**Sharan Patil**

**(018327049)**

**Instructor: Prof. Keeyong Han**

**Homework 8: Running Pinecone Job as an Airflow Pipeline**

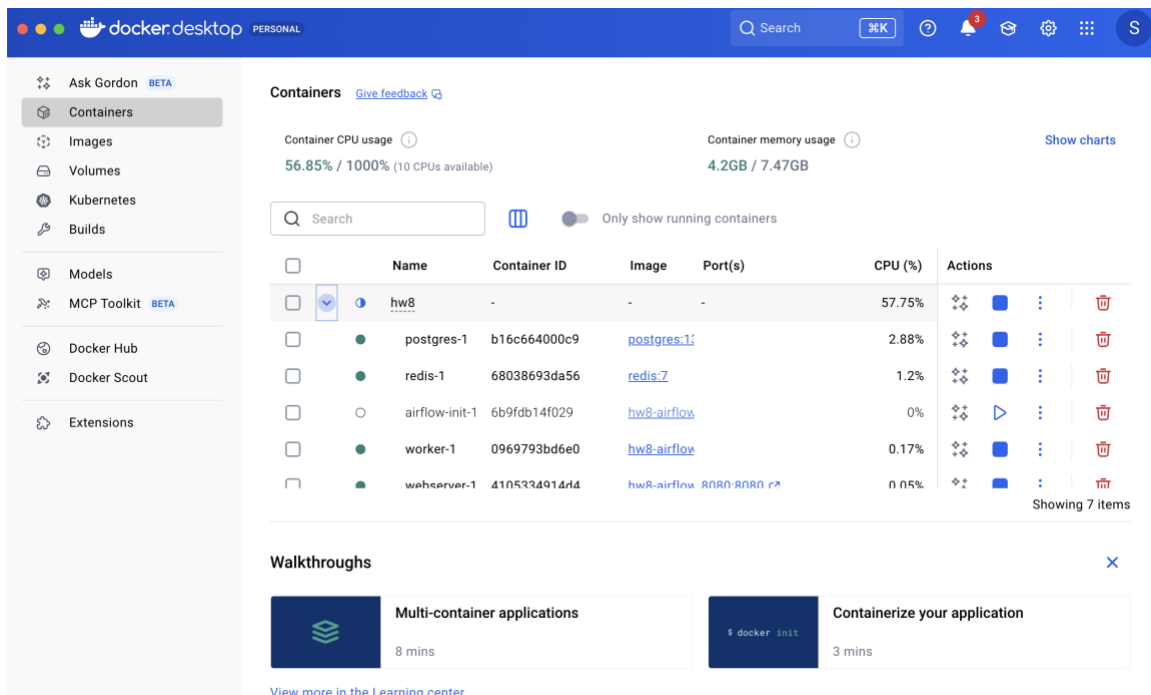**GIT- https://github.com/sharan9219790/HW8_DATAWAREHOUSE**

**Introduction**

This report documents the complete workflow and implementation steps carried out for Homework 8 of the DATA 226 course. The task involved setting up an Apache Airflow pipeline integrated with Pinecone to perform end-to-end semantic search operations. The pipeline automated the process of downloading a dataset, preprocessing it, creating a Pinecone index, converting text into embeddings, and running a semantic search query against the index. The setup was deployed using Docker and verified through successful Airflow task executions.

## 1. Environment Setup

The first step was to set up the development environment using Docker Compose. We created a complete Airflow environment that included the webserver, scheduler, worker, Redis, and Postgres containers. A custom Dockerfile was written to pre-install the required Python libraries (sentence-transformers and pinecone), ensuring that Airflow could seamlessly run tasks involving these dependencies without restart errors. The docker-compose file was configured to expose Airflow on port 8080, map the required volumes, and initialize the Airflow database with an admin account.

