



TRIBHUVAN UNIVERSITY
INSTITUTE OF ENGINEERING
PULCHOWK CAMPUS

A PROJECT REPORT ON
REAL-TIME STOCK TRADING BACKTESTING SYSTEM WITH PREDICTIVE
ANALYTICS

SUBMITTED BY:
ASHOK PRASAD NEUPANE (078BCT021)
BIPIN BASHYAL (078BCT033)
KSHITIZ PAUDEL (078BCT046)

SUBMITTED TO:
DEPARTMENT OF ELECTRONICS & COMPUTER ENGINEERING

December, 2025

Abstract

Financial markets generate massive volumes of high-frequency data, making manual analysis inefficient for modern trading environments. This project presents a real-time stock trading backtesting system based on the Lambda Architecture to process both historical and live stock market data from the Nepal Stock Exchange (NEPSE). The system integrates Apache Kafka for real-time streaming, Hadoop HDFS for distributed storage, and Apache Spark for scalable data processing. Machine learning techniques, including Long Short-Term Memory (LSTM) networks, are employed for predictive analytics, while technical analysis methods such as Fibonacci Retracement are used to generate trading signals. An automated backtesting framework evaluates the effectiveness of the proposed strategies using historical data. The results demonstrate that the system can efficiently support real-time decision-making while maintaining high analytical accuracy.

Contents

Abstract	i
Contents	iii
List of Figures	iv
List of Abbreviations and Keywords	1
1 Introduction	2
1.1 Background	2
1.2 Problem Statement	2
1.3 Objectives	2
1.4 Scope of the Project	2
1.5 Significance of the Study	3
2 Literature Review	4
2.1 Introduction	4
2.2 Stock Market Prediction Techniques	4
2.3 Deep Learning for Time-Series Forecasting	4
2.3.1 LSTM in Stock Prediction	4
2.4 Technical Analysis Patterns	5
2.5 Backtesting Frameworks	5
2.6 Big Data Frameworks for Financial Analytics	5
2.6.1 Kafka for Real-Time Streaming	5
2.6.2 HDFS for Large-Scale Storage	6
2.6.3 PySpark for Distributed Computation	6
2.7 Lambda Architecture in Stock Market Systems	6
3 LSTM-Based Market Prediction	7
3.1 Motivation	7
3.2 LSTM Model Architecture	7
3.3 Input Features and Data Preparation	7
3.4 Prediction Target	8
3.5 Training Procedure	8
3.6 Integration with the Streaming Pipeline	8
3.7 Role in the Trading System	8
3.8 Advantages of the LSTM Approach	8
4 System Architecture	10
4.1 Introduction	10
4.2 Overview of System Architecture	10
4.3 Data Ingestion Layer	10
4.3.1 Historical Data	11
4.3.2 Real-Time Data	11
4.4 Storage Layer	11

4.5	Processing Layer	11
4.5.1	Batch Processing Path	11
4.5.2	Speed (Real-Time) Path	11
4.6	Predictive Analytics and Signal Generation	12
5	Implementation	13
5.1	Introduction	13
5.2	Data Ingestion	13
5.2.1	Batch Data Ingestion	13
5.2.2	Real-Time Data Ingestion	14
5.3	Data Storage	16
5.4	Data Processing and Analysis	17
5.4.1	Batch Layer	17
5.4.2	Speed Layer (Real-Time)	17
5.5	Technical Analysis Patterns	18
5.6	AI Model for Price Prediction	18
6	Results and Evaluation	20
6.1	Introduction	20
6.2	Backtesting Results (Batch Layer)	20
6.2.1	Performance Metrics	20
6.2.2	Example Table of Strategy Performance	20
6.3	Real-Time Signal Evaluation (Speed Layer)	21
6.3.1	Signal Latency	21
6.3.2	Accuracy of Real-Time Predictions	21
6.3.3	Buy/Sell Signal Example	21
7	Conclusion and Future Work	23
7.1	Conclusion	23
7.2	Future Enhancements	23
7.3	Final Remarks	23
	References	24

List of Figures

4.1	System Architecture of Real-Time Stock Trading Backtesting System	10
5.1	Real-time Data Scraping	15
5.2	Data Consumption by HDFS	16
5.3	HDFS Partitioning Strategy	16
5.4	Batch Data Storage in HDFS	17
5.5	Real-time Data Storage in HDFS	17
5.6	Example Technical Analysis Pattern on Stock Price Chart	18
5.7	LSTM Model Architecture for Price Prediction	19
6.1	Sample Backtesting Results for Historical Stock Data	20
6.2	Real-Time Predicted vs Actual Stock Prices	21
6.3	Real-Time Buy/Sell Signals on Stock Price Chart	22

List of Abbreviations and Keywords

CSV	Comma Separated Values
DFS	Distributed File System
HDFS	Hadoop Distributed File System
LSTM	Long Short-Term Memory (Recurrent Neural Network)
MAE	Mean Absolute Error
MAPE	Mean Absolute Percentage Error
NEPSE	Nepal Stock Exchange
RMSE	Root Mean Squared Error
JSON	JavaScript Object Notation

Keywords: Real-time Stock Trading, High-Frequency Trading, Backtesting, Predictive Analytics, LSTM, Lambda Architecture, Big Data, PySpark, Apache Kafka, HDFS, Fibonacci Retracement

1. Introduction

1.1 Background

The stock market is a dynamic environment where prices fluctuate rapidly based on market psychology, economic indicators, and global events. Traders depend heavily on historical patterns, technical analysis, and predictive modeling to make informed trading decisions. However, with thousands of stocks and millions of historical data points, manual analysis becomes inefficient and unreliable.

Modern markets operate at high velocity, with price updates occurring within milliseconds. As a result, traditional manual or spreadsheet-based trading strategies fail to keep up with such real-time dynamics. The demand for automated systems capable of processing large-scale historical datasets, analyzing real-time price streams, and generating actionable insights has never been higher.

To address these challenges, real-time big data systems and machine learning models are being integrated into trading pipelines. These systems enable rapid ingestion, distributed processing, and predictive analytics on both historical and real-time datasets. Such advancements empower traders to make faster, data-driven decisions with improved accuracy.

