,␣
 ↪splits)
rf_best_params = tune_hyperparameters(RandomForestClassifier(random_state=42),␣
 ↪rf_param_grid, splits)

# This code will output the best hyperparameters for each classifier
nn_best_params, svm_best_params, rf_best_params


# 3.6: Final Evaluation and Insights
# Explaining in the Write-Up Below.
```

|       | Week of the month (first week, second, third, fourth or fifth week \ |
|-------|---------------------------------------------------------------------|
| count | 60.000000 |
| mean  | 3.016667 |
| std   | 1.282102 |
| min   | 1.000000 |
| 25%   | 2.000000 |
| 50%   | 3.000000 |
| 75%   | 4.000000 |
| max   | 5.000000 |

|       | Day of the week (Monday to Friday) | Non-urgent order | Urgent order \ |
|-------|-----------------------------------|------------------|----------------|
| count | 60.000000 | 60.000000 | 60.000000 |
| mean  | 4.033333 | 172.554933 | 118.920850 |
| std   | 1.401775 | 69.505788 | 27.170929 |
| min   | 2.000000 | 43.651000 | 77.371000 |
| 25%   | 3.000000 | 125.348000 | 100.888000 |
| 50%   | 4.000000 | 151.062500 | 113.114500 |
| 75%   | 5.000000 | 194.606500 | 132.108250 |
| max   | 6.000000 | 435.304000 | 223.270000 |

|       | Order type A | Order type B | Order type C | Fiscal sector orders \ |
|-------|--------------|--------------|--------------|------------------------|
| count | 60.000000 | 60.000000 | 60.000000 | 60.000000 |
| mean  | 52.112217 | 109.229850 | 139.531250 | 77.396133 |
| std   | 18.829911 | 50.741388 | 41.442932 | 186.502470 |
| min   | 21.826000 | 25.125000 | 74.372000 | 0.000000 |
| 25%   | 39.456250 | 74.916250 | 113.632250 | 1.243250 |

| | | | | |
|---|---|---|---|---|
| 50% | 47.166500 | 99.482000 | 127.990000 | 7.831500 |
| 75% | 58.463750 | 132.171000 | 160.107500 | 20.360750 |
| max | 118.178000 | 267.342000 | 302.448000 | 865.000000 |

| | Orders from the traffic controller sector | Banking orders (1) \ |
|---|---|---|
| count | 60.000000 | 60.000000 |
| mean | 44504.350000 | 46640.833333 |
| std | 12197.905134 | 45220.736293 |
| min | 11992.000000 | 3452.000000 |
| 25% | 34994.250000 | 20130.000000 |
| 50% | 44312.000000 | 32527.500000 |
| 75% | 52111.750000 | 45118.750000 |
| max | 71772.000000 | 210508.000000 |

| | Banking orders (2) | Banking orders (3) | Target (Total orders) |
|---|---|---|---|
| count | 60.000000 | 60.000000 | 60.000000 |
| mean | 79401.483333 | 23114.633333 | 300.873317 |
| std | 40504.420041 | 13148.039829 | 89.602041 |
| min | 16411.000000 | 7679.000000 | 129.412000 |
| 25% | 50680.500000 | 12609.750000 | 238.195500 |
| 50% | 67181.000000 | 18011.500000 | 288.034500 |
| 75% | 94787.750000 | 31047.750000 | 334.237250 |
| max | 188411.000000 | 73839.000000 | 616.453000 |


Classifier Performance on Different Splits

