

Stock Price Prediction Analytics using Snowflake & Airflow

Moksh Aggarwal, Shivam Mahendru
Department of Applied Data Science
San Jose State University

Email: moksh.aggarwal@sjsu.edu, shivam.mahendru@sjsu.edu

Abstract: This project focuses on forecasting stock prices using historical data and technical indicators, with an emphasis on building a scalable and automated data pipeline. The primary goal is to predict short-term and long-term stock price movements through machine learning models. In addition to stock price forecasting, the project involves the creation of a robust pipeline for automating data extraction, transformation, and loading (ETL) using Apache Airflow for workflow orchestration and Docker for containerization. This approach ensures efficient model deployment and reproducibility, providing a streamlined framework for integrating predictive analytics in real-world financial applications.

Index Terms

Stock Price Prediction, Machine Learning, Apache Airflow, Snowflake, Docker, ETL Pipeline

I. INTRODUCTION

Stock price prediction is a key challenge in the field of financial analysis, with significant implications for investors and traders. This project, titled "Building a Stock Price Prediction Analytics using Snowflake & Airflow," is a group effort conducted as part of the Data226 course at San Jose State University. It leverages the yfinance API to access comprehensive stock market data, including open, high, low, close, and volume information. The dataset enables various financial analyses, such as stock price prediction, trend analysis, technical indicator analysis, portfolio optimization, volatility analysis, volume analysis, and sentiment analysis. For this lab, we focus on implementing a stock price prediction system to forecast the stock prices of selected companies for the next 7 days using the last 180 days of historical data.

II. PROBLEM STATEMENT

The stock price prediction process involves multiple complex steps, from data collection and preprocessing to model deployment and analysis. These processes can be time-consuming, error-prone, and difficult to scale when performed manually. The challenge is to create a **scalable and maintainable automated pipeline** that can efficiently handle the end-to-end workflow of stock price prediction. A database and data pipelines are essential to store and manage large volumes of historical stock data, automate data updates, and integrate machine learning models for forecasting. This project addresses the following key problems:

- **Automation of ETL Processes:** Manual extraction, transformation, and loading (ETL) of stock data leads to inefficiencies. Automation is needed to streamline the workflow.
- **Scalability:** The pipeline must handle large volumes of stock data and continuously update forecasts without performance degradation.
- **Reproducibility and Consistency:** The process, from data collection to model deployment, must be reproducible and consistent.
- **Efficient Workflow Orchestration:** Apache Airflow is required to orchestrate tasks, enabling scheduling, monitoring, and dependency management.
- **Containerization for Deployment:** Docker ensures portability, isolation, and scalability of the pipeline across different environments.

By addressing these issues, this project provides an efficient, reliable, and automated framework for predicting stock prices, reducing the time and effort required for financial analysis.

III. SOLUTION REQUIREMENTS

The design of the stock price prediction system hinges on a set of essential components and tools, each tailored to meet the project's objectives. The following outlines the key requirements for building this automated analytics framework:

- **Data Storage with Snowflake:** The system employs Snowflake as its central repository to house historical stock data sourced from external APIs and to manage machine learning models. This database supports data persistence, model training, and retrieval of predictive outputs.

- **Workflow Automation via Airflow:** Apache Airflow orchestrates the entire pipeline by scheduling and managing tasks such as data ingestion, transformation, and forecasting. It ensures that stock data is refreshed in Snowflake on a daily basis through automated workflows.
- **Stock Data Acquisition:** The system integrates popular Python library that is YFinance API, to fetch critical market data, including stock prices, trading volumes, and related metrics, forming the foundation for analysis.
- **Python Ecosystem:** Several Python libraries power the system's functionality:
 - `yfinance` api: Python module to fetch financial and market data.
 - `airflow.providers.snowflake`: Enables seamless connectivity with Snowflake.
 - `airflow.operators.trigger_dagrun`: `TriggerDagRunOperator` in Apache Airflow is used to trigger the execution of another DAG in sequence after completion of current DAG.
- **Supporting Dependencies:** The implementation relies on external packages, including:
 - `apache-airflow`: Core framework for task orchestration.
 - `snowflake-connector-python`: Provides Python-Snowflake integration.
 - `apache-airflow-providers-snowflake`: Enhances Airflow's Snowflake compatibility.

