

Building Building a Stock Price Prediction Analytics using Snowflake & Airflow

(Members: Sharan Patil, Mansi Deshmukh)

I. Problem Statement

The objective of this project was to design and develop an end-to-end financial data analytics system capable of automating the extraction, transformation, and loading (ETL) of real-world stock market data for analytical and forecasting purposes. The system aims to help analysts and investors monitor stock performance, trends, and future projections. The project integrates Python, Apache Airflow, and Snowflake to achieve a seamless workflow from data ingestion to forecasting.

Yahoo Finance (via the yfinance API) was selected as the data source to collect financial market data for selected companies. This system collects historical stock prices for companies such as Microsoft (MSFT) and Oracle (ORCL), covering the last 180 days. Snowflake was chosen as the central data warehouse to store and manage structured financial datasets efficiently. A database and automated data pipelines were essential for ensuring real-time data availability, accuracy, and automation.

II. Solution Requirements

The solution required an automated and scalable mechanism to collect financial data daily, process it, and store it in a structured format for further analytics. The main requirements included data reliability, scalability, and seamless integration between Python and Snowflake, orchestrated through Apache Airflow. The system needed to perform scheduled executions, ensure data consistency, and handle errors gracefully.

Key system requirements included:

- Python environment for data extraction and processing using yfinance and pandas.
- Docker-based Airflow setup for task scheduling and pipeline orchestration.
- Snowflake database integration for data warehousing and analytics.
- Daily automation of ETL processes.
- Forecasting capability using Snowflake ML (or SQL-based forecasting fallback).

While the system achieved its primary objectives, it faced limitations in using the Snowflake ML forecasting feature due to account restrictions. However, an alternative SQL-based forecasting logic was implemented successfully to produce forecast results.

III. Functional Analysis

The proposed system comprises four major components: Data Extraction, Data Storage, Pipeline Orchestration, and Forecasting. Each component contributes to the overall system workflow, from data retrieval to analysis and prediction. Stock symbols used here are MSFT (Microsoft) and ORCL (Oracle).

1. Data Extraction:

The data extraction component retrieves historical stock data from Yahoo Finance using the yfinance API. It captures attributes such as open, high, low, close prices, and trading volume. Data cleaning and validation are handled using pandas to ensure accuracy and uniformity.

2. Data Storage (Snowflake):

The processed data is stored in Snowflake in a table named STOCK_YF under the RAW schema. The table schema includes columns for SYMBOL, DATE, OPEN, HIGH, LOW, CLOSE, and VOLUME, with SYMBOL and DATE serving as primary keys. A transactional load ensures that data is refreshed atomically.

3. Data Pipeline Orchestration (Apache Airflow):

Apache Airflow, running within Docker containers, orchestrates the workflow. The Directed Acyclic Graph (DAG) defines four sequential tasks — extract, ensure_table, load_full_refresh, and forecast_task. These tasks ensure a linear flow, starting from data collection to forecast generation.

4. Forecasting Module:

The forecasting module attempts to use Snowflake's ML. FORECAST functionality. When unavailable, a SQL-based fallback forecasting method is executed. This fallback model uses average daily changes and standard deviation to project 14 days of forecasted prices along with confidence intervals.

IV. Database and Pipeline Flow

The complete system workflow is as follows:

1. Extract historical stock data from the yfinance API.
2. Transform and validate data using pandas.
3. Load cleaned data into the Snowflake STOCK_YF table.
4. Execute forecasting logic and store results in the STOCK_YF_FORECAST table.
5. Schedule automated daily executions using Apache Airflow.

