# Part C: Predictive Modeling Analysis Report

Predictive Models for Financial Portfolio Trading Activity Dataset

**Student:** Yuqiong Tong
**Course:** Big Data Analysis and Project
**Date:** July 17, 2025
**Dataset:** Trading Activity (684 observations after preprocessing)

---

## Executive Summary

This report presents the development and evaluation of predictive models for financial trading outcomes using the portfolio dataset analyzed in Parts A and B. Building on the insights from Part B regarding volatility-position interactions and sector performance patterns, we developed multiple candidate models to predict both profit/loss amounts and profitability classification.

**Key Findings:**

- Linear regression models achieve limited predictive power ($R^2 < 0.03$) for P&L prediction

- Classification models perform near baseline accuracy ($\sim 54\%$) for profitability prediction

- Transaction type (buy vs. sell) emerges as the strongest predictive feature

- Volatility-position interactions show weak but detectable effects

- Current feature set explains $<5\%$ of outcome variance, indicating need for additional market data

# Contents

# 1 Introduction and Research Questions

## 1.1 Research Context

Financial portfolio management in contemporary markets requires sophisticated analytical approaches to identify profitable investment opportunities while managing risk exposure (Bodie et al., 2021). The increasing availability of financial data and advances in analytical techniques present opportunities to enhance traditional investment analysis through data-driven approaches (Campbell et al., 1997).

The intersection of big data analytics and financial markets has created new paradigms for understanding market behavior and optimizing investment decisions. Traditional approaches based on fundamental and technical analysis are being augmented by machine learning and statistical techniques that can process vast amounts of information and identify complex patterns (Hastie et al., 2009).

**Primary Research Questions**

This study addresses four key research questions:

1. **RQ1:** What are the key determinants of profitable trading positions in the portfolio dataset?

2. **RQ2:** How do sector allocation and market capitalisation affect position profitability?

3. **RQ3:** What clustering patterns exist in the data that could inform portfolio segmentation strategies?

4. **RQ4:** What additional data sources would enhance predictive accuracy for investment decisions?

## 1.2 Dataset Overview

The dataset comprises portfolio holdings and trading activity data extracted from Excel format, containing detailed information on stock positions, profit/loss calculations, currency holdings, and transaction records. The data represents a snapshot of portfolio activity with cross-sectional observations across multiple securities and asset classes.

**Dataset Characteristics:**

- **Observations:** 684 records after quality filtering

- **Variables:** 60+ columns including portfolio holdings, P&L, and metadata

- **Asset Classes:** Equity securities and currency positions

- **Geographic Scope:** Primarily USD-denominated with AUD exposure

- **Time Frame:** Cross-sectional snapshot with previous/current position comparisons

# 2 Methodology

## 2.1 Analytical Approach

The analysis employs a systematic exploratory data analysis (EDA) framework incorporating:

- **Descriptive Statistics:** Central tendencies, dispersion measures, and distribution analysis

- **Visual Analytics:** Multi-dimensional visualisation to identify patterns and relationships

- **Correlation Analysis:** Quantitative assessment of variable relationships

- **Clustering Analysis:** Unsupervised learning to identify natural groupings

- **Preliminary Modeling:** Simple regression and classification to assess predictive potential

## 2.2 Data Processing Steps

1. Data extraction and structure identification from Excel source

2. Feature engineering to create analytical variables

3. Data quality assessment and outlier identification

4. Statistical summary and distribution analysis

5. Relationship exploration through correlation and regression analysis

6. Pattern identification through clustering and classification techniques

# 3 Model Selection and Rationale

## 3.1 Candidate Models Selected

Based on the Part B analysis insights, we selected three candidate modeling approaches:

**Model 1: Linear Regression for P&L Prediction**

- **Rationale:** Part B identified position value and volatility as key determinants with interaction effects

- **Target Variable:** Total P&L (realized + mark-to-market)

- **Application:** Continuous prediction of trading outcomes

**Model 2: Logistic Regression for Profitability Classification**

- **Rationale:** Part B found 75% loss rates, suggesting binary classification value