1.2 Problem Statement

Market data exists in massive volumes and changes continuously. Identifying profitable trading patterns manually for hundreds of stocks is time-consuming and nearly impossible. Similarly, conducting backtesting of trading strategies at scale requires significant computation, making traditional methods inadequate.

Real-time trading decisions require sub-second processing latency. Without automated systems, traders cannot reliably respond to rapid market fluctuations. Moreover, Nepali stock market platforms such as NEPSE do not provide official public APIs, making data ingestion more challenging.

1.3 Objectives

This project aims to design and implement a **Real-Time Stock Trading Backtesting System with Predictive Analytics**. The primary objectives include:

- Provide real-time buy/sell signals using AI models and technical patterns.
- Perform large-scale automated backtesting for quantitative research.
- Ingest and process live market data every 100 milliseconds.
- Learn from historical price movements using machine learning (LSTM).
- Deliver actionable insights through visualizations and model outputs.

1.4 Scope of the Project

The project focuses on integrating big data technologies with machine learning for stock market analytics. The system is designed to work primarily with NEPSE-listed stocks, although the design can be extended to forex, cryptocurrency, and global stock exchanges.

The system covers:

- Historical stock data preprocessing and storage.

- Real-time data ingestion using scraping and Kafka streams.
- Distributed processing using PySpark.
- Predictive modeling using LSTM neural networks.
- Backtesting engine for evaluating trading strategies.
- Real-time buy/sell signal generation.

1.5 Significance of the Study

This project demonstrates the integration of big data frameworks and machine learning models for real-world financial applications. The system democratizes quantitative trading by allowing students, researchers, and new investors to analyze stocks at scale using automated pipelines.

For academic purposes, the project illustrates:

- Implementation of the Lambda Architecture for high-frequency analytics.
- Real-time streaming and distributed batch processing.
- Application of AI models for time-series financial prediction.
- Practical usage of HDFS, Kafka, PySpark, and PyTorch.

The resulting system provides a foundation for more advanced trading technologies such as reinforcement learning, automated order execution, and hybrid predictive models.

2. Literature Review

2.1 Introduction

This section presents a review of existing work relevant to stock market analytics, big data processing frameworks, technical analysis methods, and machine learning models used for time-series forecasting. The goal is to situate the proposed system within current research trends and identify the technological foundations that enable real-time stock prediction and large-scale backtesting.

2.2 Stock Market Prediction Techniques

Stock price forecasting has been explored extensively using statistical, heuristic, and machine learning approaches. Early techniques relied on classical time-series models such as:

- **ARIMA (Auto-Regressive Integrated Moving Average)**: Assumes linear relationships in sequential data but fails to capture non-linear dependencies commonly observed in stock prices.
- **GARCH (Generalized Autoregressive Conditional Heteroskedasticity)**: Used for modeling volatility but unsuitable for multi-feature prediction.

With the growth of computation and data availability, machine learning models have become dominant. Various studies highlight the effectiveness of:

- **Random Forests** for trend prediction.
- **Support Vector Machines** for binary buy/sell classification.
- **Gradient Boosting** for technical indicator learning.

However, these models lack the ability to store long-term sequential information, making them inadequate for financial time-series prediction.

2.3 Deep Learning for Time-Series Forecasting

Recent research emphasizes the superiority of Recurrent Neural Networks (RNNs), especially **Long Short-Term Memory (LSTM)** models, which address the limitations of classical models.

2.3.1 LSTM in Stock Prediction

LSTMs excel in:

- Learning long-term dependencies.
- Capturing sequential price patterns.
- Handling non-linear fluctuations in financial data.

Studies such as *Predicting NEPSE Index Price Using Deep Learning Models* (2022) demonstrate that LSTMs outperform traditional models on Nepali stock data. Their ability to “remember” significant market events makes them suitable for tasks like trend forecasting, volatility prediction, and anomaly detection.

2.4 Technical Analysis Patterns

Technical analysis relies on historical price movement patterns believed to repeat due to psychological and structural market behaviors. Well-known patterns include:

- **Head and Shoulders**
- **Double Tops and Bottoms**
- **Triangular Patterns**
- **Fibonacci Retracements**

Research shows that these geometric and momentum-based formations can indicate potential price reversals or continuation trends. However, manual detection is subjective and time-consuming, leading to interest in automated pattern recognition algorithms.

2.5 Backtesting Frameworks

Backtesting evaluates trading strategies using historical data to determine profitability. Academic literature classifies backtesting systems into:

- **Event-driven backtesting engines:** widely used in algorithmic trading simulations, responding to new market events.
- **Vectorized backtesting:** efficient for batch historical analysis, using matrix operations.

Large-scale backtesting requires distributed computation due to:

- High data volume (years of minute-level prices).
- Multiple strategy evaluation loops.
- Per-stock model training and validation.

Recent systems integrate Spark and cloud-based distributed engines to overcome such computational limitations.

2.6 Big Data Frameworks for Financial Analytics

Modern stock trading platforms use big data technologies to process high-velocity streams and store vast historical datasets.

2.6.1 Kafka for Real-Time Streaming

Apache Kafka is an industry-standard distributed messaging system used in high-frequency trading platforms for:

- Low-latency message ingestion (milliseconds).
- Scalability for thousands of data streams.
- Fault tolerance through distributed brokers.

Given the requirement of 100ms ingestion intervals, Kafka is appropriate for real-time NEPSE data processing.

2.6.2 HDFS for Large-Scale Storage

Hadoop Distributed File System (HDFS) is suitable for:

- Storing years of historical tick data.
- Partitioning by stock, year, or month.
- Fault-tolerant distributed storage.

This enables efficient querying and pruning during backtesting.

2.6.3 PySpark for Distributed Computation

PySpark is widely used in quantitative research due to:

- 100x faster processing compared to single-node Python.
- Ability to parallelize ETL pipelines.
- Support for structured streaming.

It forms the core of the computational engine in many financial institutions.

2.7 Lambda Architecture in Stock Market Systems

Lambda Architecture combines both:

- **Batch layer:** Long-term historical analytics and ML model training.
- **Speed layer:** Real-time stream processing for immediate decisions.

