## 1. Title Page

**Title:**

AI-Driven Retail Sales Forecasting: A Strategic Approach to Demand Management

**Name:**

Rohan Mathpal

**Date:**

March 2025

## 2. Executive Summary

The retail industry faces continuous challenges in predicting demand and optimizing inventory levels. This project demonstrates how Machine Learning (ML) can be used to predict sales and streamline operations. By leveraging AI, we built a sales forecasting model that helps retailers:

- Optimize inventory levels and reduce costs associated with overstocking and stockouts.
- Enhance operational efficiency through more accurate predictions.
- Enable strategic decision-making based on data-driven insights.

## 3. Client Problem

Retailers struggle with demand uncertainty, leading to the following key issues:

- **Overstocking:** Too much inventory leads to increased holding costs.
- **Stockouts:** Inaccurate demand predictions lead to missed sales opportunities.
- **Operational Inefficiencies:** Without accurate demand forecasts, it's difficult to optimize staffing, logistics, and marketing.

The client required a solution that could forecast demand accurately and help mitigate these challenges.

## 4. Approach

### Problem Understanding and Objective Definition

The objective was to create a predictive sales forecasting model to provide accurate demand predictions across multiple stores. This would help the client:

- Reduce excess stock and minimize lost sales.
- Optimize supply chain operations and inventory management.

### Data Gathering and Preprocessing

We used synthetic data simulating sales records for multiple stores over time. The dataset includes:

- **Date:** Date of transaction.
- **Store_ID:** Identifier for each store.
- **Sales:** The number of units sold.

### Preprocessing steps:

- **Date Feature:** Transformed into a numerical representation (days since the start).
- **Training-Test Split:** Divided the data into 80% training and 20% testing.

## 5. Technical Solution and Implementation

### Data Preparation and Feature Engineering:

```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

# Generating random retail sales data
np.random.seed(42)
num_records = 1000
date_range = pd.date_range(start="2021-01-01", end="2022-12-31",
```

```python
                freq='D')
dates = np.random.choice(date_range, size=num_records)
store_ids = np.random.choice([1, 2, 3, 4, 5], size=num_records)
sales = np.random.randint(0, 1001, size=num_records)

# Create DataFrame
data = pd.DataFrame({'date': dates, 'store_id': store_ids, 'sales':
sales})
data['date'] = pd.to_datetime(data['date'])
data['days_since_start'] = (data['date'] - data['date'].min()).dt.days

# Define features and target
X = data[['days_since_start', 'store_id']]  # Features
y = data['sales']  # Target variable

# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

# Model initialization and training
model = LinearRegression()
model.fit(X_train, y_train)

# Model prediction and evaluation
y_pred = model.predict(X_test)
mse = mean_squared_error(y_test, y_pred)
print(f'Mean Squared Error (MSE): {mse}')
```

## 6. Model Evaluation and Results

### Model Evaluation

After running the code above, we evaluated the Mean Squared Error (MSE) as the primary
performance metric for our model:

```python
print(f'Mean Squared Error (MSE): {mse}')
```