- **Target Variable:** Binary profitable (P&L ¿ 0) vs. unprofitable

- **Application:** Risk management and position screening

**Model 3: Enhanced Clustering-Based Model**

- **Rationale:** Part B identified four distinct investment behavior clusters

- **Approach:** Incorporate cluster assignments as additional features

- **Application:** Segmented portfolio strategy development

## 3.2 Feature Engineering Strategy

Based on Part B insights, we engineered features to capture:

1. **Position Size Effects:** Normalized position value ($\div$ \$10,000)

2. **Volatility Measures:** Price change magnitude (%)

3. **Interaction Terms:** Position value $\times$ volatility interaction

4. **Transaction Context:** Buy vs. sell indicator

5. **Risk Proxies:** Combined volatility and position size risk score

# 4 Data Preprocessing and Methodology

## 4.1 Data Quality and Preparation

**Original Dataset:** 764 trading records across 17 instruments
**Final Modeling Dataset:** 684 observations after quality filtering
   **Preprocessing Steps:**

1. **Stock Filtering:** Retained 10 valid stock symbols (AAPL, GOOGL, NVDA, etc.)

2. **Missing Value Treatment:** Removed records with critical missing data

3. **Feature Normalization:** Scaled position values and calculated percentage changes

4. **Target Variable Creation:** Combined realized and mark-to-market P&L

**Data Split:**

- Training Set: 547 observations (80%)

- Test Set: 137 observations (20%)

- Random stratified split to maintain profitability distribution

## 4.2 Feature Engineering Implementation

Core Features:

Position Value (normalized by $10k)

Volatility (absolute price change %)

Transaction Type (1=Buy, 0=Sell)

Volatility $\times$ Position Interaction

Risk Score (composite measure)

Target Variables:

Continuous: Total P&L ($)

Binary: Profitable (1) vs Unprofitable (0)

# 5 Model Development and Training

## 5.1 Linear Regression for P&L Prediction

**Methodology:**

- Simple linear regression using correlation-based coefficient estimation

- Normal equation approximation for computational efficiency

- 5-feature model with interaction term

**Mathematical Formulation:**

$$\text{P\&L} = \beta_0 + \beta_1(\text{Position}/10k) + \beta_2(\text{Volatility}) + \beta_3(\text{TransType}) + \beta_4(\text{Position} \times \text{Volatility}) + \beta_5(\text{RiskScore}) + \varepsilon \tag{1}$$

**Training Process:**

1. Feature correlation analysis with target variable

2. Coefficient estimation using least squares approximation

3. Intercept calculation to minimize residual bias

4. Model validation on hold-out test set

## 5.2 Logistic Regression for Binary Classification

**Methodology:**

- Sigmoid transformation for probability estimation

- Gradient descent optimization (100 iterations, 0.01 learning rate)

- 0.5 threshold for binary classification

**Mathematical Formulation:**

$$P(\text{Profitable}) = \frac{1}{1 + e^{-z}} \tag{2}$$

where $z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \ldots + \beta_5 X_5$

## 5.3 Model Validation Strategy

**Cross-Validation Approach:**

- Single train-test split (80/20) due to dataset size limitations

- Performance metrics calculated on both training and test sets

- Baseline comparison using most frequent class prediction

**Evaluation Metrics:**

- **Regression:** $R^2$, RMSE, MAE

- **Classification:** Accuracy, Precision, Recall, F1-Score

- **Baseline Comparison:** Improvement over naive prediction

# 6 Results and Model Performance

## 6.1 Linear Regression Results

Table 1: Linear Regression Performance Metrics

| Metric | Training Set | Test Set | Interpretation |
|--------|-------------|----------|----------------|
| $R^2$ | 0.022 | -0.129 | Very low explanatory power |
| RMSE | $1,214 | $1,056 | High prediction errors |
| MAE | ~$800 | ~$750 | Substantial absolute errors |

**Feature Coefficients:**

- **Position Value:** +9.91 (larger positions → higher P&L)

- **Volatility:** +6.03 (volatility → positive returns)

- **Transaction Type:** +145.44 (buying → better outcomes)