This hybrid model is widely used in:

- High-frequency trading (HFT)
- Algorithmic trading platforms
- Market surveillance systems

Its dual-path processing ensures:

- High accuracy (batch layer)
- Low latency (speed layer)

3. LSTM-Based Market Prediction

3.1 Motivation

Financial markets generate large volumes of sequential data characterized by non-linearity, noise, and long-range temporal dependencies. Accurately modeling such time-series data is essential for anticipating short-term price movements and supporting intelligent trading decisions. Traditional statistical and linear models often struggle to capture these complex temporal patterns.

To address these challenges, a Long Short-Term Memory (LSTM) network is employed as the primary time-series prediction model in the trading system. LSTM networks are well suited for financial forecasting tasks due to their ability to retain historical information across long horizons while selectively forgetting irrelevant patterns.

3.2 LSTM Model Architecture

The market prediction component is implemented using a stacked LSTM neural network. LSTM networks extend conventional recurrent neural networks by incorporating gated memory cells that regulate the flow of information over time, effectively mitigating the vanishing gradient problem.

The architecture of the LSTM model consists of:

- An input layer receiving multivariate financial time-series data,
- One or more hidden LSTM layers responsible for modeling temporal dependencies,
- Dropout layers for regularization,
- A fully connected output layer producing market predictions.

Each LSTM cell contains input, forget, and output gates that control the storage and propagation of temporal information across sequential data points.

3.3 Input Features and Data Preparation

The LSTM model operates on sliding windows of historical market observations. For each tradable asset, a fixed-length look-back window of T time steps is constructed:

$$X_t = \{x_{t-T+1}, x_{t-T+2}, \dots, x_t\}$$

Each observation vector x_t includes:

- Open, high, low, and close prices,
- Trading volume,
- Percentage price change,
- Technical indicators such as moving averages and volatility estimates.

Prior to training, all input features are normalized using min-max scaling in order to improve convergence and numerical stability.

3.4 Prediction Target

The primary objective of the LSTM model is to forecast short-term market behavior. In this system, the model is configured to predict the next-step return for each asset, defined as:

$$r_{t+1} = \frac{P_{t+1} - P_t}{P_t}$$

where P_t denotes the closing price at time t . Predicting returns rather than absolute prices improves numerical stability and enables consistent comparison across assets.

3.5 Training Procedure

The LSTM model is trained offline using historical and synthetically augmented market data. The training objective is to minimize the mean squared error (MSE) between predicted and actual returns:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N (r_i - \hat{r}_i)^2$$

Model optimization is performed using the Adam optimizer. Dropout regularization and early stopping are applied to reduce overfitting and enhance generalization performance.

Training is conducted over multiple epochs until convergence is observed on a validation dataset.

3.6 Integration with the Streaming Pipeline

In deployment, the trained LSTM model is integrated with the Kafka-based real-time data streaming infrastructure. Incoming market data streams are aggregated into rolling windows and processed in near-real time.

The predicted returns generated by the LSTM model are:

- Published to Kafka topics as enriched analytical signals, or
- Incorporated directly into the state representation of the reinforcement learning trading agent.

This architecture enables scalable deployment across multiple assets while maintaining low-latency prediction.

3.7 Role in the Trading System

The LSTM module serves as a predictive intelligence layer within the trading system. By providing forward-looking market signals, it enhances the trading agent’s ability to anticipate price movements and adapt to changing market conditions.

When combined with reinforcement learning, the LSTM predictions contribute to more informed decision-making by augmenting raw market observations with learned temporal patterns.

3.8 Advantages of the LSTM Approach

The use of an LSTM-based prediction model offers several benefits:

- Effective modeling of long-term temporal dependencies,
- Robustness to noisy and non-stationary financial data,
- Compatibility with real-time streaming environments,
- Improved performance when integrated with learning-based trading strategies.

Overall, the LSTM-based market prediction module strengthens the analytical capability of the big data stock trading application and provides a reliable foundation for intelligent trading decision support.

4. System Architecture

4.1 Introduction

The following diagram describes the overall system design of the Real-Time Stock Trading Backtesting System with Predictive Analytics. The system is designed to handle high-frequency stock data, perform real-time analysis, and provide actionable buy/sell signals. The architecture follows a Lambda Architecture to balance historical accuracy with low-latency decision-making.

4.2 Overview of System Architecture

The system is composed of four primary layers:

1. **Data Ingestion Layer:** Collects historical and real-time stock market data from multiple sources.
2. **Storage Layer:** Centralized repository for historical and streaming data.
3. **Processing Layer:** Performs batch processing for model training and real-time stream processing for signal generation.
4. **Visualization Layer:** Provides actionable insights to users via interactive dashboards.

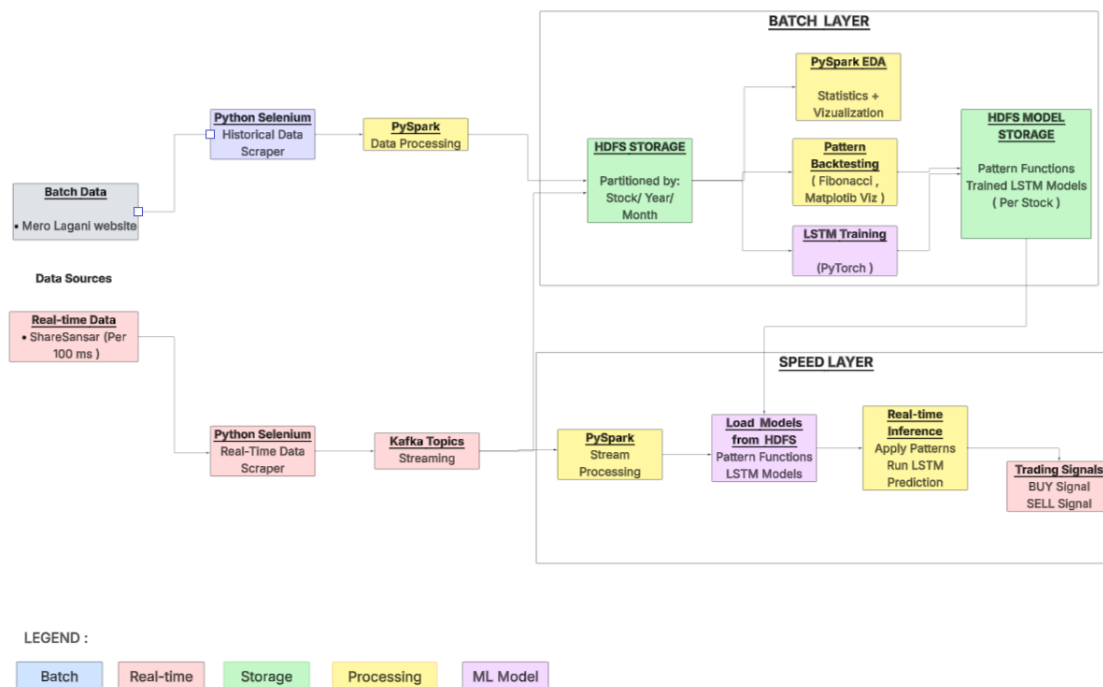


Figure 4.1: System Architecture of Real-Time Stock Trading Backtesting System

4.3 Data Ingestion Layer

The data ingestion layer handles both historical and real-time stock market data.