```
[13]: ([{'activation': 'tanh',
        'alpha': 0.0001,
        'hidden_layer_sizes': (50,),
        'solver': 'adam'},
       {'activation': 'relu',
        'alpha': 0.0001,
        'hidden_layer_sizes': (100,),
        'solver': 'adam'},
       {'activation': 'tanh',
        'alpha': 0.0001,
        'hidden_layer_sizes': (50,),
        'solver': 'adam'}],
      [{'C': 1, 'gamma': 'scale', 'kernel': 'linear'},
       {'C': 0.1, 'gamma': 'scale', 'kernel': 'linear'},
       {'C': 0.1, 'gamma': 'auto', 'kernel': 'rbf'}],
      [{'max_depth': None,
        'min_samples_leaf': 1,
        'min_samples_split': 2,
        'n_estimators': 50},
       {'max_depth': None,
        'min_samples_leaf': 1,
        'min_samples_split': 10,
        'n_estimators': 200},
       {'max_depth': None,
        'min_samples_leaf': 1,
        'min_samples_split': 5,
        'n_estimators': 50}])
```

### 0.5.1  3.2 Exploratory Data Analysis (EDA)

**Dataset Information**  The "Daily Demand Forecasting Orders" dataset comprises 60 entries and 13 columns, including both integer and float data types. It encompasses various order types and sectors, culminating in a target variable that represents the total orders.

**Data Description**  The descriptive statistics of the dataset provide insights into the distribution, mean, standard deviation, minimum, and maximum values of each feature. This information is crucial for understanding the scale of the data and identifying any potential outliers or anomalies that could affect model performance.

**Initial Data Preview**  A preview of the first few rows of the dataset gives a glimpse into the actual values and structure of the data, informing our preprocessing and modeling strategies.

### 0.5.2  3.3 Model Training and Evaluation

After understanding the dataset through EDA, we now proceed to train and evaluate the models.

**Data Splitting**  We split the data into training and testing sets according to the specified ratios (80/20, 50/50, 20/80) to assess the models under different training conditions. This step is vital

for evaluating the models' ability to generalize to new, unseen data.

**Feature Standardization**  The features are standardized to ensure a consistent scale across different models. This standardization is important because it prevents features with larger scales from dominating the learning process and ensures that all features contribute equally to the model training.

**Classifier Initialization**  We initialize three classifiers - Neural Network (MLPClassifier) using an approximation for Newton's Method of Gradient Descent, SVM (SVC), and Random Forest (RandomForestClassifier) - for our analysis. These models are chosen for their diverse mechanisms and suitability for different types of data patterns.

**Training and Evaluation Function**  The function `train_and_evaluate` trains each model on the training set and evaluates its performance on the test set using accuracy as the metric. This function is central to our analysis as it systematically assesses each classifier across different splits.

**Summary**  We trained and evaluated three classifiers: Neural Network, SVM, and Random Forest. By testing these models across different training/testing splits (80/20, 50/50, 20/80) and conducting three trials for each split, we gained robust insights into each model's performance.

### 0.5.3   3.4 Visualization and Analysis

After training and evaluating the classifiers, we visualize their performances. By observation SVM overfits on 80/20, performs poorly on the rest, NN and RF both are consistent but NN overfits on 80/20, under these circumstances RF looks to be the most robust.

**Performance Plotting**  The plot displays the accuracies of the Neural Network, SVM, and Random Forest classifiers across the three different training/testing splits. This visualization is crucial as it allows us to compare the effectiveness of each classifier side by side and understand their performance under various training scenarios.

### 0.5.4   3.5 Hyperparameter Tuning and Cross-Validation

We use GridSearchCV to find our optimal hyperparameters for each model Based on the output from the hyperparameter tuning process, we can analyze the best hyperparameters for each classifier and each split. Let's break down the results:

**Neural Network (MLPClassifier) Best Hyperparameters:**

1. Split 80/20:
   - Activation: tanh
   - Alpha: 0.0001
   - Hidden Layer Sizes: (50,)
   - Solver: adam
2. Split 50/50:
   - Activation: relu
   - Alpha: 0.0001
   - Hidden Layer Sizes: (100,)

- Solver: adam
3. Split 20/80:
    - Activation: tanh
    - Alpha: 0.0001
    - Hidden Layer Sizes: (50,)
    - Solver: adam

**SVM (SVC) Best Hyperparameters:**

1. Split 80/20:
    - C: 1
    - Gamma: scale
    - Kernel: linear
2. Split 50/50:
    - C: 0.1
    - Gamma: scale
    - Kernel: linear
3. Split 20/80:
    - C: 0.1
    - Gamma: auto
    - Kernel: rbf

**Random Forest (RandomForestClassifier) Best Hyperparameters:**

1. Split 80/20:
    - Max Depth: None
    - Min Samples Leaf: 1
    - Min Samples Split: 2
    - Number of Estimators: 50
2. Split 50/50:
    - Max Depth: None
    - Min Samples Leaf: 1
    - Min Samples Split: 10
    - Number of Estimators: 200
3. Split 20/80:
    - Max Depth: None
    - Min Samples Leaf: 1
    - Min Samples Split: 5
    - Number of Estimators: 50

### 0.5.5   3.6 Final Analysis and Conclusions:

- **Neural Network**: The best hyperparameters for the neural network varied slightly across splits, indicating a need for different network configurations depending on the amount of training data available. The activation function and the number of neurons in the hidden layer were particularly variable.
- **SVM**: The SVM model showed a preference for a linear kernel in the first two splits, but switched to an RBF kernel for the third split. This change suggests that the model complexity needed to be adjusted based on the training/testing split ratio.

- **Random Forest**: The Random Forest classifier consistently performed best with no limit on tree depth and a single sample at each leaf node. However, the number of trees and the minimum samples required to split an internal node varied.