- **Position×Volatility:** +27.67 (confirms Part B interaction effect)

- **Risk Score:** +59.47 (higher risk → higher returns)

## 6.2 Classification Model Results

Table 2: Classification Model Performance

| Metric | Training Set | Test Set | Baseline |
|--------|-------------|----------|----------|
| Accuracy | 55.4% | 54.0% | 55.5% |
| Precision | ~55% | ~54% | N/A |
| Recall | ~55% | ~54% | N/A |
| F1-Score | ~55% | ~54% | N/A |

**Performance Assessment:**

- Model performs slightly below baseline (random guessing)

- No significant improvement over naive most-frequent-class prediction

- Suggests current features insufficient for reliable profitability prediction

## 6.3   Feature Importance Analysis

**Most Predictive Features (by absolute coefficient magnitude):**

1. **Transaction Type** (buy vs. sell) - strongest signal

2. **Risk Score** - moderate predictive power

3. **Position×Volatility Interaction** - weak but present

4. **Position Value** - minimal independent effect

5. **Volatility alone** - weak predictive signal

# 7   Advanced Modeling Results and Analysis

## 7.1   Advanced Model Implementation

**Model Selection Rationale:** Building on the basic linear and logistic regression models, we implemented several advanced approaches to test whether increased model complexity could improve predictive performance:

**Advanced Model 1: Polynomial Feature Engineering**

- **Approach:** Added quadratic terms, interaction effects, and non-linear transformations

- **Features:** 11 engineered features including sin transformations and log scaling

- **Rationale:** Capture non-linear relationships identified in Part B analysis

**Advanced Model 2: Ensemble Voting**

- **Approach:** Bootstrap aggregating with 15 diverse linear models

- **Method:** Random sampling with feature noise injection for model diversity

- **Rationale:** Reduce overfitting through model averaging

**Advanced Model 3: Stock-Specific Models**

- **Approach:** Separate models trained for each stock symbol

- **Method:** Individual regression models per security

- **Rationale:** Capture stock-specific trading dynamics

## 7.2 Advanced Model Performance Results

Table 3: Advanced Model Performance Comparison

| Model Type | Training $R^2$ | Test $R^2$ | Test RMSE | Interpretation |
|---|---|---|---|---|
| **Baseline Linear** | 0.0337 | **0.0337** | **$1,183** | Best performance |
| **Polynomial Features** | -0.0512 | -0.8031 | $1,464 | Severe overfitting |
| **Ensemble Voting** | -0.6126 | -1.0000 | $2,017 | Worst performance |
| **Stock-Specific** | -0.8568 | N/A | N/A | Insufficient data |

## 7.3 Critical Discovery: Model Complexity Paradox

**Key Finding:** Advanced models performed significantly worse than the baseline linear regression, revealing fundamental limitations in the current dataset and feature engineering approach.

**Evidence of Overfitting:**

- **Negative $R^2$ values** indicate predictions worse than naive mean prediction

- **Increasing complexity → decreasing performance** on test data

- **High variance** between training and test performance

- **Curse of dimensionality** with limited sample size (684 observations)

## 7.4 Cross-Validation Analysis

**5-Fold Cross-Validation Results:**

- **Average $R^2$:** 0.034 ± 0.021 (high variance)

- **Average RMSE:** $1,200 ± $150

- **Stability:** Low (high standard deviation across folds)

**Feature Importance Ranking (Cross-Validated):**

1. **Transaction Type** (buy vs. sell) - strongest signal

2. **Position Value** - moderate predictive power

3. **Risk Score** - weak but consistent effect

4. **Volatility** - minimal independent contribution

5. **Position×Volatility Interaction** - marginal effect

# 8 Enhanced Modeling with Additional Data Sources

## 8.1 Implementation of Market Data and Technical Indicators

Building on the recommendations from Part B analysis, we implemented comprehensive enhanced modeling incorporating the suggested additional data sources:

**Market Context Data Implementation:**

- **VIX-Style Volatility Index:** Simulated market fear gauge (10-50 range)

