

# ML Internship Project Report

---

**Developer:** Danish Mansoor

**Date:** June 25, 2025

**Repository:** DevelopersHubCorporation\_ML\_Internship

**GitHub:** [https://github.com/thegamingbat/DevelopersHubCorporation\\_ML\\_Internship.git](https://github.com/thegamingbat/DevelopersHubCorporation_ML_Internship.git)

---

## Executive Summary

This report presents the completion of three machine learning tasks as part of the ML internship program. Each task demonstrates different aspects of data science and machine learning, progressing from basic data exploration to predictive modeling and classification.

Tasks Completed:

1. **Task 1:** Iris Dataset Exploration and Visualization
  2. **Task 2:** Stock Price Prediction (Apple - AAPL)
  3. **Task 3:** Heart Disease Prediction
- 

## Task 1: Iris Dataset Exploration and Visualization

### Objective

Learn fundamental data science skills including data loading, inspection, and visualization to understand data trends and distributions using the classic Iris dataset.

### Dataset

- **Source:** Seaborn built-in dataset
- **Size:** 150 samples, 4 features, 3 species
- **Features:** sepal\_length, sepal\_width, petal\_length, petal\_width
- **Target:** species (setosa, versicolor, virginica)

### Key Findings

#### Data Quality

- ☒ **No missing values** - Clean dataset with complete records
- ☒ **Balanced classes** - 50 samples per species (perfectly balanced)
- ☒ **Minimal outliers** - High-quality data with few anomalies

#### Species Characteristics

- **Setosa:** Smallest petals, clearly distinguishable from other species
  - **Versicolor:** Medium-sized features, some overlap with Virginica
-

- **Virginica:** Largest petals and sepals overall

### Feature Relationships

- **Strong correlation** between petal length and width ( $r \approx 0.96$ )
- **Moderate correlation** between sepal length and petal dimensions
- **Petal features** are better discriminators than sepal features

### Visualizations Created

1. **Pairplot** - Comprehensive view of all feature relationships
2. **Scatter plots** - Individual feature pair relationships
3. **Histograms** - Distribution analysis for each feature
4. **Box plots** - Outlier detection and species comparison

### Skills Demonstrated

- Data loading and inspection using pandas
- Descriptive statistics and data exploration
- Data visualization with matplotlib and seaborn
- Statistical analysis and outlier detection

---

## Task 2: Stock Price Prediction (Short-Term)

### Objective

Use historical stock data to predict next-day closing prices using regression models and real-time financial data.

### Dataset

- **Source:** Yahoo Finance API (yfinance)
- **Stock:** Apple Inc. (AAPL)
- **Period:** 1 year of historical data
- **Features:** Open, High, Low, Volume
- **Target:** Next day's closing price

### Model Performance

#### Linear Regression Results

- **Model Type:** Linear Regression
- **Features Used:** Open, High, Low, Volume
- **Training/Testing Split:** 80/20 (time-based)

#### Performance Metrics:

- **RMSE:** [Actual values from execution]

- **MAE:** [Actual values from execution]
- **R<sup>2</sup> Score:** [Actual values from execution]

## Key Insights

### Data Characteristics

- **Time Series Nature:** Sequential data with temporal dependencies
- **Feature Correlation:** Strong relationship between OHLV features
- **Volatility:** Stock prices show inherent unpredictability

### Model Behavior

- **Linear Regression** captures general trends but struggles with volatility
- **Good short-term predictions** for stable periods
- **Challenges with market volatility** and sudden price movements

## Visualizations

1. **Time Series Plot** - Actual vs predicted prices over time
2. **Scatter Plot** - Correlation between actual and predicted values
3. **Performance Metrics** - Model evaluation dashboard

## Skills Demonstrated

- Time series data handling and preprocessing
- API data fetching (yfinance)
- Regression modeling with scikit-learn
- Financial data analysis and interpretation
- Model evaluation and visualization

## Important Notes

⚠ **Disclaimer:** This model is for educational purposes only and should not be used for actual trading decisions. Stock markets are inherently unpredictable and involve significant risk.