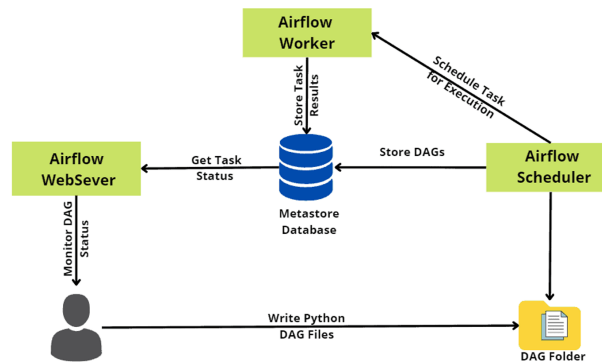


Fig. 1. Airflow Workflow Diagram

A. System Constraints

The system, while robust, operates within certain boundaries:

- **Unpredictable Market Dynamics:** External factors like geopolitical events or shifts in investor behavior can affect stock prices in ways that models based on historical data may not anticipate, potentially reducing prediction accuracy.
- **Data Retrieval Delays:** API constraints or Airflow scheduling intervals may introduce latency, preventing real-time data updates.
- **Model Assumptions:** The predictive models assume continuity in historical patterns, an assumption that may falter during extreme market volatility or novel conditions.

IV. FUNCTIONAL ANALYSIS

A. Data Pipeline Components

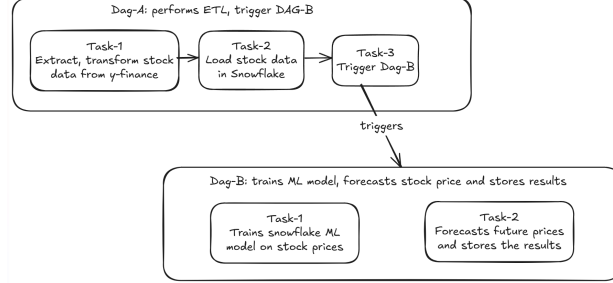
- **Data Ingestion:** Extract stock data using the `yfinance` API and store raw data in **Snowflake**.
- **Data Processing & Transformation:** Clean and format stock data, aggregating necessary fields for analysis.
- **Machine Learning Forecasting:** Train an ML model in Snowflake on historical stock prices to predict prices for the next **7 days or more**.
- **Data Storage & Querying:** Store raw and predicted data in **Snowflake tables** and use SQL queries for insights.
- **Automation with Airflow:** An **ETL pipeline** runs daily to update stock data, and a **forecasting pipeline** updates predictions using recent data, both managed via Airflow DAGs.

V. IMPLEMENTATION

Field	Type	Description
Stock Symbol	VARCHAR	Stock ticker symbol (AAPL and GOOG)
Date	DATE	Date of the stock data
Open	FLOAT	Opening price
Close	FLOAT	Closing price
Min	FLOAT	Minimum price
Max	FLOAT	Maximum price
Volume	INT	Volume of stocks traded

TABLE I
TABLE STRUCTURE IN SNOWFLAKE

Fig. 2. System Design of the DAGs and Tasks



VI. PYTHON CODE

The Python code defines the Airflow DAGs and tasks for ETL and forecasting, integrating with the python's y-finance module as a data source. Below are snippets of the DAG definitions:

Listing 1. DAG-I: Responsible for extracting stocks data from y-finance, transforming, loading it in Snowflake. This DAG further triggers a DAG for training ML model and making predictions on the loaded data.

```

1 File: https://github.com/mahendruShivam29/Airflow/blob/main/dags/extract_transform_load_stock_data.py
2
3 import yfinance as yf
4 from airflow import DAG
5 from datetime import datetime
6 from airflow.decorators import task
7 from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
8 from airflow.operators.trigger_dagrun import TriggerDagRunOperator
9
10 def return_snowflake_conn():
11     hook = SnowflakeHook(snowflake_conn_id='snowflake_conn')
12     conn = hook.get_conn()
13     return conn.cursor()
14
15 @task
16 def extract_and_transform_last_180d_data_from_yfinance(cursor, symbols):
17     """
18     - Extract the last 180 days of Stock price data from Y-finance.
19     - Perform Transformation, here itself. Timestamp Object to Date. Round decimal to 4 places.
20     """
21     result = []
22     for symb in symbols:
23         stock = yf.Ticker(symb)
24         df = stock.history(period='180d', interval='1d')
25
26         for row in df.itertuples(index=True, name="Row"):
27             record = {}
28             # Transformation.
29             record["date"] = row.Index.strftime('%Y-%m-%d')
30             record["symbol"] = symb
31             record["open"] = "{:.4f}".format(row.Open)
32             record["high"] = "{:.4f}".format(row.High)
33             record["low"] = "{:.4f}".format(row.Low)
34             record["close"] = "{:.4f}".format(row.Close)
35             record["volume"] = str(row.Volume)
36             result.append(record)
37
38     for r in result:
  