We also created an .env file to store environment variables such as AIRFLOW_UID and AIRFLOW_GID to ensure proper file permissions on macOS. Once everything was configured, the setup was initialized and verified using the commands:
docker compose build --no-cache
docker compose up airflow-init
docker compose up -d

## 2. Airflow and Pinecone Configuration

Once the environment was live, we accessed the Airflow UI via http://localhost:8080 and created four Airflow Variables under Admin → Variables. These variables stored Pinecone credentials and configuration details:
- PINECONE_API_KEY
- PINECONE_INDEX_NAME
- PINECONE_CLOUD
- PINECONE_REGION

This approach allowed us to keep sensitive API credentials secure while making them accessible to all Airflow tasks via the Variable.get() method.
A free Pinecone account was used, and when the default region was unsupported, we updated the variables to use gcp/us-central1,
which is compatible with free-tier accounts.

## 3. DAG Creation and Task Implementation

The Airflow DAG file pinecone_medium_pipeline.py was created to define the workflow. The DAG contained five PythonOperator tasks,
each performing a specific stage of the process:

• download_csv → Downloaded the dataset from a public URL (Medium articles dataset).
• preprocess → Processed the CSV file by cleaning, merging title and subtitle, and generating a metadata column.
• create_index → Created a Pinecone serverless index (dimension 384, dotproduct metric). It automatically retried with alternate
regions if the user's free plan didn't support the default region.
• embed_and_upsert → Loaded the SentenceTransformer model (all-MiniLM-L6-v2), generated text embeddings, and uploaded
them to Pinecone using upsert_from_dataframe().
• test_query → Queried Pinecone with a sample semantic search ("what is ethics in AI") and logged the top results.

The DAG was defined with a retry mechanism, two-minute retry delay, and a clear dependency chain ensuring linear execution.
Each task was implemented as a standalone Python function with proper logging for debugging and monitoring.

## 4. Execution and Results

After configuring everything, we triggered the DAG manually from the Airflow UI. Each task executed successfully in sequence.
Initially, we encountered an error when creating the Pinecone index due to unsupported region configuration, which was resolved
by modifying the region to us-central1. Once the pipeline was rerun, all tasks turned green (success).

The logs from the test_query task confirmed that the embeddings and index were functioning correctly. The log entries showed
top-ranked results for the semantic search query, validating that the pipeline successfully encoded, ingested, and queried data
from Pinecone. The final DAG graph view showed all tasks completed successfully.

## LOGS SCREENSHOT FOR EVERY DAG TASK:

## (Next page)

Airflow — localhost:8080/dags/pinecone_medium_pipeline/grid?dag_run_id=manual__2025-11-18T22%3A12%3A00.434537%2B00%3A00&tab=logs&task_id=download_csv

DAGs   Cluster Activity   Datasets   Security ⌄   Browse ⌄   Admin ⌄   Docs ⌄        22:55 UTC ⌄   AU ⌄

Press shift + / for Shortcuts

deferred  failed  queued  removed  restarting  running  scheduled  shutdown  skipped  success  up_for_reschedule  up_for_retry  upstream_failed  no_status

DAG
pinecone_medium_pipeline / ▶ 2025-11-18, 22:12:00 UTC / download_csv

Clear task   Mark state as... ⌄   Filter DAG by task ⌄

⚠ Details   🔲 Graph   📊 Gantt   <> Code   📄 Audit Log   ≡ Logs   ⇄ XCom   ⏱ Task Duration

(by attempts)
[ 1 ]