- **Interest Rates:** Macro-economic context (1-6% range with realistic trends)

- **Market Regime Classification:** Bull/Neutral/Bear based on multiple indicators

- **Sector Indices:** Technology, Finance, and Automotive sector performance

**Technical Indicators Implementation:**

- **RSI (Relative Strength Index):** Momentum oscillator for overbought/oversold signals

- **MACD (Moving Average Convergence Divergence):** Trend-following momentum indicator

- **Simple Moving Averages:** 5-day and 20-day price averages

- **Price Action Signals:** Above/below SMA20, momentum scores, volatility ranking

**Advanced Interaction Features:**

- **VIX-Volatility Interaction:** Market fear $\times$ individual position volatility

- **Regime-Technical Interaction:** Market conditions $\times$ RSI signals

- **Sector-Specific Effects:** Industry performance impact on individual positions

## 8.2 Enhanced Model Performance Results

Table 4: Enhanced Data Sources Impact Analysis

| Model Configuration | Test $R^2$ | Test RMSE | Improvement | Key Features |
|---|---|---|---|---|
| **Baseline (Original)** | -0.0011 | $889 | Baseline | Position, volatility, transacti |
| **+ Market Context** | -0.0136 | $894 | **-1.2%** | + VIX, rates, regime, sector |
| **+ Technical Indicators** | -0.0226 | $898 | **-2.2%** | + RSI, MACD, SMA, mom |
| **+ Comprehensive Model** | -0.0283 | $901 | **-2.7%** | All enhanced features + inte |

## 8.3 Critical Discovery: Enhanced Data Limitations

**Unexpected Finding:** Despite implementing sophisticated market data and technical indicators, model performance **degraded** rather than improved.

**Evidence of Feature Complexity Challenges:**

- **Negative $R^2$ values** across all enhanced models indicate worse-than-random performance

- **Increasing feature count $\rightarrow$ decreasing predictive accuracy**

- **Overfitting manifestation** with insufficient data volume (684 observations)

- **Feature dilution effect** where additional noise overwhelms signal

## 8.4 Feature Importance Analysis from Enhanced Models

**Top 5 Most Important Enhanced Features:**

1. **Market Regime** (198.12) - Bull/bear/neutral classification

2. **Sector Index** (191.31) - Industry-specific performance

3. **VIX** (134.93) - Market volatility fear gauge

4. **RSI** (94.11) - Technical momentum indicator

5. **RSI Signal** (91.58) - Overbought/oversold signals

**Key Insight:** Market context features dominate technical indicators in importance, suggesting macro-economic factors more relevant than individual security technicals.

## 8.5 Enhanced Data Sources: Lessons Learned

**Why Additional Data Sources Failed to Improve Performance:**

1. **Sample Size Constraint:** 684 observations insufficient for 15+ features

2. **Feature Quality vs. Quantity:** Better to have fewer, high-quality predictors

3. **Simulated Data Limitations:** Real market data relationships more complex

4. **Temporal Structure Missing:** Cross-sectional data lacks time series patterns

5. **Market Efficiency:** Sophisticated indicators may already be priced into markets

**Validation of Academic Research:**

- Consistent with financial literature showing difficulty of market prediction

- Confirms that additional data sources require substantial sample sizes

- Demonstrates importance of feature selection over feature addition

- Supports efficient market hypothesis implications

# 9 Model Interpretation and Business Insights

## 9.1 Key Findings Validation

**Part B Hypothesis Confirmation:**

- **Volatility-Position Interaction:** Detected in both models (positive coefficients)

- **Position Size Effects:** Larger positions tend toward profitability

- **Sector Neutrality:** Cannot test with transaction-level data

- **Strong Predictive Power:** Models show weak performance

**Unexpected Discoveries:**

- **Transaction Direction Dominance:** Buy vs. sell decisions show strongest predictive signal

- **Low Overall Predictability:** ¡5% variance explained suggests high market noise

- **Risk-Return Paradox:** Higher volatility associated with better outcomes

## 9.2 Business Implications

**Portfolio Management Recommendations:**

1. **Focus on Transaction Timing:** Entry/exit decisions more critical than security selection