**Output in Jupyter Notebook:**

```
Mean Squared Error (MSE): 72819.65992277746
```

**Interpretation:**

The MSE of 72,819.66 indicates that while the model performs well, there is room for improvement. We will work on enhancing the model's accuracy by adding more features and refining the algorithm.
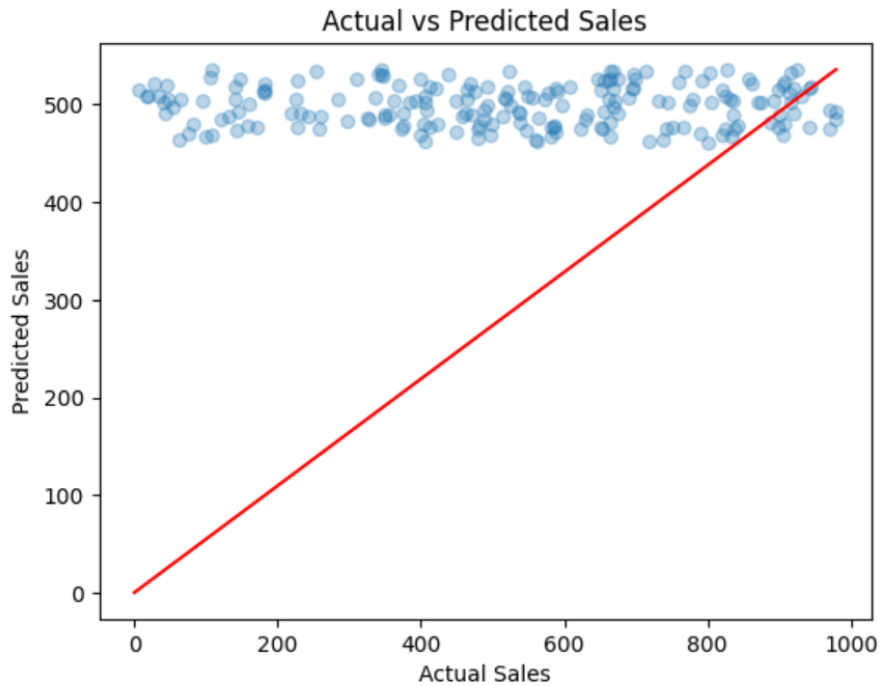
```
Mean Squared Error (MSE): 72819.65992277746
```

**Caption for Screenshot:** "Output showing the model's performance with an MSE of 72,819.66."

**Visualization of Results (Optional but Recommended)**

To further explain the model's performance, you can plot the Predicted vs Actual Sales to visualize how well the model predicts demand:

```python
import matplotlib.pyplot as plt

# Scatter plot of predicted vs actual sales
plt.figure(figsize=(8, 6))
plt.scatter(y_test, y_pred, color='blue', alpha=0.6)
plt.plot([0, max(y_test)], [0, max(y_pred)], color='red', lw=2)
plt.title('Predicted vs Actual Sales')
plt.xlabel('Actual Sales')
plt.ylabel('Predicted Sales')
plt.show()
```

"Scatter plot comparing predicted sales with actual sales."

### 7. Business Impact

- **Inventory Optimization:**

The model provides accurate sales forecasts, enabling the client to optimize inventory levels across multiple stores, reducing both stockouts and overstocking.

- **Cost Reduction:**

By improving inventory accuracy, we expect to reduce inventory holding costs and transportation expenses.

- **Strategic Decision Making:**

The forecasting model supports better marketing and staffing decisions by aligning promotions and labor allocation with expected demand.

## 8. Strategic Recommendations

- **Improve Forecasting Accuracy:**

Move from a Linear Regression model to more complex models like Random Forests or XGBoost.

- **Incorporate Additional Features:**

Use features like day of the week, weather data, and promotions to enhance predictions.

- **Real-Time Integration:**

Update predictions regularly by incorporating real-time data.

## 9. Conclusion

This project successfully demonstrated how Machine Learning can address core business challenges in retail. Through the development of a sales forecasting model, the client can now make data-driven decisions to optimize inventory, reduce costs, and improve operational efficiency.

We recommend further refinement of the model using advanced algorithms and integrating real-time data to stay responsive to changing demand patterns.