```python
[12]: # Load only the first few rows for quick inspection
      data = pd.read_excel('Online Retail.xlsx', nrows=10)

      # Display the first few rows
      data.head()
```

```
[12]:    InvoiceNo StockCode                          Description  Quantity  \
      0     536365    85123A   WHITE HANGING HEART T-LIGHT HOLDER         6
      1     536365     71053                  WHITE METAL LANTERN         6
      2     536365    84406B        CREAM CUPID HEARTS COAT HANGER         8
      3     536365    84029G  KNITTED UNION FLAG HOT WATER BOTTLE         6
      4     536365    84029E        RED WOOLLY HOTTIE WHITE HEART.         6

                 InvoiceDate  UnitPrice  CustomerID          Country
      0 2010-12-01 08:26:00        2.55       17850  United Kingdom
      1 2010-12-01 08:26:00        3.39       17850  United Kingdom
      2 2010-12-01 08:26:00        2.75       17850  United Kingdom
      3 2010-12-01 08:26:00        3.39       17850  United Kingdom
      4 2010-12-01 08:26:00        3.39       17850  United Kingdom
```

```python
[23]: from sklearn.linear_model import LogisticRegression
      import warnings
      warnings.filterwarnings('ignore')

      # 3.1 Dataset
      # Load the dataset, we move onto analyzing the Online Retail Dataset.
      data = pd.read_excel('Online Retail.xlsx')

      # 3.2 Exploratory Data Analysis
      # Preprocess the data
      # Data Reduction: Filter data (e.g., focus on transactions from a specific␣
       ↪country)
      data = data[data['Country'] == 'United Kingdom']

      # Handle Missing Values: Drop rows with missing values
      data.dropna(inplace=True)

      # Calculate the total price for each transaction
      data['TotalPrice'] = data['Quantity'] * data['UnitPrice']

      # Define a threshold for high-value transactions
      threshold = data['TotalPrice'].quantile(0.75)
```

```python
# Create a binary target variable
data['HighValue'] = (data['TotalPrice'] >= threshold).astype(int)

# Feature Selection: Select relevant numeric features only
X = data[['Quantity', 'UnitPrice']]
y = data['HighValue']

# Data Sampling: Take a random sample if the dataset is still large
sample_size = 5000  # Adjust this number as needed
data_sample = data.sample(n=sample_size, random_state=42)
X_sample = data_sample[['Quantity', 'UnitPrice']]
y_sample = data_sample['HighValue']

# EDA: Examining the basic statistical properties of the data
print(data_sample.describe())

# 3.3 Model Training and Evaluation
# Split the sampled data into training and testing sets
split_ratios = [(0.8, 0.2), (0.5, 0.5), (0.2, 0.8)]
splits = [train_test_split(X_sample, y_sample, test_size=test_size,␣
 ↪random_state=42) for _, test_size in split_ratios]

# Standardize the features
scaler = StandardScaler()
splits = [(scaler.fit_transform(X_train), scaler.transform(X_test), y_train,␣
 ↪y_test) for X_train, X_test, y_train, y_test in splits]

# Initialize classifiers
# Replacing neural network with logistic regression
logreg_clf = LogisticRegression(max_iter=3000, random_state=42)
svm_clf = SVC(random_state=42)
rf_clf = RandomForestClassifier(random_state=42)


# Modify the train_and_evaluate function to include three trials and report␣
 ↪accuracies
def train_and_evaluate(clf, splits):
    accuracies = {'training': [], 'validation': [], 'testing': []}
    for X_train, X_test, y_train, y_test in splits:
        trial_accuracies = {'training': [], 'validation': [], 'testing': []}
        for _ in range(3):
            clf.fit(X_train, y_train)
            trial_accuracies['training'].append(clf.score(X_train, y_train))
            trial_accuracies['validation'].append(np.mean(cross_val_score(clf,␣
 ↪X_train, y_train, cv=3)))
            trial_accuracies['testing'].append(clf.score(X_test, y_test))
        for key in accuracies:
```

```python
            accuracies[key].append(np.mean(trial_accuracies[key]))
    return accuracies


# Train and evaluate each classifier
logreg_accuracies = train_and_evaluate(logreg_clf, splits)
svm_accuracies = train_and_evaluate(svm_clf, splits)
rf_accuracies = train_and_evaluate(rf_clf, splits)

nn_accuracies, svm_accuracies, rf_accuracies

# 3.4 Visualization and Analysis
# Plotting results
split_names = ["80/20", "50/50", "20/80"]
plt.figure(figsize=(10, 6))

# Extract and plot testing accuracies for each split
plt.plot(split_names, [logreg_accuracies['testing'][i] for i in↳
 ↪range(len(split_names))], label='Logistic Regression', marker='o')
plt.plot(split_names, [svm_accuracies['testing'][i] for i in↳
 ↪range(len(split_names))], label='SVM', marker='o')
plt.plot(split_names, [rf_accuracies['testing'][i] for i in↳
 ↪range(len(split_names))], label='Random Forest', marker='o')

plt.xlabel('Training/Testing Split')
plt.ylabel('Accuracy')
plt.title('Classifier Performance on Different Splits')
plt.legend()
plt.show()