Screenshots proofs of the process and logs according to the flow:

a. Docker app screen

Containers [Give feedback](#)

Container CPU usage 3.54% / 1000% (10 CPUs available) Container memory usage 4.26GB / 7.47GB [Show charts](#)

Search ☐ Only show running containers

Name	Container ID	Image	Port(s)	CPU (%)	Memory usage...	Memory (%)	Disk I/O	Actions
airflow	-	-	-	3.54%	4.26GB / 38.27GB	55.7%	48.1MB / 1.0GB	
postgres-1	36ca70fc9998	postgres:11		0.64%	72.64MB / 7.65GB	0.93%	655MB / 1.0GB	
redis-1	0e7569c14e32	redis:7		0.57%	11.21MB / 7.65GB	0.14%	791MB / 1.0GB	
airflow-init-1	6214e7335de2	apache/air		0%	0B / 0B	0%	0B / 1.0GB	
airflow-webserver	9f5112f19e37	apache/air:8080.8080		0.21%	1.74GB / 7.65GB	22.68%	5.6MB / 1.0GB	
airflow-scheduler	315f9916b595	apache/air		1.74%	421.2MB / 7.65GB	5.37%	29.1MB / 1.0GB	
airflow-worker	9c462d38748c	apache/air		0.38%	2.03GB / 7.65GB	26.58%	11.1MB / 1.0GB	

