# Chocolate Sales Forecasting Project Documentation

## Project Overview

This project implements a comprehensive machine learning pipeline for forecasting chocolate sales revenue using historical sales data. The analysis includes data preprocessing, exploratory data analysis, feature engineering, model development, evaluation, and forecasting capabilities.

## Dataset Description

### Original Data Structure

The dataset contains chocolate sales transaction records with the following columns:

- **Sales Person**: Name of the salesperson handling the transaction
- **Country**: Country where the sale occurred (UK, India, Australia, New Zealand)
- **Product**: Type of chocolate product sold (e.g., Mint Chip Choco, 85% Dark Bars, Peanut Butter Cubes, etc.)
- **Date**: Transaction date
- **Amount**: Sales revenue in USD
- **Boxes Shipped**: Number of product boxes shipped

### Data Characteristics

- Time period: 2022 (January to August based on visible data)
- Geographic coverage: Multiple countries (UK, India, Australia, New Zealand)
- Product variety: Multiple chocolate product types
- Transaction-level granularity with daily sales records

## Methodology

### 1. Data Loading and Initial Inspection

- Import necessary libraries (pandas, numpy, matplotlib, seaborn, scikit-learn)
- Load the chocolate sales dataset

- Perform initial data exploration and structure analysis

- Display basic statistical information about the dataset

## 2. Data Cleaning and Preprocessing

- **Date Conversion**: Convert date strings to datetime format for time series analysis

- **Amount Cleaning**: Remove currency symbols and convert amount to numeric values

- **Duplicate Removal**: Identify and remove duplicate records

- **Data Validation**: Check for missing values and data inconsistencies

## 3. Exploratory Data Analysis (EDA)

The EDA phase includes comprehensive analysis across multiple dimensions:

### 3.1 Distribution Analysis

- Sales amount distribution analysis

- Boxes shipped distribution patterns

- Statistical summary of key metrics

### 3.2 Time Series Analysis

- Monthly sales trends and patterns

- Revenue evolution over time

- Seasonal pattern identification

### 3.3 Dimensional Analysis

- **Sales Performance by Person**: Individual salesperson performance metrics

- **Geographic Analysis**: Sales performance by country

- **Product Analysis**: Revenue contribution by product type

### 3.4 Seasonal Patterns

- Monthly sales pattern identification

- Quarterly performance analysis

- Day-of-week effects on sales

## 3.5 Correlation Analysis

- Relationship between boxes shipped and revenue amount

- Feature correlation matrix analysis

## 4. Feature Engineering

Advanced feature creation for improved model performance:

## Time-Based Features

- **Month**: Extracted from transaction date

- **Year**: Calendar year of transaction

- **Quarter**: Quarterly categorization (Q1-Q4)

- **Day_of_week**: Day of the week (0-6)

- **Day_of_month**: Day within the month (1-31)

- **Week_of_year**: Week number within the year (1-52)

## Lag Features for Time Series

- **Amount_lag1**: Sales amount from 1 day ago

- **Amount_lag7**: Sales amount from 7 days ago

## Rolling Window Features

- **Amount_rolling_7**: 7-day rolling average of sales amount

- **Amount_rolling_30**: 30-day rolling average of sales amount

## Categorical Encoding

- One-hot encoding for categorical variables (Sales Person, Country, Product)

## 5. Data Preparation for Modeling

- **Temporal Ordering**: Maintain chronological order for time series modeling

- **Train-Test Split**: Use 80% of data for training, 20% for testing (temporal split)

- **Feature Selection**: Prepare feature matrix and target variable

- **Data Scaling**: Normalize features for model compatibility

## 6. Model Development

Implementation of multiple machine learning algorithms:

## Models Implemented

1. **Linear Regression**: Baseline model for linear relationships

2. **Random Forest**: Ensemble method for capturing complex patterns

3. **Gradient Boosting**: Advanced boosting algorithm

4. **XGBoost**: Optimized gradient boosting framework

## Model Configuration

- Cross-validation for hyperparameter tuning

- Grid search for optimal parameter selection

- Ensemble methods for improved accuracy

## 7. Model Evaluation

Comprehensive evaluation using multiple metrics:

## Evaluation Metrics

- **MAE (Mean Absolute Error)**: Average absolute difference between predicted and actual values

- **RMSE (Root Mean Square Error)**: Square root of average squared differences

- **MAPE (Mean Absolute Percentage Error)**: Percentage-based error metric

## Performance Results

Based on the evaluation results:

- **Linear Regression**: MAE: 3,210.97, RMSE: 3,918.15, MAPE: 3.40%

- **Random Forest**: MAE: 3,291.78, RMSE: 4,049.39, MAPE: 3.66%

- **Gradient Boosting**: MAE: 3,366.94, RMSE: 4,161.39, MAPE: 3.69%

- **XGBoost**: MAE: 3,539.12, RMSE: 4,341.09, MAPE: 3.96%

**Best Performing Model**: Linear Regression (lowest MAE and MAPE)

## 8. Model Interpretation

- **Feature Importance Analysis**: Identify key drivers of sales performance

- **Model Explainability**: Understanding which features contribute most to predictions

- **Business Insights**: Extract actionable insights from model results

## 9. Forecasting Implementation

### Forecast Configuration

- **Forecast Period**: September 2022 (30-day forecast)

- **Prediction Intervals**: Confidence bounds for predictions

- **Model Selection**: Best performing model used for forecasting

### Forecast Results

Daily sales predictions for September 2022:

- **Range**: $3,987.97 - $5,411.70

- **Average Daily Forecast**: ~$4,800

- **Confidence Intervals**: Upper and lower bounds provided for each prediction

### Technical Implementation

### Libraries and Dependencies

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from datetime import datetime
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
```

```
from sklearn.ensemble import RandomForestRegressor, GradientBoostingRegressor
from xgboost import XGBRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error
import warnings
warnings.filterwarnings('ignore')
```

## Data Processing Pipeline

1. **Data Loading**: CSV import and initial structure analysis

2. **Cleaning**: Date parsing, amount conversion, duplicate removal

3. **Feature Engineering**: Time features, lag variables, rolling statistics

4. **Encoding**: Categorical variable transformation

5. **Modeling**: Multiple algorithm implementation and evaluation

6. **Forecasting**: Future period prediction with confidence intervals

## Key Findings and Insights

## Model Performance

- Linear Regression achieved the best performance with 3.40% MAPE

- All models demonstrate reasonable accuracy for business forecasting needs

- Feature engineering significantly improved model performance

## Forecast Insights

- September 2022 daily sales forecast averages around $4,800

- Sales show typical business patterns with weekday/weekend variations

- Confidence intervals provide risk assessment for business planning

## Business Implications

- Reliable forecasting capability enables better inventory planning

- Sales performance analysis identifies top-performing products and regions

- Time series patterns support seasonal business strategy development

**Future Enhancements**

**Model Improvements**

- Deep learning models for complex pattern recognition

- External factors integration (holidays, marketing campaigns)

- Multi-step ahead forecasting capabilities

**Feature Engineering**

- Advanced lag features and seasonal decomposition

- External economic indicators integration

- Customer behavior and market trend features

**Deployment Considerations**

- Model versioning and monitoring system

- Automated retraining pipeline

- Real-time prediction API development

- Business intelligence dashboard integration

**Conclusion**

This chocolate sales forecasting project demonstrates a comprehensive approach to time series forecasting using machine learning. The implemented solution provides accurate sales predictions with confidence intervals, enabling data-driven business decision making. The modular approach allows for easy enhancement and adaptation to changing business requirements.

The linear regression model's superior performance suggests that the underlying relationships in chocolate sales data follow relatively linear patterns, making it an efficient and interpretable choice for business forecasting applications.