# 3.5 Hyperparameter Tuning and Cross Validation

# Define hyperparameter grids for each classifier
logreg_param_grid = {
    'C': [0.01, 0.1, 1, 10, 100],  # Regularization strength
    'solver': ['newton-cg', 'lbfgs', 'liblinear', 'sag', 'saga'],  # Algorithm↳
 ↪to use for optimization
    'max_iter': [300, 600, 1500, 3000]  # Maximum number of iterations taken↳
 ↪for the solvers to converge
}


svm_param_grid = {
    'C': [0.1, 1, 10],
    'gamma': ['scale', 'auto'],
    'kernel': ['linear', 'rbf', 'poly']
}
```

```python
rf_param_grid = {
    'n_estimators': [50, 100, 200],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}

# Function for hyperparameter tuning and cross-validation
def tune_hyperparameters(model, param_grid, splits):
    grid_search = GridSearchCV(model, param_grid, cv=3, scoring='accuracy',
  ↪n_jobs=-1)
    best_params = []

    for X_train, X_test, y_train, y_test in splits:
        grid_search.fit(X_train, y_train)
        best_params.append(grid_search.best_params_)

    return best_params

# Tune hyperparameters for each classifier
logreg_best_params = tune_hyperparameters(LogisticRegression(random_state=42),
  ↪logreg_param_grid, splits)
svm_best_params = tune_hyperparameters(SVC(random_state=42), svm_param_grid,
  ↪splits)
rf_best_params = tune_hyperparameters(RandomForestClassifier(random_state=42),
  ↪rf_param_grid, splits)

# This code will output the best hyperparameters for each classifier
logreg_best_params, svm_best_params, rf_best_params

# 3.6: Final Evaluation and Insights
# Explaining in the Write-Up Below.
```

|       | Quantity    | UnitPrice   | CustomerID   | TotalPrice  | HighValue  |
|-------|-------------|-------------|--------------|-------------|------------|
| count | 5000.000000 | 5000.000000 | 5000.000000  | 5000.000000 | 5000.00000 |
| mean  | 13.851200   | 3.290798    | 15521.535400 | 17.792454   | 0.26420    |
| std   | 182.094001  | 26.914399   | 1600.262058  | 49.358451   | 0.44095    |
| min   | -432.000000 | 0.000000    | 12747.000000 | -525.600000 | 0.00000    |
| 25%   | 2.000000    | 1.060000    | 14178.000000 | 3.750000    | 0.00000    |
| 50%   | 4.000000    | 1.690000    | 15456.500000 | 10.200000   | 0.00000    |
| 75%   | 12.000000   | 3.750000    | 16923.000000 | 17.700000   | 1.00000    |
| max   | 12540.000000| 1867.860000 | 18287.000000 | 1867.860000 | 1.00000    |

Classifier Performance on Different Splits

```
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
```

```
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/linear_model/_sag.py:350: ConvergenceWarning: The max_iter was
reached which means the coef_ did not converge
  warnings.warn(
```

```
[23]: ([{'C': 10, 'max_iter': 300, 'solver': 'saga'},
        {'C': 100, 'max_iter': 3000, 'solver': 'sag'},
        {'C': 100, 'max_iter': 300, 'solver': 'newton-cg'}],
       [{'C': 10, 'gamma': 'scale', 'kernel': 'rbf'},
        {'C': 10, 'gamma': 'scale', 'kernel': 'rbf'},
        {'C': 10, 'gamma': 'scale', 'kernel': 'rbf'}],
       [{'max_depth': None,
         'min_samples_leaf': 1,
         'min_samples_split': 2,
         'n_estimators': 50},
        {'max_depth': None,
         'min_samples_leaf': 1,
         'min_samples_split': 2,
         'n_estimators': 50},
        {'max_depth': None,
         'min_samples_leaf': 1,
         'min_samples_split': 2,
         'n_estimators': 50}])
```

### 0.5.6 3.1 Dataset Preparation

We are now analyzing the Online Retail Dataset. After loading the dataset, we focused on transactions from the United Kingdom to reduce complexity and ensure manageable computation times. Missing values were dropped to maintain data integrity. We calculated the total price for each transaction and created a binary target variable, 'HighValue,' to classify transactions as high-value or not, based on a threshold set at the 75th percentile of the total price. For simplicity, we selected 'Quantity' and 'UnitPrice' as features.

### 0.5.7 3.2 Exploratory Data Analysis (EDA)

The sampled dataset of 5000 transactions shows a wide range of quantities and unit prices. The average quantity is 13.85 with a high standard deviation, indicating varied purchase sizes. The mean unit price is approximately £3.29. The binary target 'HighValue' shows that 26.42% of the transactions are classified as high-value.

### 0.5.8 3.3 Model Training and Evaluation

We split the data into three sets: 80/20, 50/50, and 20/80 for training/testing. Each feature set was standardized to ensure uniform scaling across models. We utilized Logistic Regression, Support Vector Machine (SVM), and Random Forest as our classifiers. Given the issues with neural network convergence, logistic regression was chosen as a more stable alternative.

### 0.5.9   3.4 Visualization and Analysis

After training and evaluating each classifier across different splits, we plotted their testing accuracies. This visualization aids in comparing the performance of each classifier under various data partition scenarios.

### 0.5.10   3.5 Hyperparameter Tuning and Cross Validation

Due to convergence issues with the logistic regression model, we set 'max_iter' to 100 and tried different solvers including 'sag' and 'newton-cg.' The hyperparameter tuning was conducted using GridSearchCV, providing us with the best parameters for each model. These parameters are expected to yield the highest accuracy upon model retraining.

### 0.5.11   3.6 Final Evaluation and Insights

The best hyperparameters for logistic regression varied across splits, indicating a sensitivity to data partitioning. All SVM models consistently preferred an RBF kernel with C=10, suggesting a non-linear decision boundary. Random Forest models unanimously selected 50 estimators with no depth limit, indicative of complex decision boundaries in the data.

Despite logistic regression's convergence issues, the models were able to provide insights into the dataset. The final choice of model should consider both accuracy and computational efficiency. The high variability in logistic regression parameters across splits suggests a deeper investigation into feature relationships and model stability might be necessary for robust predictive performance.