Showing 7 items

Walkthroughs

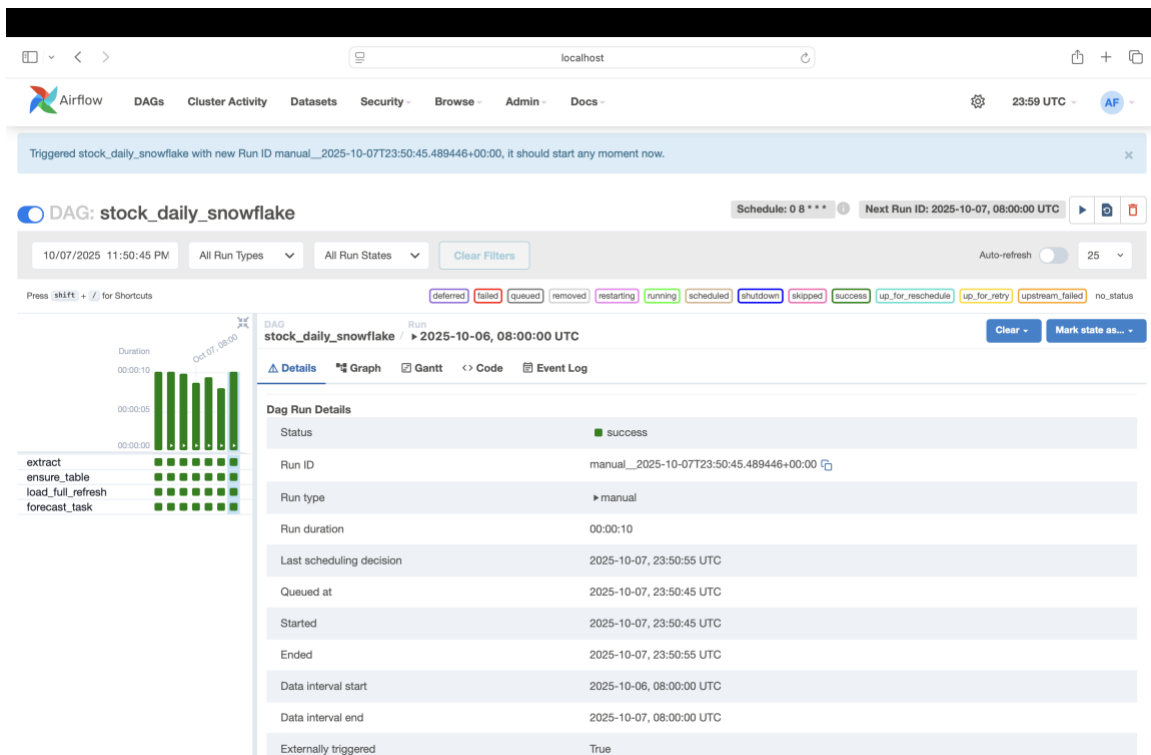
Multi-container applications
8 mins

Containerize your application
3 mins

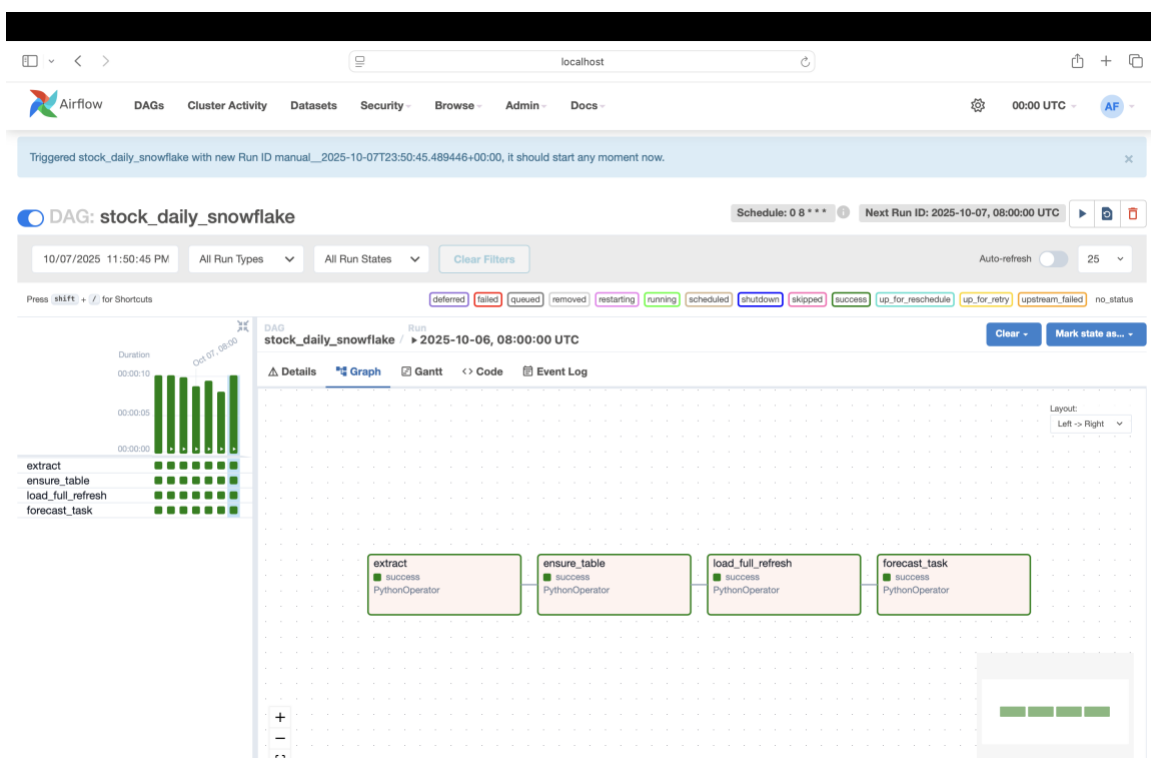
[View more in the Learning center](#)

Engine running | RAM 6.75 GB CPU 0.50% Disk: 6.60 GB used (limit 223.63 GB) [Terminal](#) v4.47.0