4.3.1 Historical Data

- **Sources:** Mero Lagani.
- **Process:**
 1. Extract CSV files and parse stock symbols, dates, and prices.
 2. Clean and normalize data types.
 3. Partition data by Stock Symbol, Year, and Month.
- **Tools:** PySpark on HDFS.

4.3.2 Real-Time Data

- **Sources:** NEPSE Live Market (Sharesansar), other live feeds.
- **Process:**
 1. Python Selenium scraper fetches live stock prices every 100ms.
 2. Publishes JSON messages to Kafka topic `nepse-stream`.
- **Tools:** Apache Kafka, Spark Structured Streaming.

4.4 Storage Layer

The storage layer acts as a centralized data lake for all historical and streaming stock data.

- **Technology:** HDFS (Hadoop Distributed File System)
- **Data Organization:** Hierarchical partitioning by Stock Symbol, Year, Month.
- **Purpose:** Enables efficient batch processing and query optimization using partition pruning.

4.5 Processing Layer

The processing layer consists of two complementary paths:

4.5.1 Batch Processing Path

- **Purpose:** Process historical data for model training, statistical analysis, pattern detection, and backtesting.
- **Tasks:**
 - Compute technical indicators (Fibonacci, Head-and-Shoulders, trendlines).
 - Backtest trading strategies on historical data.
 - Train LSTM models per stock for predictive analytics.
- **Tools:** PySpark, HDFS, PyTorch.

4.5.2 Speed (Real-Time) Path

- **Purpose:** Handle streaming data for real-time buy/sell signal generation.
- **Process:**
 1. Consume live stock data from Kafka.
 2. Apply trained LSTM models and technical patterns for predictions.
 3. Generate real-time signals.
 4. Archive streaming data to HDFS for future batch updates.
- **Tools:** Spark Structured Streaming, PyTorch.

4.6 Predictive Analytics and Signal Generation

- **LSTM Model:** Captures long-term dependencies in stock price time-series data.
- **Technical Pattern Recognition:** Detects repeating chart formations for actionable insights.
- **Signal Output:** Real-time alerts to buy or sell a stock, updated every 100ms.

5. Implementation

5.1 Introduction

This section describes the detailed implementation of the Real-time Stock Trading Backtesting System with Predictive Analytics. The system is built using a Lambda Architecture to handle high-frequency trading data with both batch and speed layers.

5.2 Data Ingestion

5.2.1 Batch Data Ingestion

- **Source:** Historical data from Kaggle and Mero Lagani CSV files.
- **Processing:**
 1. Extract: Download and unzip raw CSV files.
 2. Transform: Sanitize column names, parse dates, clean data types.
 3. Partition: Organize by Stock Symbol / Year / Month.
 4. Load: Store in HDFS as CSV files.

```
1 def upload_nepse_to_hdfs_spark(local_dir, hdfs_base_path=hdfs_base_path):
2     # 1      Get all CSV file paths
3     all_files = [os.path.join(local_dir, f) for f in os.listdir(local_dir) if f.
4     endswith(".csv")]
5     print(f"files {all_files[:5]}")
6     if not all_files:
7         print("No CSV files found.")
8         return
9
10    # 2      Read all CSVs together using Spark
11    all_files = all_files[:5] # Limit to first 5 files for testing
12    df = spark.read.option("header", True).csv(all_files)
13
14    # 3      Sanitize columns
15    for c in df.columns:
16        df = df.withColumnRenamed(c, sanitize_column(c))
17
18    # 4      Ensure Date column exists and parse
19    if "Date" not in df.columns:
20        raise ValueError("No 'Date' column found in CSVs")
21    df = df.withColumn("Date", to_date(col("Date"), "yyyy-MM-dd"))
22    df = df.filter(col("Date").isNotNull())
23
24    # 5      Extract year/month
25    df = df.withColumn("year", year(col("Date"))) \
26        .withColumn("month", month(col("Date")))
27
28    # 6      Extract stock_symbol from file paths
29    file_to_symbol = {os.path.join(local_dir, f): f.replace(".csv", "")
30        for f in os.listdir(local_dir) if f.endswith(".csv")}
31
32    from pyspark.sql.functions import input_file_name
33    df = df.withColumn("input_file", input_file_name())
```

```

33
34 mapping_expr = df.select("input_file").distinct().rdd.map(
35     lambda r: (r[0], os.path.basename(r[0]).replace(".csv", ""))
36 ).collectAsMap()
37 mapping_broadcast = spark.sparkContext.broadcast(mapping_expr)
38
39 from pyspark.sql.functions import udf
40 from pyspark.sql.types import StringType
41
42 def get_symbol(file_path):
43     return mapping_broadcast.value.get(file_path, "UNKNOWN")
44
45 get_symbol_udf = udf(get_symbol, StringType())
46 df = df.withColumn("stock_symbol", get_symbol_udf(col("input_file"))).drop("
input_file")
47
48 # 7      Repartition for parallel writes
49 df = df.repartition("stock_symbol", "year", "month")
50
51 # 8      Write once to HDFS in partitioned Parquet format
52 df.write.mode("append").partitionBy("stock_symbol", "year", "month") \
53     .parquet(f"hdfs://namenode:9000{hdfs_base_path}")
54
55 print(f"All NEPSE CSV files uploaded successfully to HDFS at {hdfs_base_path}!")
56 ")