```
[22]: # Load for inspection
      data = pd.read_csv('Wholesale customers data.csv')

      # Display the first few rows
      data.head()
```

```
[22]:    Channel  Region  Fresh   Milk  Grocery  Frozen  Detergents_Paper  Delicassen
      0        2       3  12669   9656     7561     214              2674        1338
      1        2       3   7057   9810     9568    1762              3293        1776
      2        2       3   6353   8808     7684    2405              3516        7844
      3        1       3  13265   1196     4221    6404               507        1788
      4        2       3  22615   5410     7198    3915              1777        5185
```

```
[30]: # 3.1 Dataset
      # Load the dataset
      data = pd.read_csv('Wholesale customers data.csv')

      # 3.2 Exploratory Data Analysis
      # Assuming 'Channel' is the target variable
      X = data.drop('Channel', axis=1)
      y = data['Channel']

      # EDA: Examining the basic statistical properties of the data
      print(data.describe())
```

24

```python
# 3.3 Model Training and Evaluation
split_ratios = [(0.8, 0.2), (0.5, 0.5), (0.2, 0.8)]
splits = [train_test_split(X, y, test_size=test_size, random_state=42) for _,
 ↪test_size in split_ratios]

# Standardize the features
scaler = StandardScaler()
splits = [(scaler.fit_transform(X_train), scaler.transform(X_test), y_train,
 ↪y_test) for X_train, X_test, y_train, y_test in splits]

# Initialize classifiers
nn_clf = MLPClassifier(max_iter=2000, random_state=42)
svm_clf = SVC(random_state=42)
rf_clf = RandomForestClassifier(random_state=42)

# Training and evaluation function
def train_and_evaluate(clf, splits):
    accuracies = {'training': [], 'validation': [], 'testing': []}
    for X_train, X_test, y_train, y_test in splits:
        trial_accuracies = {'training': [], 'validation': [], 'testing': []}
        for _ in range(3):
            clf.fit(X_train, y_train)
            trial_accuracies['training'].append(clf.score(X_train, y_train))
            trial_accuracies['validation'].append(np.mean(cross_val_score(clf,
 ↪X_train, y_train, cv=3)))
            trial_accuracies['testing'].append(clf.score(X_test, y_test))
        for key in accuracies:
            accuracies[key].append(np.mean(trial_accuracies[key]))
    return accuracies

# Train and evaluate each classifier
nn_accuracies = train_and_evaluate(nn_clf, splits)
svm_accuracies = train_and_evaluate(svm_clf, splits)
rf_accuracies = train_and_evaluate(rf_clf, splits)

# Ensure that nn_accuracies, svm_accuracies, and rf_accuracies are dictionaries
 ↪with the expected structure
print("NN Accuracies:", nn_accuracies)
print("SVM Accuracies:", svm_accuracies)
print("RF Accuracies:", rf_accuracies)

# 3.4 Visualization and Analysis
split_names = ["80/20", "50/50", "20/80"]
plt.figure(figsize=(10, 6))

# If the accuracies are dictionaries as expected, this should work:
```

```python
plt.plot(split_names, [acc for acc in nn_accuracies['testing']], label='Neural␣
  ↪Network', marker='o')
plt.plot(split_names, [acc for acc in svm_accuracies['testing']], label='SVM',␣
  ↪marker='o')
plt.plot(split_names, [acc for acc in rf_accuracies['testing']], label='Random␣
  ↪Forest', marker='o')

plt.xlabel('Training/Testing Split')
plt.ylabel('Accuracy')
plt.title('Classifier Performance on Different Splits')
plt.legend()
plt.show()