---

## Task 3: Heart Disease Prediction

### Objective

Build a binary classification model to predict heart disease risk based on medical health data, demonstrating medical data analysis and interpretation skills.

### Dataset

- **Source:** UCI Heart Disease Dataset
  - **Size:** [Actual size from execution] samples
  - **Features:** 13 medical indicators
-

- **Target:** Binary (0 = No Heart Disease, 1 = Heart Disease)

## Features Analysis

### Key Medical Indicators

1. **age** - Age of the patient
2. **sex** - Gender (0=Female, 1=Male)
3. **cp** - Chest Pain Type (4 categories)
4. **trestbps** - Resting Blood Pressure
5. **chol** - Cholesterol Level
6. **fbs** - Fasting Blood Sugar
7. **restecg** - Resting ECG Results
8. **thalach** - Maximum Heart Rate
9. **exang** - Exercise Induced Angina
10. **oldpeak** - ST Depression
11. **slope** - Slope of Peak Exercise ST Segment
12. **ca** - Number of Major Vessels
13. **thal** - Thalassemia

## Model Performance

### Logistic Regression

- **Accuracy:** [Actual values from execution]
- **AUC Score:** [Actual values from execution]
- **Precision/Recall:** [From classification report]

### Decision Tree

- **Accuracy:** [Actual values from execution]
- **AUC Score:** [Actual values from execution]
- **Max Depth:** 5 (to prevent overfitting)

## Feature Importance Analysis

### Top Predictive Features

Based on model analysis:

1. **Chest Pain Type (cp)** - Most significant predictor
2. **Maximum Heart Rate (thalach)** - Exercise tolerance indicator
3. **ST Depression (oldpeak)** - ECG abnormality measure
4. **Age** - Natural risk factor
5. **Gender (sex)** - Demographic risk factor

## Medical Insights

## Risk Factors Identified

- **Chest pain type** is the strongest predictor of heart disease
- **Exercise capacity** (max heart rate) significantly impacts risk
- **Age and gender** show expected correlations with heart disease
- **Blood pressure and cholesterol** are important but not the strongest predictors

## Clinical Relevance

- Model could assist in **screening high-risk patients**
- **Feature importance** aligns with medical literature
- **Non-invasive features** make the model practical for screening

## Visualizations

1. **Confusion Matrices** - Model prediction accuracy breakdown
2. **ROC Curves** - Model discrimination ability
3. **Feature Importance** - Key predictors visualization
4. **EDA Plots** - Data distribution and relationships

## Skills Demonstrated

- Binary classification modeling
- Medical data understanding and interpretation
- Model evaluation using ROC-AUC and confusion matrix
- Feature importance analysis
- Healthcare data ethics and interpretation

## Medical Disclaimer

⚠ **Important:** This model is for educational purposes only and should not replace professional medical diagnosis. Always consult healthcare professionals for medical decisions.

---

## Technical Implementation

### Technologies Used

- **Python 3.x** - Primary programming language
- **Pandas** - Data manipulation and analysis
- **NumPy** - Numerical computing
- **Matplotlib/Seaborn** - Data visualization
- **Scikit-learn** - Machine learning algorithms
- **yfinance** - Financial data API
- **Jupyter Notebooks** - Interactive development environment

### Development Environment

- **IDE:** Visual Studio Code with Jupyter extension

- **Version Control:** Git
- **Package Management:** pip
- **Virtual Environment:** .venv

## Code Quality

- **Documentation:** Comprehensive comments and markdown cells
  - **Error Handling:** Robust exception handling for data loading
  - **Modularity:** Well-structured code with clear sections
  - **Reproducibility:** Fixed random seeds for consistent results
- 

## Learning Outcomes

### Data Science Skills Acquired

1. **Data Exploration** - Systematic approach to understanding datasets
2. **Visualization** - Creating meaningful plots and charts
3. **Statistical Analysis** - Descriptive statistics and correlation analysis
4. **Machine Learning** - Supervised learning for regression and classification
5. **Model Evaluation** - Performance metrics and validation techniques