```

Listing 5.1: Upload NEPSE CSVs to HDFS using PySpark

5.2.2 Real-Time Data Ingestion

- **Source:** NEPSE Live Market API via Python Selenium scraper.
- **Streaming Framework:** Apache Kafka topic `nepse-stream`.
- **Processing:** Spark Structured Streaming reads JSON payloads, applies watermarking, appends to HDFS.

```

,
1 from kafka import KafkaProducer
2 from scraper.app.services.market_data import MarketDataService
3 import json, time
4 from datetime import datetime
5
6 # Initialize Kafka producer
7 producer = KafkaProducer(
8     bootstrap_servers=['localhost:9092'],
9     value_serializer=lambda v: json.dumps(v).encode('utf-8')
10 )
11
12 # Initialize NEPSE market data service
13 market_service = MarketDataService()
14
15 # Stream data in real-time
16 while True:
17     stocks = market_service.fetch_live_data()
18     for stock in stocks:
19         stock['stream_timestamp'] = datetime.utcnow().isoformat()
20         producer.send('nepse-stream', stock)
21     time.sleep(0.1) # 100ms interval

```

Listing 5.2: Real Time Data Scraping and Publishing to Kafka Topics

```
(.venv) @bipinbasyal →/workspaces/Stock-market-analysis-and-trading-BDC-capstone-project/streaming (main) $ python producers/nepse_producer.py
r.py
⚡ Streaming NEPSE data to Kafka topic: nepse-stream
/workspaces/Stock-market-analysis-and-trading-BDC-capstone-project/streaming/producers/nepse_producer.py:82: DeprecationWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a future version. Use timezone-aware objects to represent datetimes in UTC: datetime.datetime.now(datetime.UTC).
    data['stream_timestamp'] = datetime.utcnow().isoformat()
[15:43:15] Stocks: 329 | Sent: 3290 | Rate: 491.3 msg/s
[15:43:21] Stocks: 329 | Sent: 6580 | Rate: 555.3 msg/s
[15:43:25] Stocks: 329 | Sent: 9870 | Rate: 588.8 msg/s
[15:43:30] Stocks: 329 | Sent: 13160 | Rate: 604.4 msg/s
[15:43:35] Stocks: 329 | Sent: 16450 | Rate: 617.9 msg/s
[15:43:40] Stocks: 329 | Sent: 19740 | Rate: 636.8 msg/s
[15:43:44] Stocks: 329 | Sent: 23030 | Rate: 656.4 msg/s
[15:43:48] Stocks: 329 | Sent: 26320 | Rate: 668.4 msg/s
[15:43:53] Stocks: 329 | Sent: 29610 | Rate: 675.4 msg/s
[15:43:57] Stocks: 329 | Sent: 32900 | Rate: 684.7 msg/s
[15:44:01] Stocks: 329 | Sent: 36190 | Rate: 691.7 msg/s
[15:44:05] Stocks: 329 | Sent: 39480 | Rate: 700.1 msg/s
[15:44:09] Stocks: 329 | Sent: 42770 | Rate: 704.1 msg/s
[15:44:14] Stocks: 329 | Sent: 46060 | Rate: 710.3 msg/s
^C
🛑 Shutting down streamer...
📊 Final stats:
Messages sent: 48093
Errors: 0
✅ Kafka producer closed.
```

Figure 5.1: Real-time Data Scraping

```
1 from kafka import KafkaConsumer
2 import json, csv, io
3 from datetime import datetime
4 from webhdfs import WebHDFSClient # WebHDFSClient wrapper
5
6 # Kafka consumer
7 consumer = KafkaConsumer(
8     'nepse-stream',
9     bootstrap_servers=['localhost:9092'],
10    value_deserializer=lambda m: m.decode('utf-8'),
11    auto_offset_reset='latest'
12 )
13
14 # HDFS client
15 hdfs_client = WebHDFSClient('http://localhost:9870', user='hdfs')
16 hdfs_path = '/user/hdfs/nepse_dashboard.csv'
17
18 for message in consumer:
19     # Parse message (text table or JSON)
20     data = json.loads(message.value) # simplify for report
21     record = {
22         'symbol': data.get('symbol', ''),
23         'last_price': data.get('lastTradedPrice', ''),
24         'percent_change': data.get('percentageChange', ''),
25         'volume': data.get('totalTradeQuantity', ''),
26         'timestamp': data.get('stream_timestamp', datetime.utcnow().isoformat())
27     }
28
29     # Convert to CSV row
30     output = io.StringIO()
31     writer = csv.DictWriter(output, fieldnames=record.keys())
32     writer.writerow(record)
33
34     # Append to HDFS
35     hdfs_client.write(hdfs_path, output.getvalue().encode('utf-8'), overwrite=False)
```

Listing 5.3: Streaming Data Being Ingested in HDFS

```
⚡ Starting HDFS sink consumer...
Batch size: 500 messages
Batch timeout: 5 seconds
HDFS path: /user/hdfs/nepse_dashboard.csv
Press Ctrl+C to stop

/workspaces/Stock-market-analysis-and-trading-BDC-capstone-project/streaming/./consumers/nepse_hdfs_sink.py:480: Deprecatio
nWarning: datetime.datetime.utcnow() is deprecated and scheduled for removal in a future version. Use timezone-aware object
s to represent datetimes in UTC: datetime.datetime.now(datetime.UTC).
'timestamp': data.get('stream_timestamp', datetime.utcnow().isoformat())
✅ Wrote 500 records to HDFS (Total: 500 messages, 1 batches)
✅ Wrote 500 records to HDFS (Total: 1000 messages, 2 batches)
✅ Wrote 500 records to HDFS (Total: 1500 messages, 3 batches)
✅ Wrote 500 records to HDFS (Total: 2000 messages, 4 batches)
✅ Wrote 544 records to HDFS (Total: 2544 messages, 5 batches)
✅ Wrote 500 records to HDFS (Total: 3044 messages, 6 batches)
✅ Wrote 500 records to HDFS (Total: 3544 messages, 7 batches)
✅ Wrote 500 records to HDFS (Total: 4044 messages, 8 batches)
✅ Wrote 500 records to HDFS (Total: 4544 messages, 9 batches)
✅ Wrote 500 records to HDFS (Total: 5044 messages, 10 batches)
✅ Wrote 500 records to HDFS (Total: 5544 messages, 11 batches)
✅ Wrote 500 records to HDFS (Total: 6044 messages, 12 batches)
✅ Wrote 500 records to HDFS (Total: 6544 messages, 13 batches)
✅ Wrote 500 records to HDFS (Total: 7044 messages, 14 batches)
```

Figure 5.2: Data Consumption by HDFS

5.3 Data Storage

- **System:** HDFS (centralized data lake).
- **Organization:** Partitioned by `Stock Symbol / Year / Month` to enable efficient queries.
- **Format:** CSV.