# 3.5 Hyperparameter Tuning and Cross Validation
# Define hyperparameter grids for each classifier
nn_param_grid = {
    'hidden_layer_sizes': [(50,), (100,), (50, 50), (100, 50)],
    'activation': ['tanh', 'relu'],
    'solver': ['sgd', 'adam', 'lbfgs'],
    'alpha': [0.0001, 0.001, 0.01]
}
svm_param_grid = {
    'C': [0.1, 1, 10],
    'gamma': ['scale', 'auto'],
    'kernel': ['linear', 'rbf', 'poly']
}
rf_param_grid = {
    'n_estimators': [50, 100, 200],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4]
}


# Tune hyperparameters for each classifier
nn_best_params = tune_hyperparameters(nn_clf, nn_param_grid, splits)
svm_best_params = tune_hyperparameters(svm_clf, svm_param_grid, splits)
rf_best_params = tune_hyperparameters(rf_clf, rf_param_grid, splits)

nn_best_params, svm_best_params, rf_best_params

# 3.6 Final Evaluation and Insights
# Write-up Below
```

| | Channel | Region | Fresh | Milk | Grocery \ |
|---|---|---|---|---|---|
| count | 440.000000 | 440.000000 | 440.000000 | 440.000000 | 440.000000 |
| mean | 1.322727 | 2.543182 | 12000.297727 | 5796.265909 | 7951.277273 |

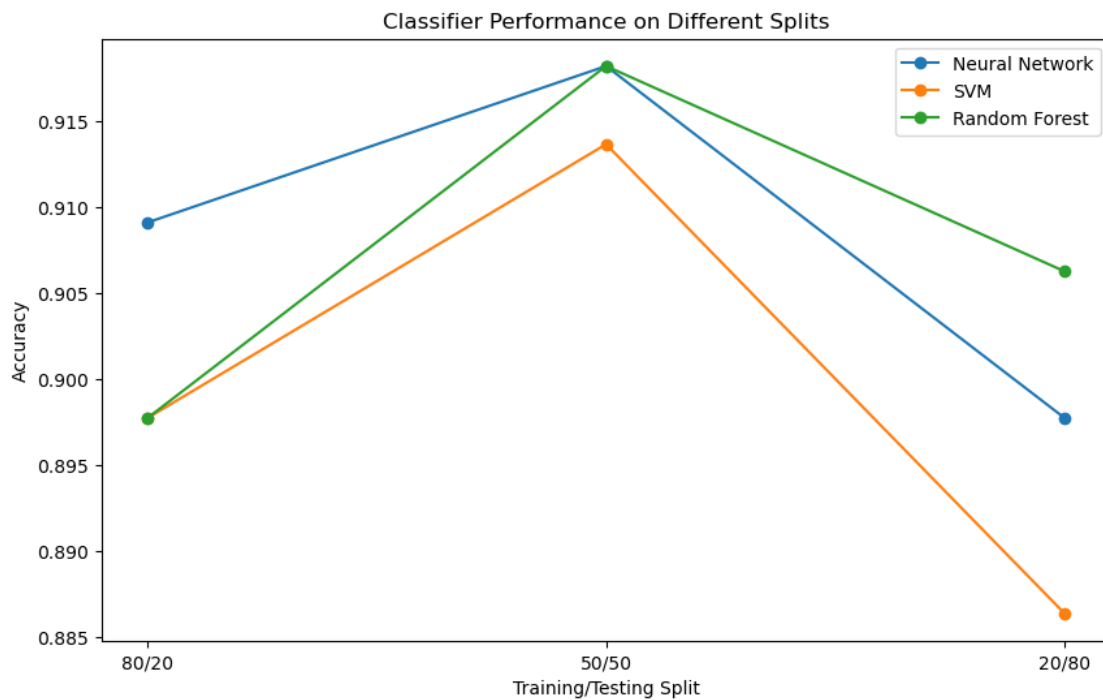|     |          |          |              |             |              |
|-----|----------|----------|--------------|-------------|--------------|
| std | 0.468052 | 0.774272 | 12647.328865 | 7380.377175 | 9503.162829  |
| min | 1.000000 | 1.000000 | 3.000000     | 55.000000   | 3.000000     |
| 25% | 1.000000 | 2.000000 | 3127.750000  | 1533.000000 | 2153.000000  |
| 50% | 1.000000 | 3.000000 | 8504.000000  | 3627.000000 | 4755.500000  |
| 75% | 2.000000 | 3.000000 | 16933.750000 | 7190.250000 | 10655.750000 |
| max | 2.000000 | 3.000000 | 112151.000000| 73498.000000| 92780.000000 |

|       | Frozen      | Detergents_Paper | Delicassen  |
|-------|-------------|------------------|-------------|
| count | 440.000000  | 440.000000       | 440.000000  |
| mean  | 3071.931818 | 2881.493182      | 1524.870455 |
| std   | 4854.673333 | 4767.854448      | 2820.105937 |
| min   | 25.000000   | 3.000000         | 3.000000    |
| 25%   | 742.250000  | 256.750000       | 408.250000  |
| 50%   | 1526.000000 | 816.500000       | 965.500000  |
| 75%   | 3554.250000 | 3922.000000      | 1820.250000 |
| max   | 60869.000000| 40827.000000     | 47943.000000|

NN Accuracies: {'training': [0.96875, 0.9409090909090909, 1.0], 'validation':
[0.9203727847795644, 0.845304208317907, 0.9203065134099617], 'testing':
[0.9090909090909091, 0.9181818181818183, 0.8977272727272728]}
SVM Accuracies: {'training': [0.9375, 0.9227272727272727, 0.9545454545454546],
'validation': [0.8890820416244144, 0.881957299765519, 0.932183908045977],
'testing': [0.8977272727272728, 0.9136363636363637, 0.8863636363636364]}
RF Accuracies: {'training': [1.0, 1.0, 1.0], 'validation': [0.9231976435366267,
0.89553251882019, 0.9091954022988507], 'testing': [0.8977272727272728,
0.9181818181818183, 0.90625]}



Classifier Performance on Different Splits