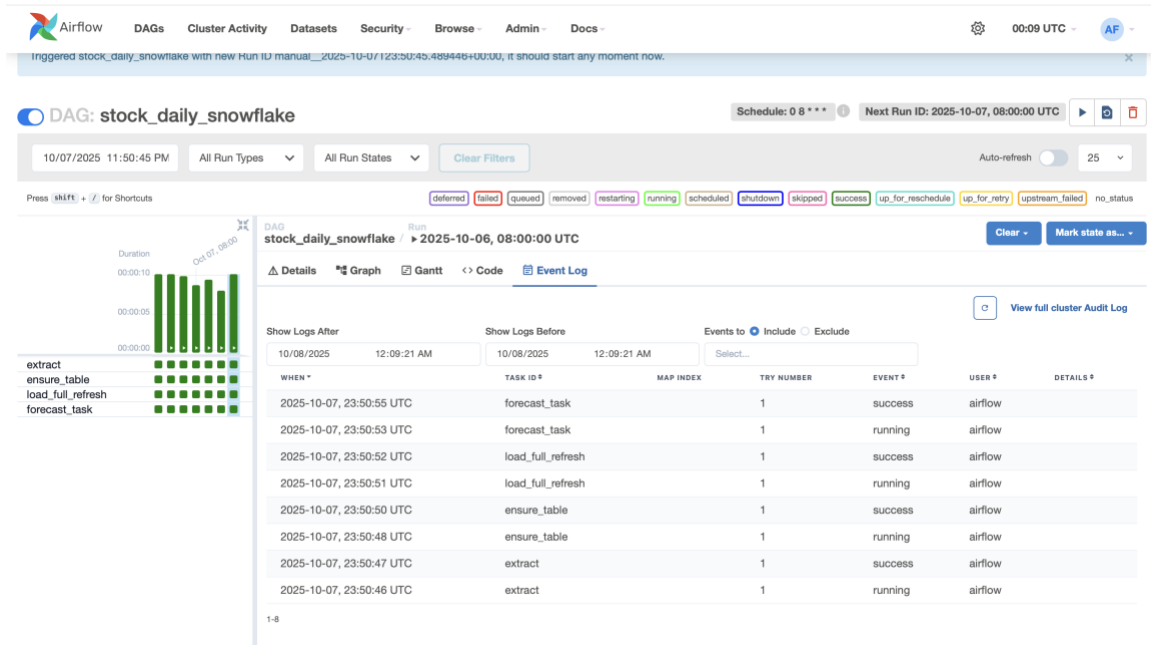
b. Airflow screenshot



c. Airflow DAG



d. Event logs



e. Extract (code)(stock_daily_snowflake.py)

1. Extract

```
def fetch_last_180_days() -> pd.DataFrame:
    """Fetch ~180 days per ticker as a flat DataFrame with clean dtypes."""
    end = datetime.now(timezone.utc).date()
    start = end - timedelta(days=200)

    frames = []
    for t in SYMBOLS:
        # Download per symbol so columns are flat
        df = yf.download(t, start=start, end=end, progress=False).reset_index()

        # ---- Transform (light): normalize columns & dtypes ----
        df.columns = [c.lower() for c in df.columns]
        keep = ["date", "open", "high", "low", "close", "volume"]
        df = df[[c for c in keep if c in df.columns]]

        if "date" in df:
            df["date"] = pd.to_datetime(df["date"]).dt.date # plain date for Snowflake binding

        for col in ["open", "high", "low", "close"]:
            if col in df:
                df[col] = pd.to_numeric(df[col], errors="coerce")

        if "volume" in df:
            df["volume"] = pd.to_numeric(df["volume"], errors="coerce").astype("Int64")

        df["symbol"] = t
        df = df[["symbol", "date", "open", "high", "low", "close", "volume"]]
        frames.append(df)
```

f. Transform (code) (stock_daily_snowflake.py)

2. Transform

```
def transform_prices(df: pd.DataFrame) -> pd.DataFrame:
    # Ensure schema & dtypes exactly as Snowflake table
    cols = ["symbol", "date", "open", "high", "low", "close", "volume"]
    for c in cols:
        if c not in df.columns:
            df[c] = pd.NA

    df["date"] = pd.to_datetime(df["date"]).dt.date
    for col in ["open", "high", "low", "close"]:
        df[col] = pd.to_numeric(df[col], errors="coerce")
    df["volume"] = pd.to_numeric(df["volume"], errors="coerce").astype("Int64")

    return df[cols].sort_values(["symbol", "date"]).reset_index(drop=True)
```

g. Load (code) (stock_daily_snowflake.py)