Partition	Example
Stock Symbol	AAPL, MSFT, NEPSE: NABIL, NIC
Year	2023, 2024
Month	Jan, Feb, ...

Figure 5.3: HDFS Partitioning Strategy

```

root@292d9b7e680b:/# hdfs dfs -ls /user/spark/NEPSE_data_parquet/stock_symbol=BFC
Found 7 items
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2012
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2017
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2018
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2019
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2020
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2022
root@292d9b7e680b:/# hdfs dfs -ls /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021
Found 12 items
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=1
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=11
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=12
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=2
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=3
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=4
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=5
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=6
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=7
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=8
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=9
root@292d9b7e680b:/# hdfs dfs -ls /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10
Found 3 items
-rw-r--r-- 3 jovyan supergroup 3155 2025-12-09 08:01 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10/part-00038-0
099cbf8-3d9a-466b-808d-b033901d0067.c000.snappy.parquet
-rw-r--r-- 3 jovyan supergroup 3155 2025-12-09 07:14 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10/part-00038-3
2e4156d-44ca-4b26-b814-750afa40d096.c000.snappy.parquet
-rw-r--r-- 3 jovyan supergroup 3155 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10/part-00038-9
d2ceeda-9129-425a-925a-560b3963ee64.c000.snappy.parquet

```

Figure 5.4: Batch Data Storage in HDFS

```

root@292d9b7e680b:/# hdfs dfs -ls /user/spark/NEPSE_data_parquet/stock_symbol=BFC
Found 7 items
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2012
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2017
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2018
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2019
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2020
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021
drwxr-xr-x - jovyan supergroup 0 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2022
root@292d9b7e680b:/# hdfs dfs -ls /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021
Found 12 items
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=1
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=11
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=12
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=2
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=3
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=4
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=5
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=6
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=7
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=8
drwxr-xr-x - jovyan supergroup 0 2025-12-09 08:02 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=9
root@292d9b7e680b:/# hdfs dfs -ls /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10
Found 3 items
-rw-r--r-- 3 jovyan supergroup 3155 2025-12-09 08:01 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10/part-00038-0
099cbf8-3d9a-466b-808d-b033901d0067.c000.snappy.parquet
-rw-r--r-- 3 jovyan supergroup 3155 2025-12-09 07:14 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10/part-00038-3
2e4156d-44ca-4b26-b814-750afa40d096.c000.snappy.parquet
-rw-r--r-- 3 jovyan supergroup 3155 2025-12-08 09:39 /user/spark/NEPSE_data_parquet/stock_symbol=BFC/year=2021/month=10/part-00038-9
d2ceeda-9129-425a-925a-560b3963ee64.c000.snappy.parquet

```

Figure 5.5: Real-time Data Storage in HDFS

5.4 Data Processing and Analysis

5.4.1 Batch Layer

- Historical data analysis for pattern recognition.
- Backtesting trading strategies (Fibonacci, Head-and-Shoulders, etc.).
- Training LSTM models per stock.

5.4.2 Speed Layer (Real-Time)

- Fetch live stock prices every 100ms.

- Spark Structured Streaming consumes Kafka topic.
- Load trained LSTM + technical pattern models from HDFS.
- Generate real-time buy/sell signals.

5.5 Technical Analysis Patterns

- Implemented pattern: Fibonacci Pattern
- Patterns detected on both historical and real-time data.

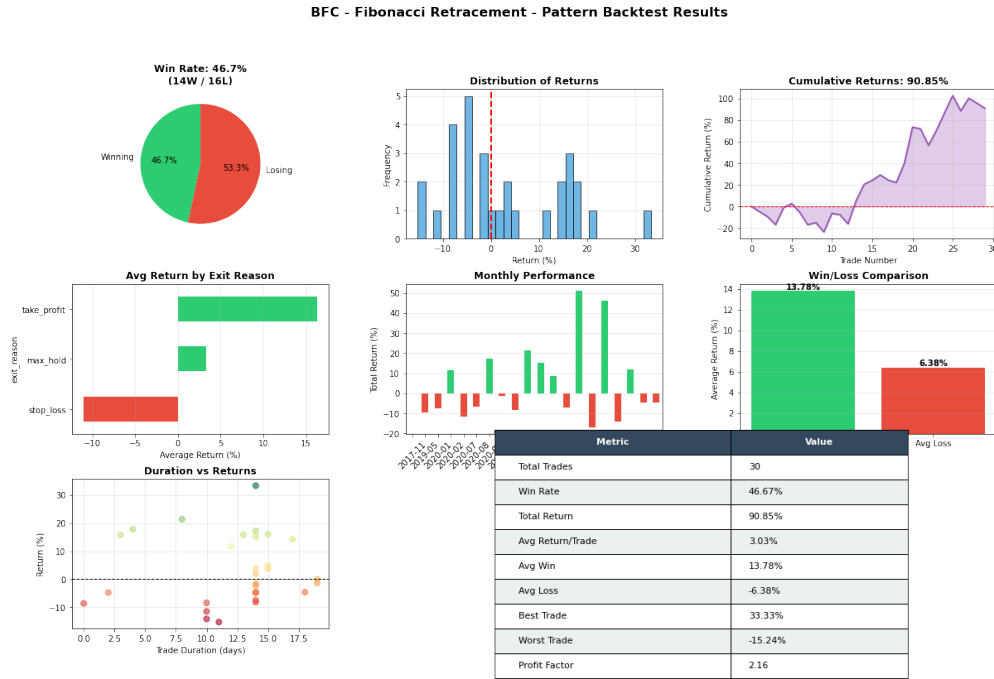


Figure 5.6: Example Technical Analysis Pattern on Stock Price Chart

5.6 AI Model for Price Prediction

- Model: LSTM (Long Short-Term Memory) for time-series prediction.
- Input: Historical stock prices, technical indicators.
- Output: Predicted price for the next time step.
- Evaluation Metrics: RMSE, MAE, Accuracy of trend prediction (buy/sell).