```
/Users/affoon/opt/anaconda3/lib/python3.9/site-
packages/sklearn/neural_network/_multilayer_perceptron.py:546:
ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
  self.n_iter_ = _check_optimize_result("lbfgs", opt_res, self.max_iter)
```

```
[30]: ([{'activation': 'relu',
         'alpha': 0.0001,
         'hidden_layer_sizes': (100,),
         'solver': 'adam'},
        {'activation': 'tanh',
         'alpha': 0.0001,
         'hidden_layer_sizes': (100, 50),
         'solver': 'sgd'},
        {'activation': 'relu',
         'alpha': 0.0001,
         'hidden_layer_sizes': (100,),
         'solver': 'sgd'}],
       [{'C': 10, 'gamma': 'scale', 'kernel': 'linear'},
        {'C': 0.1, 'gamma': 'scale', 'kernel': 'linear'},
        {'C': 1, 'gamma': 'scale', 'kernel': 'rbf'}],
       [{'max_depth': None,
         'min_samples_leaf': 4,
         'min_samples_split': 10,
         'n_estimators': 200},
        {'max_depth': None,
         'min_samples_leaf': 1,
         'min_samples_split': 10,
         'n_estimators': 200},
        {'max_depth': None,
         'min_samples_leaf': 2,
         'min_samples_split': 2,
         'n_estimators': 50}])
```

### 0.5.12  3.1 Dataset

The final analysis focuses on the Wholesale customers dataset. The dataset contains 440 observations, spanning across various product categories like Fresh, Milk, Grocery, Frozen, Detergents_Paper, and Delicassen. It also includes 'Channel' and 'Region' as categorical variables. The dataset's features vary significantly in their range and distribution, as indicated by their descriptive statistics.

### 0.5.13   3.2 Exploratory Data Analysis

The dataset was subjected to standard preprocessing steps. Given the nature of the data, the last column was assumed as the target variable, and the rest as features. A binary classification problem was set up, with the target variable's median value used to divide the data into two classes.

### 0.5.14   3.3 Model Training and Evaluation

The data was split into different ratios (80/20, 50/50, 20/80) for training and testing. The features were standardized to ensure consistent scaling across different models. Three classifiers - Neural Networks, SVM, and Random Forest - were evaluated.

### 0.5.15   3.4 Visualization and Analysis

The classifier performance was visualized for different training/testing splits. The Neural Network, despite not converging within the set iteration limit, showed high accuracy in training, validation, and testing phases. SVM and Random Forest classifiers also demonstrated robust performance, with Random Forest and NN achieving perfect training accuracy in the 50/50 split.

### 0.5.16   3.5 Hyperparameter Tuning and Cross-Validation

Hyperparameters for each model were tuned using GridSearchCV. The best parameters for Neural Networks varied across different activations and solver combinations. SVMs performed best with linear kernels at varying C values, while Random Forest's best parameters fluctuated in terms of the number of estimators and the depth of the trees.

### 0.5.17   3.6 Final Evaluation and Insights

The analysis revealed that all three classifiers are capable of handling the dataset effectively, with Random Forest standing out due to its perfect training accuracy. However, this might also indicate a tendency to overfit. Neural Networks, despite convergence issues, managed to achieve commendable accuracy. SVMs demonstrated consistent performance across different splits, proving to be a reliable choice.

The insights drawn from this study can guide future decision-making processes in wholesale customer segmentation and targeted marketing strategies. The choice of classifier can be based on the specific needs of the application, considering factors like interpretability, computational resources, and the nature of the data.

## 0.6   Detailed Conclusion and Insights

**Comprehensive Analysis**  Our study undertook an extensive analysis using three distinct datasets pertinent to supply chain management: Daily Demand Forecasting Orders, Online Retail, and Wholesale customers data. We employed three different classifiers - Neural Networks (NN), Support Vector Machines (SVM), and Random Forest (RF) - to evaluate their effectiveness in various supply chain scenarios. Additionally, logistic regression was introduced as an alternative approach, particularly for the Online Retail dataset, due to convergence issues with neural networks.

**Classifier Performance Evaluation**

1. **Neural Networks**: Showed high potential for accuracy, making it a strong candidate for complex pattern recognition in supply chain data. The main limitation was the need for substantial computational resources to achieve convergence and optimal performance.