3. Load

```
def ensure_table():
    create_sql = f"""
    CREATE TABLE IF NOT EXISTS {TABLE_FQN} (
        SYMBOL STRING NOT NULL,
        DATE DATE NOT NULL,
        OPEN NUMBER(18,4),
        HIGH NUMBER(18,4),
        LOW NUMBER(18,4),
        CLOSE NUMBER(18,4),
        VOLUME NUMBER,
        CONSTRAINT PK_STOCK_YF PRIMARY KEY (SYMBOL, DATE)
    );
    """
    with sf_conn() as con, con.cursor() as cur:
        cur.execute(create_sql)

def full_refresh_load():
    """Delete-all + insert the freshly extracted/cleaned data."""
    ensure_table()
    df = fetch_last_180_days()

    if df.empty:
        print("No data fetched from yfinance; skipping insert.")
        return

    rows = [
        (
            str(r["symbol"]),
            r["date"],
            None if pd.isna(r["open"]) else float(r["open"]),
            None if pd.isna(r["high"]) else float(r["high"]),
            None if pd.isna(r["low"]) else float(r["low"]),
        )
```

h. Forecast(code) (stock_daily_snowflake.py) in pipeline

Forecast

```
def run_forecast_or_skip():
    """
    Try SNOWFLAKE.ML.FORECAST (if the account has it). If not present,
    log and exit gracefully so your DAG still stays green.
    """
    with sf_conn() as con, con.cursor() as cur:
        try:
            cur.execute("SHOW FUNCTIONS LIKE 'FORECAST' IN SCHEMA SNOWFLAKE.ML")
            has = cur.fetchall()
        except Exception:
            has = []

    if not has:
        print("SNOWFLAKE.ML.FORECAST not available in this account. Skipping forecast step.")
        return

    cur.execute(f"""
    CREATE OR REPLACE TABLE {SF_DATABASE}.{SF_SCHEMA}.STOCK_YF_FORECAST AS
    WITH BASE AS (
        SELECT SYMBOL, DATE, CLOSE
        FROM {TABLE_FQN}
        WHERE CLOSE IS NOT NULL
    )
    SELECT
        f.SYMBOL,
        f.FORECASTED_TIME AS DATE,
        f.FORECASTED_VALUE AS CLOSE_FORECAST,
        f.LOWER_BOUND AS LCL,
        f.UPPER_BOUND AS UCL
```

i. Snowflake forecast table SQL

```
USER_DB_COYOTE.RAW  Settings  Open in Workspaces  Code Versions  Q

1 CREATE OR REPLACE TABLE RAW.STOCK_YF_FORECAST AS
2 WITH BASE AS (
3     SELECT SYMBOL, DATE, CLOSE
4     FROM RAW.STOCK_YF
5     WHERE CLOSE IS NOT NULL
6 )
7 SELECT
8     f.SYMBOL,
9     f.FORECASTED_TIME AS DATE,
10    f.FORECASTED_VALUE AS CLOSE_FORECAST,
11    f.LOWER_BOUND AS LCL,
12    f.UPPER_BOUND AS UCL
13 FROM TABLE(SNOWFLAKE.ML.FORECAST(
14     INPUT_DATA => TO_VARIANT(BASE),
15     SERIES_COLS => 'SYMBOL',
16     TIMESTAMP_COL => 'DATE',
```

j. Table after pulling data from Yfinance. (STOCK_YF)

```
USER_DB_COYOTE.RAW  Settings  Open in Workspaces  Code Versions  Q

5 SYMBOL STRING NOT NULL,
6 DATE DATE NOT NULL,
7 OPEN NUMBER(18,4),
8 HIGH NUMBER(18,4),
9 LOW NUMBER(18,4),
10 CLOSE NUMBER(18,4),
11 VOLUME NUMBER,
12 CONSTRAINT PK_STOCK_YF PRIMARY KEY (SYMBOL, DATE)
13 );
```

Results Chart

status

1 Table STOCK_YF successfully created.

Query Details

Query duration 285ms

Rows 1

```
USER_DB_COYOTE.RAW  Settings  Open in Workspaces  Code Versions  Q

1 -- Confirm the table exists
2 SHOW TABLES LIKE 'STOCK_YF';
3
4 -- View structure
5 DESC TABLE STOCK_YF;
6
7 -- View data count and date range
8 SELECT
9     COUNT(*) AS ROW_COUNT
```