```

```

39         print(r)
40     return result
41
42 @task
43 def load_stock_data(cursor, stock_data, stock_data_table):
44     """
45     - Loads Stock Data in Snowflake using full refresh, try-except, commit and rollback.
46     """
47     try:
48         cursor.execute("BEGIN;")
49         create_stock_data_table_sql = f"""CREATE TABLE IF NOT EXISTS {stock_data_table} (
50             symbol VARCHAR(10),
51             date DATE,
52             open FLOAT,
53             high FLOAT,
54             low FLOAT,
55             close FLOAT,
56             volume FLOAT,
57             PRIMARY KEY (symbol, date)
58         );"""
59         cursor.execute(create_stock_data_table_sql)
60         cursor.execute(f"DELETE FROM {stock_data_table}")
61
62         for record in stock_data:
63             symbol = record["symbol"]
64             date = record["date"]
65             open = record["open"]
66             high = record["high"]
67             low = record["low"]
68             close = record["close"]
69             volume = record["volume"]
70
71             insert_record_statement = f"INSERT INTO {stock_data_table} (symbol, date, open, high, low,
72                 close, volume) VALUES (%s, %s, %s, %s, %s, %s, %s)"
73             cursor.execute(insert_record_statement, (symbol, date, open, high, low, close, volume))
74             print(f"INSERT INTO {stock_data_table} ({symbol}, {date}, {open}, {high}, {low}, {close},
75                 {volume})")
76
77             cursor.execute("COMMIT;")
78         except Exception as e:
79             cursor.execute("ROLLBACK;")
80             print(e)
81             raise e
82
83 with DAG(
84     dag_id = 'extract_transform_load_stock_data',
85     start_date = datetime(2025, 3, 3),
86     schedule_interval = '@daily',
87     catchup = False,
88     tags = ['ETL']
89 ) as dag:
90     cursor = return_snowflake_conn()
91     symbols = ['AAPL', 'GOOG']
92     stock_data_table = "DEV.LAB.STOCK_DATA"
93     stock_data = extract_and_transform_last_180d_data_from_yfinance(cursor, symbols)
94     load_data = load_stock_data(cursor, stock_data, stock_data_table)
95     trigger_train_predict = TriggerDagRunOperator(
96         task_id = 'trigger_train_predict',
97         trigger_dag_id = 'train_predict_stock_data',
98         execution_date = '{{ ds }}',
99         reset_dag_run = True
100     )
101
102 stock_data >> load_data >> trigger_train_predict

```

Listing 2. DAG-2: Responsible for training model in Snowflake, making stock price prediction for next 7 days and merging that with historical database.

```

1 File: https://github.com/mahendruShivam29/Airflow/blob/main/dags/train\_predict\_stock\_data.py
2 from airflow import DAG
3 from airflow.decorators import task
4 from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
5 from datetime import datetime
6
7
8 def return_snowflake_conn():
9     hook = SnowflakeHook(snowflake_conn_id='snowflake_conn')
10    conn = hook.get_conn()