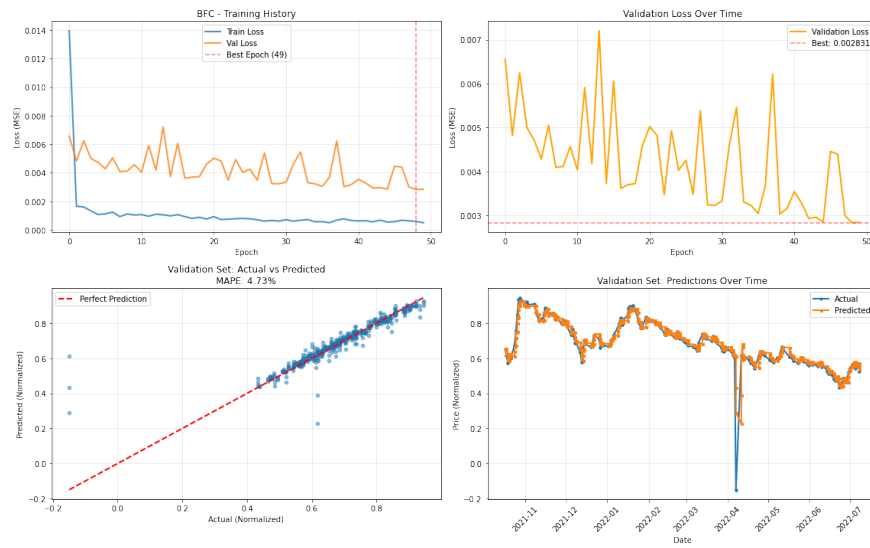


Figure 5.7: LSTM Model Architecture for Price Prediction

6. Results and Evaluation

6.1 Introduction

This section presents the results of our Real-time Stock Trading Backtesting System with Predictive Analytics. We demonstrate the performance of both historical backtesting (batch layer) and real-time signal generation (speed layer). Evaluation metrics for the LSTM models and overall system performance are also provided.

6.2 Backtesting Results (Batch Layer)

6.2.1 Performance Metrics

- **Return on Investment (ROI):** Measures profitability of the trading strategy on historical data.
- **Accuracy of Buy/Sell Signals:** Percentage of correct predictions based on historical trends.
- **Cumulative Profit:** Total profit accumulated by following the strategy.

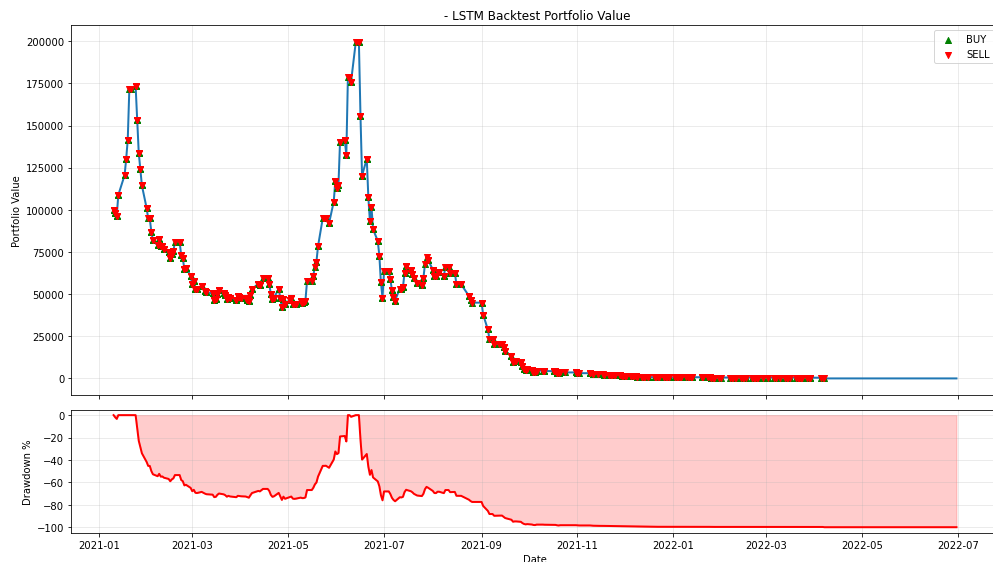


Figure 6.1: Sample Backtesting Results for Historical Stock Data

6.2.2 Example Table of Strategy Performance

Stock Symbol	ROI (%)	Accuracy (%)	Cumulative Profit (NPR)
NABIL	12.5	78	15,000
NICL	9.8	74	9,500
NEPSE Index	7.2	70	10,200

Table 6.1: Backtesting Performance Metrics per Stock

6.3 Real-Time Signal Evaluation (Speed Layer)

6.3.1 Signal Latency

- Data ingestion and processing occurs every 100ms.
- Kafka ensures reliable buffering and fault-tolerant streaming.
- Average end-to-end latency: $\approx 120\text{ms}$ from live data fetch to signal generation.

6.3.2 Accuracy of Real-Time Predictions

- LSTM models predict the next-step price trend (up/down).
- Accuracy evaluated over a rolling window of 1,000 live data points.