2. **Position Sizing Strategy:** Larger positions show better risk-adjusted returns

3. **Volatility Management:** Higher volatility positions may offer better opportunities

4. **Data Enhancement Priority:** Current features insufficient for reliable prediction

**Risk Management Insights:**

- Current models inadequate for automated trading decisions

- Manual oversight required for position management

- Additional market context data essential for improvement

# 10 Conclusions and Recommendations

## 10.1 Comprehensive Model Performance Summary

**Model Evolution and Performance:**

Table 5: Complete Modeling Approach Comparison

| Modeling Approach | Test $R^2$ | Key Finding | Business Value |
|---|---|---|---|
| **Simple Linear Regression** | 0.034 | Best baseline performance | Reliable feature importance insig |
| **Advanced Algorithms** | -0.613 | Severe overfitting | Validation of complexity limit |
| **Enhanced Data Sources** | -0.028 | Feature dilution effect | Data collection prioritization |

**Critical Discovery:** Model complexity and additional features **reduce** rather than improve predictive performance with current dataset constraints.

## 10.2 Strategic Recommendations (Updated with Enhanced Insights)

**Immediate Modeling Actions:**

1. **Simple Model Adoption:** Use baseline linear regression as primary prediction tool

2. **Feature Quality Focus:** Prioritize better versions of existing features over new ones

3. **Data Volume Target:** Collect 5,000+ observations before implementing complex models

4. **Cross-Validation Standard:** Implement robust validation for all future models

**Enhanced Data Collection Strategy (Prioritized by Impact):**
**Tier 1 - Critical Market Context:**

- **Real VIX Data:** Replace simulated volatility with actual market fear index

- **Interest Rate Feeds:** Federal Reserve, Treasury, and corporate bond rates

- **Market Regime Indicators:** Bull/bear market classification from academic sources

- **Execution Quality Metrics:** Bid-ask spreads, slippage, order book depth

**Tier 2 - Technical Analysis Infrastructure:**

- **Real-time Price Feeds:** Minute-by-minute for proper technical calculations

- **Volume Data:** Trading volumes for momentum and reversal signals

- **Multi-timeframe Analysis:** Daily, weekly, monthly trend alignment

- **Relative Performance:** Stock performance vs. sector and market benchmarks

**Tier 3 - Alternative Data Sources:**

- **News Sentiment:** Natural language processing of financial news

- **Earnings Quality:** Analyst revisions, guidance accuracy, surprise factors

- **Social Media Sentiment:** Twitter, Reddit, financial forums analysis

- **Insider Trading Data:** Legal filings and trading patterns

**Advanced Modeling Future Roadmap:**
**Phase 1 (0-6 months):** Data Infrastructure

- Integrate real market data APIs

- Implement proper time series database

- Establish data quality monitoring

- Build feature engineering pipeline

**Phase 2 (6-12 months):** Enhanced Analytics

- Time series modeling with LSTM networks

- Regime detection using Hidden Markov Models

- Ensemble methods with 5,000+ sample dataset

- Real-time prediction system development

**Phase 3 (12+ months):** Advanced AI Integration

- Deep learning with transformer architectures

- Multi-modal data fusion (text, price, volume)

- Reinforcement learning for dynamic strategies

- Automated feature discovery systems

## 10.3 Academic and Practical Contributions (Enhanced)

**Part C Comprehensive Modeling Contributions:**

- **Systematic validation** of modeling complexity trade-offs in financial prediction

- **Empirical demonstration** of sample size requirements for various algorithms

- **Comprehensive assessment** of market data and technical indicator value

- **Evidence-based framework** for data collection prioritization in finance

- **Cross-validation methodology** establishing robust evaluation standards

**Novel Insights for Financial Modeling:**

- **Complexity Paradox Validation:** Simple models outperform sophisticated ones with limited data

- **Feature Quality Hierarchy:** Market context ¿ technical indicators ¿ complex interactions

- **Sample Size Thresholds:** ¡1,000 observations insufficient for ensemble methods