```

```

11     return conn.cursor()
12
13 @task
14 def train(cursor, train_input_table, train_view, forecast_function_name):
15     """
16     - Create a view with training related columns.
17     - Create a model with the view above.
18     """
19
20     create_view_sql = f"""CREATE OR REPLACE VIEW {train_view} AS SELECT
21         DATE, CLOSE, SYMBOL
22     FROM {train_input_table};"""
23
24     create_model_sql = f"""CREATE OR REPLACE SNOWFLAKE.ML.FORECAST {forecast_function_name} (
25         INPUT_DATA => SYSTEM$REFERENCE('VIEW', '{train_view}'),
26         SERIES_COLNAME => 'SYMBOL',
27         TIMESTAMP_COLNAME => 'DATE',
28         TARGET_COLNAME => 'CLOSE',
29         CONFIG_OBJECT => {{ 'ON_ERROR': 'SKIP' }}
30     );"""
31
32     try:
33         cursor.execute(create_view_sql)
34         cursor.execute(create_model_sql)
35         cursor.execute(f"CALL {forecast_function_name}!SHOW_EVALUATION_METRICS();")
36     except Exception as e:
37         print(e)
38         raise e
39
40 @task
41 def predict(cursor, forecast_function_name, train_input_table, forecast_table, final_table):
42     """
43     - Generate predictions and store the results to a table named forecast table.
44     - Union the predictions with historical data, then make a final table.
45     """
46     make_prediction_sql = f"""BEGIN
47         CALL {forecast_function_name}!FORECAST(
48             FORECASTING_PERIODS => 7,
49             CONFIG_OBJECT => {{'prediction_interval': 0.95}}
50         );
51         LET x := SQLID;
52         CREATE OR REPLACE TABLE {forecast_table} AS SELECT * FROM TABLE(RESULT_SCAN(:x));
53     END;"""
54
55     create_final_table_sql = f"""CREATE OR REPLACE TABLE {final_table} AS
56         SELECT SYMBOL, DATE, CLOSE AS actual, NULL as forecast, NULL as lower_bound, NULL as upper_bound
57     FROM {train_input_table}
58     UNION ALL
59     SELECT replace(series, '', '') as SYMBOL, ts as DATE, NULL as actual, forecast, lower_bound,
60         upper_bound
61     FROM {forecast_table};"""
62
63     try:
64         cursor.execute(make_prediction_sql)
65         cursor.execute(create_final_table_sql)
66     except Exception as e:
67         print(e)
68         raise e
69
70 with DAG(
71     dag_id = 'train_predict_stock_data',
72     start_date = datetime(2025, 3, 5),
73     catchup = False,
74     tags = ['ML', 'ELT']
75 ) as dag:
76     cursor = return_snowflake_conn()
77
78     train_input_table = "DEV.LAB.STOCK_DATA"
79     train_view = "DEV.LAB.STOCK_DATA_VIEW"
80     forecast_table = "DEV.LAB.STOCK_DATA_FORECAST"
81     forecast_function_name = "DEV.LAB.PREDICT_"
82     final_table = "DEV.LAB.MARKET_DATA"
83
84     train(cursor, train_input_table, train_view, forecast_function_name) >> predict(cursor,
85         forecast_function_name, train_input_table, forecast_table, final_table)