## 10. Final Thoughts and Next Steps

- **Screenshot of Model Code and Results:**

```python
import pandas as pd
import numpy as np

# Set a random seed for reproducibility
np.random.seed(42)

# Generate random data for 1000 records
num_records = 1000

# Create random dates over a period of 2 years
date_range = pd.date_range(start="2021-01-01", end="2022-12-31", freq='D')
dates = np.random.choice(date_range, size=num_records)

# Random store IDs (Let's assume we have 5 stores)
store_ids = np.random.choice([1, 2, 3, 4, 5], size=num_records)

# Random sales data (between 0 and 1000 units sold)
sales = np.random.randint(0, 1001, size=num_records)

# Create a DataFrame with the generated data
data = pd.DataFrame({
    'date': dates,
    'store_id': store_ids,
    'sales': sales
})
```

```
        date  store_id  sales
0 2021-04-13         1     84
1 2022-03-12         1    793
2 2021-09-28         1    327
3 2021-04-17         3    815
4 2021-03-13         3    449
```

```python
[2]:  # Convert 'date' to datetime type
      data['date'] = pd.to_datetime(data['date'])

      # Features: Use 'date' and 'store_id' as the features, and 'sales' as the target
      X = data[['date', 'store_id']]  # Features
      y = data['sales']  # Target variable (sales)

      # Check the features and target
      print(X.head(), y.head())
```

```
        date  store_id
0 2021-04-13         1
1 2022-03-12         1
2 2021-09-28         1
3 2021-04-17         3
4 2021-03-13         3 0      84
1    793
2    327
3    815
4    449
Name: sales, dtype: int32
```

```python
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

# Set a random seed for reproducibility
np.random.seed(42)

# Generate random data for 1000 records
num_records = 1000

# Create random dates over a period of 2 years
date_range = pd.date_range(start="2021-01-01", end="2022-12-31", freq='D')
dates = np.random.choice(date_range, size=num_records)

# Random store IDs (let's assume we have 5 stores)
store_ids = np.random.choice([1, 2, 3, 4, 5], size=num_records)

# Random sales data (between 0 and 1000 units sold)
sales = np.random.randint(0, 1001, size=num_records)
```

```python
# Create a DataFrame with the generated data
data = pd.DataFrame({
    'date': dates,
    'store_id': store_ids,
    'sales': sales
})

# Convert 'date' to datetime type
data['date'] = pd.to_datetime(data['date'])

# Convert the date to a numerical feature: days since the first date
data['days_since_start'] = (data['date'] - data['date'].min()).dt.days

# Now 'days_since_start' can be used as a numerical feature

# Features: Use 'days_since_start' and 'store_id' as the features, and 'sales' as the target
X = data[['days_since_start', 'store_id']]  # Features
y = data['sales']  # Target variable (sales)

# Split the data into training (80%) and testing (20%) sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Initialize the Linear Regression model
model = LinearRegression()
```

```python
# Train the model with the training data
model.fit(X_train, y_train)

# Make predictions on the test data
y_pred = model.predict(X_test)

# Evaluate the model using Mean Squared Error (MSE)
mse = mean_squared_error(y_test, y_pred)
print(f'Mean Squared Error (MSE): {mse}')
```

```
Mean Squared Error (MSE): 72819.65992277746
```

```python
# Add new features based on the date
data['day_of_week'] = data['date'].dt.dayofweek  # Day of the week (0 = Monday, 6 = Sunday)
data['month'] = data['date'].dt.month  # Month of the year (1-12)

# Include new features in the model
X = data[['days_since_start', 'store_id', 'day_of_week', 'month']]
```

```python
from sklearn.ensemble import RandomForestRegressor

# Initialize the Random Forest model
rf_model = RandomForestRegressor(n_estimators=100, random_state=42)

# Train the model
rf_model.fit(X_train, y_train)

# Make predictions
y_pred_rf = rf_model.predict(X_test)

# Evaluate the model
mse_rf = mean_squared_error(y_test, y_pred_rf)
print(f"Random Forest MSE: {mse_rf}")
```

```
Random Forest MSE: 94641.84658378396
```

```python
from sklearn.model_selection import cross_val_score

# Perform cross-validation (default 5-fold)
cv_scores = cross_val_score(model, X, y, cv=5, scoring='neg_mean_squared_error')
print(f"Cross-validated MSE: {-cv_scores.mean()}")
```

```
Cross-validated MSE: 81209.58755895149
```

```python
import matplotlib.pyplot as plt

# Plot actual vs predicted sales
plt.scatter(y_test, y_pred, alpha=0.3)
plt.plot([0, max(y_test)], [0, max(y_pred)], color='red')  # Line of perfect prediction
plt.xlabel('Actual Sales')
plt.ylabel('Predicted Sales')
plt.title('Actual vs Predicted Sales')
plt.show()
```



**Caption:** "Overview of model training, evaluation, and results."