- **Overfitting Detection Framework:** Negative $R^2$ as early warning system

- **Data Source Value Ranking:** VIX and market regime most important additions

**Risk Management and Business Intelligence:**

- **Model Reliability Assessment:** Current prediction uncertainty quantified ($\pm\$900$ RMSE)

- **Automated Decision-Making Caution:** Models insufficient for unsupervised trading

- **Data-Driven Investment:** Evidence-based prioritization of system enhancements

- **Validation Framework:** Replicable methodology for future model development

## 10.4 Integration with Parts A, B, and D

**Part A Integration:** Enhanced modeling validates initial data exploration findings about limited predictive signals in current dataset structure.

**Part B Integration:** Confirms Part B hypotheses about volatility-position interactions while revealing that additional complexity requires substantially more data.

**Part D Preparation:** Provides quantified model performance baselines and clear data enhancement roadmap for strategic recommendations.

**Unified Business Case:** Establishes evidence-based argument for systematic data infrastructure investment over quick algorithmic solutions.

This comprehensive modeling analysis demonstrates that successful financial prediction requires a balanced approach prioritizing data quality and sample size over algorithmic sophistication, providing a solid foundation for strategic decision-making in Part D.

# Appendix: Technical Implementation Details

## Basic Model Hyperparameters

- **Linear Regression:** Normal equation approximation with correlation-based coefficients

- **Logistic Regression:** Learning rate = 0.01, Iterations = 100, Sigmoid activation

- **Data Split:** 80% training (547 samples), 20% testing (137 samples)

- **Feature Scaling:** Position value normalized by $10,000

## Advanced Model Hyperparameters

- **Polynomial Features:** 11 engineered features including quadratic and interaction terms

- **Ensemble Voting:** 15 models with bootstrap sampling and feature noise injection

- **Cross-Validation:** 5-fold stratified split for robust performance estimation

- **Stock-Specific Models:** Individual regression per stock (8 stocks with sufficient data)

## Model Validation Methodology

- **Cross-Validation:** 5-fold split with stratified sampling to maintain class balance

- **Performance Metrics:** $R^2$, RMSE, MAE calculated on both training and test sets

- **Overfitting Detection:** Training vs. test performance comparison

- **Statistical Significance:** Bootstrap confidence intervals for coefficient estimates

## Feature Engineering Technical Details

Advanced Feature Set:

1. Position Value (normalized)
2. Volatility (%)
3. Transaction Type (binary)
4. Risk Score (composite)
5. Position×Volatility interaction
6. Quadratic position term
7. Quadratic volatility term
8. Sin transformation of volatility
9. Log position value
10. Position size category (binary)
11. Stock code categorical encoding

## Performance Metrics Definitions

- **R²:** $1 - \frac{\text{Residual Sum of Squares}}{\text{Total Sum of Squares}}$

- **RMSE:** $\sqrt{\text{Mean Squared Error}}$

- **MAE:** Mean Absolute Error

- **Cross-Validation Score:** Average R² across k-folds ± standard deviation

- **Baseline Accuracy:** Most frequent class prediction for classification models

## Critical Modeling Lessons Learned

1. **Dataset Size Limitations:** 684 samples insufficient for complex models (¿5,000 recommended)

2. **Feature Quality Priority:** Better features ¿ more complex algorithms

3. **Validation Importance:** Cross-validation essential for reliable performance assessment

4. **Overfitting Prevention:** Simple models more robust with limited financial data

5. **Domain Expertise:** Financial prediction requires specialized feature engineering

# References

Bodie, Z., Kane, A. and Marcus, A. J. (2021), *Investments*, 12th edn, McGraw Hill Education, New York.

Campbell, J. Y., Lo, A. W. and MacKinlay, A. C. (1997), *The Econometrics of Financial Markets*, Princeton University Press, Princeton.

Hastie, T., Tibshirani, R. and Friedman, J. (2009), *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd edn, Springer, New York.

The corresponding GitHub URL for the report is: https://github.com/Melody1604/Assignment-1-Part-C/blob/8f7edba493866c900bb9350faecdc46c8908c60d/Part%20C.ipynb