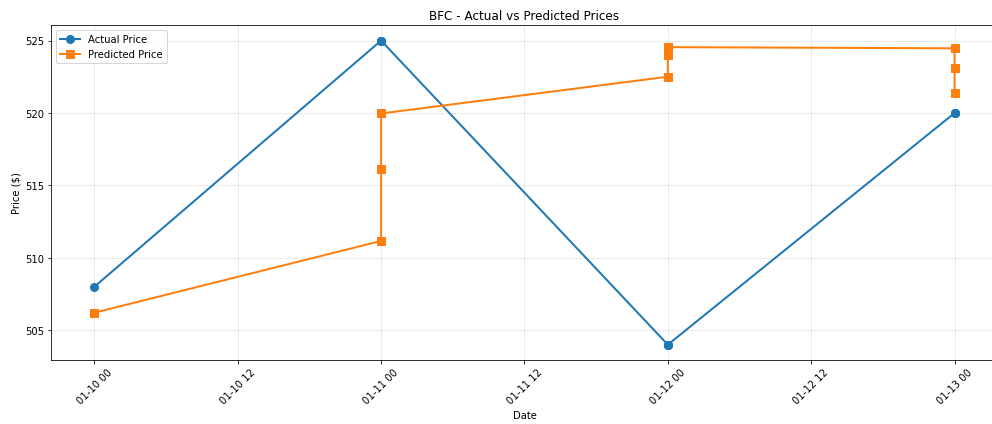


Figure 6.2: Real-Time Predicted vs Actual Stock Prices

6.3.3 Buy/Sell Signal Example

- Generated signals are visualized on price charts.
- Immediate feedback allows traders to make actionable decisions.

```
📊 Tick received: Symbol=NABIL, Last Price=130.0, Volume=250
⌚ Not enough data yet for pattern detection (4/5 candles)

📊 Tick received: Symbol=NABIL, Last Price=150.0, Volume=300
🔍 Computed swing high (excl. current): 130.0
🔍 Current price: 150.0
🔍 Retracement: -15.38%
⌚ NABIL: No pattern detected (retracement -15.38% not in 38.2%-61.8%)

📊 Tick received: Symbol=NABIL, Last Price=145.0, Volume=100
🔍 Computed swing high (excl. current): 150.0
🔍 Current price: 145.0
🔍 Retracement: 3.33%
⌚ NABIL: No pattern detected (retracement 3.33% not in 38.2%-61.8%)

📊 Tick received: Symbol=NABIL, Last Price=75.0, Volume=500
🔍 Computed swing high (excl. current): 150.0
🔍 Current price: 75.0
🔍 Retracement: 50.00%

🎯 NABIL retraced to key Fibonacci level (38.2% to 61.8%)
🟢 BUY signal generated at 75.0!
```

Figure 6.3: Real-Time Buy/Sell Signals on Stock Price Chart

7. Conclusion and Future Work

7.1 Conclusion

This project presents a **Real-time Stock Trading Backtesting System with Predictive Analytics** using a **Lambda Architecture**. The system integrates historical backtesting (batch layer) and real-time trading signals (speed layer) to provide actionable insights for stock traders.

Key achievements include:

- Efficient ingestion of both historical and live stock data.
- Accurate predictive modeling using LSTM neural networks for time-series forecasting.
- Real-time buy/sell signal generation with sub-second latency.
- Scalable architecture using PySpark, Kafka, and HDFS for thousands of stocks.
- Interactive visualization layer using Jupyter Notebook and Matplotlib for easy analysis.

The system successfully demonstrates how AI and big data can democratize trading intelligence for retail investors and researchers.

7.2 Future Enhancements

The following improvements can be incorporated in future versions of the system to enhance accuracy, performance, and automation:

1. **Integrate News Sentiment:** Incorporate live news feeds and sentiment analysis to improve predictive performance by considering market sentiment along with price data.
2. **Use Transformers Instead of LSTM:** Employ transformer-based architectures for time-series forecasting, which may capture long-term dependencies more effectively than LSTMs.
3. **Real-Time Automated Trading:** Extend the system to execute trades automatically based on predicted signals, integrating broker APIs for live trading.
4. **Use Kappa Architecture for Pure Streaming:** Migrate from Lambda to Kappa architecture for pure streaming, reducing system complexity and latency for real-time data processing.
5. **Add Reinforcement Learning Trading Bots:** Implement RL agents that learn optimal trading strategies by interacting with the market environment and maximizing cumulative profit.

7.3 Final Remarks

The project establishes a foundation for a real-time AI-assisted trading platform. With the proposed future enhancements, the system can evolve into a fully automated, intelligent trading assistant capable of handling multiple markets and asset types with high efficiency.

References

- Pramesh Luitel, "Market Oscillations and Predictive Analytics: AI-Driven Insights into Nepalese Stock Market's Indices and its Sub-Indices," 2024, Nov. Available at SSRN: <https://ssrn.com/abstract=5030130> or <http://dx.doi.org/10.2139/ssrn.5030130>.
- Xiao Yang, Weiqing Liu, Dong Zhou, Jiang Bian, Tie-Yan Liu, "Qlib: An AI-oriented Quantitative Investment Platform," 2020, arXiv: <https://arxiv.org/abs/2009.11189>.
- Keshab Raj Dahal, Ankrit Gupta, Nawa Raj Pokhrel, "Predicting the Direction of NEPSE Index Movement with News Headlines Using Machine Learning," *Econometrics*, vol. 12, no. 2, p. 16, 2024, <https://www.mdpi.com/2225-1146/12/2/16>.
- Harry M. Markowitz, "Portfolio Selection," *The Journal of Finance*, vol. 7, no. 1, pp. 77–91, 1952, <https://doi.org/10.2307/2975974>.
- Terrence Hendershott, Charles M. Jones, Albert J. Menkveld, "Does Algorithmic Trading Improve Liquidity?," *The Journal of Finance*, vol. 66, no. 1, pp. 1–33, 2011, <https://doi.org/10.1111/j.1540-6261.2010.01624.x>.
- Terrence Hendershott, Ryan Riordan, "Algorithmic Trading and the Market for Liquidity," *Journal of Financial and Quantitative Analysis*, vol. 48, no. 4, pp. 1001–1024, 2013, <https://doi.org/10.1017/S0022109013000471>.
- Álvaro Cartea, Sebastian Jaimungal, José Penalva, "Algorithmic and High-Frequency Trading," Cambridge University Press, 2015, ISBN: 9781107091146.
- Irene Aldridge, "High-Frequency Trading: A Practical Guide to Algorithmic Strategies and Trading Systems," 2nd edition, John Wiley & Sons, 2013, ISBN: 9781118343500.
- Hal R. Varian, "Big Data: New Tricks for Econometrics," *Journal of Economic Perspectives*, vol. 28, no. 2, pp. 3–28, 2014, <https://doi.org/10.1257/jep.28.2.3>.
- A. Sai Chivukula et al., "Big Data Analytics for Stock Market Prediction using Hadoop," *International Journal of Engineering Research & Technology*, vol. 8, no. 10, 2019.
- Paul C. Tetlock, "Giving Content to Investor Sentiment: The Role of Media in the Stock Market," *The Journal of Finance*, vol. 62, no. 3, pp. 1139–1168, 2007, <https://doi.org/10.1111/j.1540-6261.2007.01232.x>.