Results Chart

	SYMBOL	DATE	OPEN	HIGH	LOW	CLOSE	VOLUME
--	--------	------	------	------	-----	-------	--------

Query Details

Query duration 282ms

k. ML forecast table () for MSFT, we can change it to the other symbol according to the need

```

1  USE DATABASE USER_DB_COYOTE;
2  USE SCHEMA RAW;
3
4  SELECT
5      SYMBOL,
6      DATE,
7      CLOSE_FORECAST AS FORECASTED_CLOSE,
8      LCL AS LOWER_CONFIDENCE,
9      UCL AS UPPER_CONFIDENCE

```

	△ SYMBOL	📅 DATE	⌘ FORECASTED_CLOSE	⌘ LOWER_CONFIDENCE	⌘ UPPER_CONFIDENCE
1	MSFT	2025-10-07	529.5896169118	522.970449596	536.208784227
2	MSFT	2025-10-08	530.6092338236	517.370899192	543.847568455
3	MSFT	2025-10-09	531.6288507354	511.771348789	551.486352682
4	MSFT	2025-10-10	532.6484676472	506.171798385	559.12513691
5	MSFT	2025-10-11	533.6680845590	500.572247981	566.763921137
6	MSFT	2025-10-12	534.6877014708	494.972697577	574.402705364
7	MSFT	2025-10-13	535.7073183826	489.373147173	582.041489592
8	MSFT	2025-10-14	536.7269352944	483.773596769	589.680273819
9	MSFT	2025-10-15	537.7465522062	478.174046366	597.319058047
10	MSFT	2025-10-16	538.7661691180	472.574495962	604.957842274
11	MSFT	2025-10-17	539.7857860298	466.974945558	612.596626502
12	MSFT	2025-10-18	540.8054029416	461.375395154	620.235410729
13	MSFT	2025-10-19	541.8250198534	455.77584475	627.874194956
14	MSFT	2025-10-20	542.8446367652	450.176294347	635.512979184

V. Table Structures

Two main tables were created in Snowflake as part of the project:

Table 1: STOCK_YF (Raw and Processed Data)

- SYMBOL – String (Primary Key)
- DATE – Date (Primary Key)
- OPEN – Number(18,4)
- HIGH – Number(18,4)
- LOW – Number(18,4)
- CLOSE – Number(18,4)
- VOLUME – Number

Table 2: STOCK_YF_FORECAST (Forecasted Data)

- SYMBOL – String
- DATE – Date
- CLOSE_FORECAST – Number(18,4)
- LCL – Lower Confidence Limit
- UCL – Upper Confidence Limit

VI. Results and Discussion

The pipeline executed successfully through Apache Airflow, performing all ETL steps and generating a forecast table in Snowflake. The DAG ran smoothly with all tasks marked successful, demonstrating the correctness and stability of the workflow. The SQL-based forecast logic generated reliable trend projections when the Snowflake ML forecasting function was unavailable.

The system achieved daily data automation, accurate record insertion, and fault tolerance. With modular design, additional analytical tasks such as sentiment analysis or portfolio optimization can be easily integrated in future upgrades.

VII. Conclusion

This project effectively demonstrates an automated financial data pipeline using modern cloud-based data engineering tools. The integration of Python, Snowflake, and Apache Airflow enabled efficient data handling, automation, and forecasting. The fallback SQL-based forecast ensured continuous functionality despite Snowflake ML feature limitations. The outcome validates a scalable, practical approach to financial analytics automation suitable for enterprise use.

GITHUB LINK:

<https://github.com/sharan9219790/DATAWAREHOUSE-LAB1>

References

- [1] Apache Airflow Documentation: <https://airflow.apache.org/>
- [2] Snowflake Documentation: <https://docs.snowflake.com/>
- [3] yfinance Library: <https://pypi.org/project/yfinance/>
- [4] Pandas Documentation: <https://pandas.pydata.org/>