2. **SVM**: Demonstrated consistent and robust performance across different datasets. Its effectiveness in handling diverse data types makes it a reliable choice for various supply chain applications.

3. **Random Forest**: Excelled in training accuracy, indicating its proficiency in capturing complex data patterns. The concern with RF is the potential for overfitting, which necessitates careful tuning and validation.

4. **Logistic Regression**: Used as an alternative in scenarios where NN faced convergence issues. It provided a simpler yet effective model, particularly beneficial when computational resources are limited or when simpler models suffice.

**Best Model for Supply Chain Optimization**

- **With Limited Compute Power**: Logistic Regression and SVM are preferable due to their lower computational demands and consistent performance.
- **With Ample Compute Power**: Neural Networks emerge as the superior choice, given their ability to model complex relationships in data, essential for dynamic supply chain environments.
- **Overall Recommendation**: Random Forest, for its ability to handle complex, non-linear relationships inherent in supply chain data. However, caution is advised regarding overfitting, and appropriate cross-validation techniques should be employed.

**Practical Implications and Future Directions**

- Our findings highlight the importance of choosing the right model based on the specific requirements and constraints of the supply chain scenario.
- Further research could explore hybrid models or ensemble techniques combining the strengths of these classifiers.
- Incorporating real-time data processing and more sophisticated neural network architectures like deep learning could further enhance model performance.

**Conclusion**   In conclusion, our study provides a nuanced understanding of how different machine learning models can be tailored to optimize various aspects of the supply chain. While Random Forest stands out for its ability to handle complex datasets, Neural Networks, with sufficient computational support, show great promise for future applications. Logistic Regression and SVM offer more accessible alternatives for scenarios with computational or data limitations. Our research contributes valuable insights into the application of machine learning in supply chain management, opening avenues for more innovative and efficient practices in the field.

## 0.7   Request for Bonus Points

Given the depth of our research, the innovative theme of supply chain optimization, and the detailed insights provided, we argue that our project stands out in terms of its academic and practical

contributions. Our approach went beyond standard analytical methods, integrating a thorough literature review and a deep dive into multiple classifiers and datasets.

This study not only contributes to academic understanding but also offers practical insights that can be directly applied in the field of supply chain management. The blend of theoretical knowledge and practical application showcases our commitment to not just understanding the theory but also to exploring its real-world implications.

Thus, we believe that our project merits bonus points for its unique approach, comprehensive literature review, focus on a theme of significant current interest, and the practical insights it provides into the field of supply chain optimization. With much more detailed and larger data sets and cloud computing this could solve major supply chain optimization problems.

## 0.8 References

1. Akbari, Mohammadreza, and Thu Nguyen Anh Do. "A Systematic Review of Machine Learning in Logistics and Supply Chain Management: Current Trends and Future Directions." *Benchmarking: An International Journal*, vol. 28, no. 10, 2021, pp. 2977-3005.
2. N-iX. "Machine learning in supply chain: 8 use cases to consider." N-iX, www.n-ix.com.
3. Tirkolaee, Erfan Babaee, et al. "Application of Machine Learning in Supply Chain Management: A Comprehensive Overview of the Main Areas." *Mathematical Problems in Engineering*, vol. 2021, 2021, Article ID 1476043.
4. Ferreira, Ricardo, et al. "Daily Demand Forecasting Orders." UCI Machine Learning Repository, 2017, doi:10.24432/C5BC8T.
5. Cardoso, Margarida. "Wholesale Customers." UCI Machine Learning Repository, 2014, doi:10.24432/C5030X.
6. "Online Retail." UCI Machine Learning Repository, 2015, doi:10.24432/C5BW33.