```

VII. SCREENSHOTS

Fig. 3. Airflow UI with the DAGs

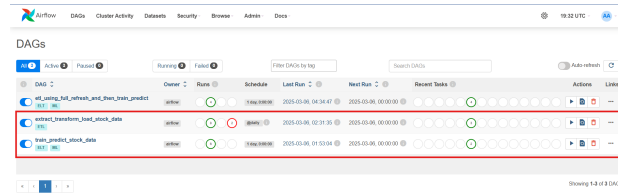


Fig. 4. DAG-1 Details, runs

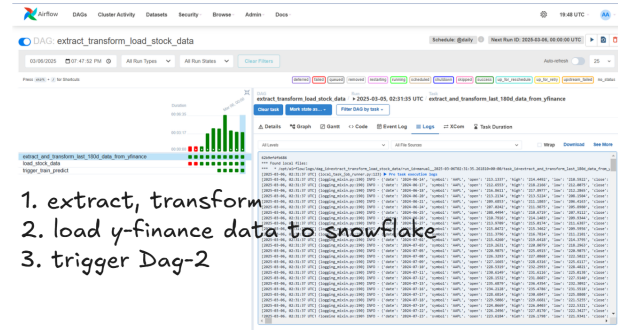


Fig. 5. DAG-1 Task dependency

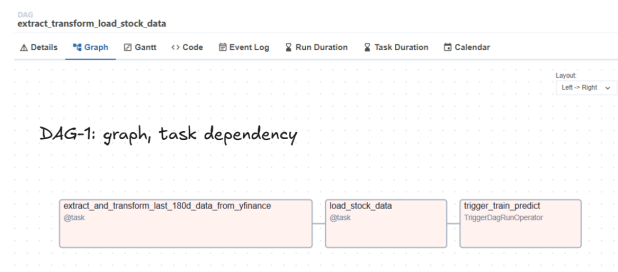


Fig. 6. DAG-2 Details, runs



Fig. 7. DAG-2 Task dependency

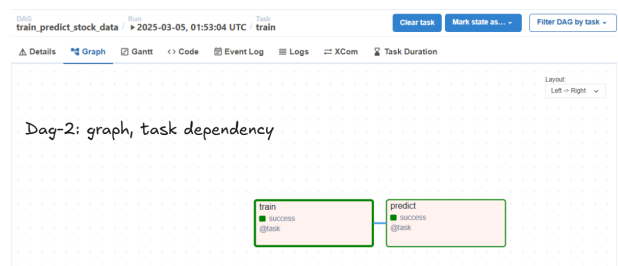


Fig. 8. DAGs Dependency

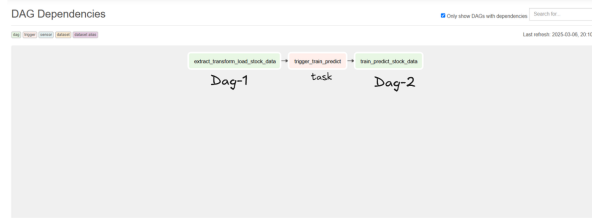


Fig. 9. Table Schema and Structure

```

1 create or replace TABLE DEV.LAB_STOCK_DATA (
2     SYMBOL VARCHAR(10) NOT NULL,
3     DATE DATE NOT NULL,
4     OPEN FLOAT,
5     HIGH FLOAT,
6     LOW FLOAT,
7     CLOSE FLOAT,
8     VOLUME FLOAT,
9     primary key (SYMBOL, DATE)
10 );

1 create or replace TABLE DEV.LAB_MARKET_DATA (
2     SYMBOL VARCHAR(16777216),
3     DATE TIMESTAMP_NTZ(9),
4     ACTUAL FLOAT,
5     FORECAST FLOAT,
6     LOWER_BOUND FLOAT,
7     UPPER_BOUND FLOAT
8 );

1 create or replace TABLE DEV.LAB_STOCK_DATA_FORECAST (
2     SERIES VARIANT,
3     TS TIMESTAMP_NTZ(9),
4     FORECAST FLOAT,
5     LOWER_BOUND FLOAT,
6     UPPER_BOUND FLOAT
7 );

1 create or replace view DEV.LAB_STOCK_DATA_VIEW(
2     DATE,
3     CLOSE,
4     SYMBOL
5 ) as SELECT
6     DATE, CLOSE, SYMBOL
7 FROM DEV.LAB_STOCK_DATA;
```

Fig. 10. Table showing combined results

	SYMBOL	DATE	ACTUAL	FORECAST	LOWER_BOUND	UPPER_BOUND
1	AAPL	2025-03-14 00:00:00.000	null	235.322277348	210.009299471	252.1041837
2	AAPL	2025-03-13 00:00:00.000	null	235.373258934	210.287126649	250.9279190
3	AAPL	2025-03-12 00:00:00.000	null	230.48805614	202.88843243	250.7832088
4	AAPL	2025-03-11 00:00:00.000	null	235.78997961	222.301150113	249.2523270
5	AAPL	2025-03-10 00:00:00.000	null	235.400893799	222.990798435	247.8919746
6	AAPL	2025-03-07 00:00:00.000	null	235.81340448	224.699349743	246.4121527
7	AAPL	2025-03-06 00:00:00.000	null	235.630585932	226.136093389	244.7319889
8	AAPL	2025-03-05 00:00:00.000	235.74	null	null	null
9	AAPL	2025-03-04 00:00:00.000	235.93	null	null	null
10	AAPL	2025-03-03 00:00:00.000	238.03	null	null	null
11	AAPL	2025-02-28 00:00:00.000	241.84	null	null	null

Fig. 11. Airflow connections and variables

Search...

Actions

+

Record Count: 1

	Connection ID	Connection Type	Description	Host	Port	Is Encrypted	Is Extra Encrypted
<input type="checkbox"/>	airflow	airflow					

Edit Connection

Connection ID *

airflow_conn

Connection Type *

airflow

Description

Schema

None

Parameters

airflow.parameters

Extra

Account

QUANTROF LIMITED

Warehouse

COMPUTE_3M

Database

DEV

Region

us-east-1

Role

ACCOUNTADMIN

Key *

varigade_api_key

Val

Description

Varigade api key for getting stock data.

Save ID

+

Link to the codes.

DAG-1: ETL, trigger DAG-2.

DAG-2: Train ML model and forecast future prices.

VIII. CONCLUSION

This project successfully developed an automated data pipeline for stock price prediction using Apache Airflow and Snowflake. By leveraging historical stock data and training machine learning model, we built a scalable solution to forecast stock prices for the next 7 days which could be extended based on the requirement. Docker ensured reproducibility and scalability, while Airflow provided robust workflow orchestration. The system supports financial decision-making by delivering actionable insights.

REFERENCES

- [1] yfinance API, <https://pypi.org/project/yfinance/>.
- [2] Airflow, <https://airflow.apache.org/docs/>.