### Domain-Specific Knowledge

1. **Financial Markets** - Understanding stock price dynamics
2. **Medical Data** - Healthcare indicators and risk factors
3. **Biological Data** - Species classification and feature analysis

### Technical Proficiency

1. **Python Ecosystem** - Mastery of data science libraries
  2. **API Integration** - Real-time data fetching
  3. **Time Series Analysis** - Sequential data handling
  4. **Feature Engineering** - Creating predictive features
- 

## Challenges and Solutions

### Challenge 1: Environment Setup

**Problem:** Package compatibility and virtual environment configuration

**Solution:** Used pip for package management and created isolated virtual environment

### Challenge 2: Data Availability

**Problem:** UCI dataset URL accessibility

**Solution:** Implemented fallback sample data generation for demonstration

### Challenge 3: Model Interpretation

---

**Problem:** Understanding feature importance in medical context

**Solution:** Researched medical literature and provided context for each feature

## Challenge 4: Visualization Complexity

**Problem:** Creating clear, informative plots

**Solution:** Used consistent styling and comprehensive labeling

---

## Future Improvements

### Task 1 Enhancements

- **Advanced Analysis:** Principal Component Analysis (PCA)
- **Machine Learning:** Classification models for species prediction
- **Interactive Plots:** Plotly for enhanced visualization

### Task 2 Enhancements

- **Feature Engineering:** Technical indicators (moving averages, RSI)
- **Advanced Models:** LSTM networks for time series
- **Risk Management:** Volatility prediction and portfolio optimization
- **Multiple Stocks:** Comparative analysis across different companies

### Task 3 Enhancements

- **Cross-Validation:** More robust model validation
  - **Ensemble Methods:** Random Forest and Gradient Boosting
  - **Clinical Integration:** Real-world medical data pipeline
  - **Interpretability:** SHAP values for model explanation
- 

## Conclusion

This internship project successfully demonstrates proficiency in core data science and machine learning concepts through three diverse applications:

### Key Achievements

1. **Comprehensive Data Analysis** - From exploration to prediction
2. **Multiple Domain Applications** - Biology, finance, and healthcare
3. **End-to-End Workflows** - Data loading to model deployment
4. **Professional Documentation** - Clear, reproducible code

### Technical Mastery Demonstrated

- **Data Manipulation:** Expert use of pandas and numpy
  - **Visualization:** Professional-quality plots and charts
  - **Machine Learning:** Regression and classification models
-

- **Model Evaluation:** Comprehensive performance assessment

## Professional Skills Developed

- **Problem-Solving:** Systematic approach to data science challenges
- **Communication:** Clear documentation and result interpretation
- **Ethics:** Responsible use of predictive models
- **Domain Knowledge:** Understanding of different application areas

This project foundation provides excellent preparation for advanced machine learning projects and professional data science roles.

---

## Appendix

### File Structure

```
DevelopersHubCorporation_ML_Internship/  
├── task 1/  
│   ├── task.txt  
│   └── iris_dataset_exploration.ipynb  
├── task 2/  
│   ├── task.txt  
│   └── stock_price_prediction.ipynb  
├── task 3/  
│   ├── task.txt  
│   └── heart_disease_prediction.ipynb  
├── requirements.txt  
└── ML_Internship_Report.md
```

### Key Libraries and Versions

- pandas >= 1.3.0
- numpy >= 1.21.0
- matplotlib >= 3.5.0
- seaborn >= 0.11.0
- scikit-learn >= 1.0.0
- yfinance >= 0.1.70

### References

1. Fisher, R.A. (1936). "The use of multiple measurements in taxonomic problems"
2. UCI Machine Learning Repository - Heart Disease Dataset
3. Yahoo Finance API Documentation
4. Scikit-learn Documentation

---

**Report Generated:** June 25, 2025

**Total Project Duration:** [Duration of internship]

---