All Levels                    All File Sources                    ☐ Wrap   Download   See More

```
0969793bd6e0
*** Found local files:
***   * /opt/airflow/logs/dag_id=pinecone_medium_pipeline/run_id=manual__2025-11-18T22:12:00.434537+00:00/task_id=download_csv/attempt=1.log
[2025-11-18, 22:12:01 UTC] {local_task_job_runner.py:120} ▼ Pre task execution logs
[2025-11-18, 22:12:01 UTC] {taskinstance.py:2076} INFO - Dependencies all met for dep_context=non-requeueable deps ti=<TaskInstance: pinecone_medium_pipeline.download_csv man
[2025-11-18, 22:12:01 UTC] {taskinstance.py:2076} INFO - Dependencies all met for dep_context=requeueable deps ti=<TaskInstance: pinecone_medium_pipeline.download_csv manual_
[2025-11-18, 22:12:01 UTC] {taskinstance.py:2306} INFO - Starting attempt 1 of 2
[2025-11-18, 22:12:01 UTC] {taskinstance.py:2330} INFO - Executing <Task(PythonOperator): download_csv> on 2025-11-18 22:12:00.434537+00:00
[2025-11-18, 22:12:01 UTC] {warnings.py:110} WARNING - /home/***/.local/lib/python3.12/site-packages/***/task/task_runner/standard_task_runner.py:61: DeprecationWarning: This
  pid = os.fork()
[2025-11-18, 22:12:01 UTC] {standard_task_runner.py:63} INFO - Started process 202 to run task
[2025-11-18, 22:12:01 UTC] {standard_task_runner.py:90} INFO - Running: ['***', 'tasks', 'run', 'pinecone_medium_pipeline', 'download_csv', 'manual__2025-11-18T22:12:00.43453
[2025-11-18, 22:12:01 UTC] {standard_task_runner.py:91} INFO - Job 19: Subtask download_csv
[2025-11-18, 22:12:01 UTC] {task_command.py:426} INFO - Running <TaskInstance: pinecone_medium_pipeline.download_csv manual__2025-11-18T22:12:00.434537+00:00 [running]> on ho
[2025-11-18, 22:12:01 UTC] {taskinstance.py:2648} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER='spartan' AIRFLOW_CTX_DAG_ID='pinecone_medium_pipeline' AIRFLOW_CTX_TASK_ID
[2025-11-18, 22:12:01 UTC] {taskinstance.py:430} ▲▲▲ Log group end
[2025-11-18, 22:12:01 UTC] {pinecone_medium_pipeline.py:28} INFO - Downloaded: /opt/***/data/medium_data.csv
[2025-11-18, 22:12:01 UTC] {python.py:237} INFO - Done. Returned value was: None
[2025-11-18, 22:12:01 UTC] {taskinstance.py:441} ▼ Post task execution logs
[2025-11-18, 22:12:01 UTC] {taskinstance.py:1206} INFO - Marking task as SUCCESS. dag_id=pinecone_medium_pipeline, task_id=download_csv, run_id=manual__2025-11-18T22:00.43
[2025-11-18, 22:12:01 UTC] {local_task_job_runner.py:240} INFO - Task exited with return code 0
[2025-11-18, 22:12:01 UTC] {taskinstance.py:3498} INFO - 1 downstream tasks scheduled from follow-on schedule check
[2025-11-18, 22:12:01 UTC] {local_task_job_runner.py:222} ▲▲▲ Log group end
```

Version: v2.9.1
Git Version: .release:2d53c1089f78d8d1416f51af60e1e0354781c661

---



Airflow — localhost:8080/dags/pinecone_medium_pipeline/grid?tab=logs&dag_run_id=manual__2025-11-18T22%3A12%3A00.434537%2B00%3A00&task_id=preprocess

DAGs   Cluster Activity   Datasets   Security ⌄   Browse ⌄   Admin ⌄   Docs ⌄        22:55 UTC ⌄   AU ⌄

Press shift + / for Shortcuts

deferred  failed  queued  removed  restarting  running  scheduled  shutdown  skipped  success  up_for_reschedule  up_for_retry  upstream_failed  no_status

DAG
pinecone_medium_pipeline / ▶ 2025-11-18, 22:12:00 UTC / preprocess

Clear task   Mark state as... ⌄   Filter DAG by task ⌄
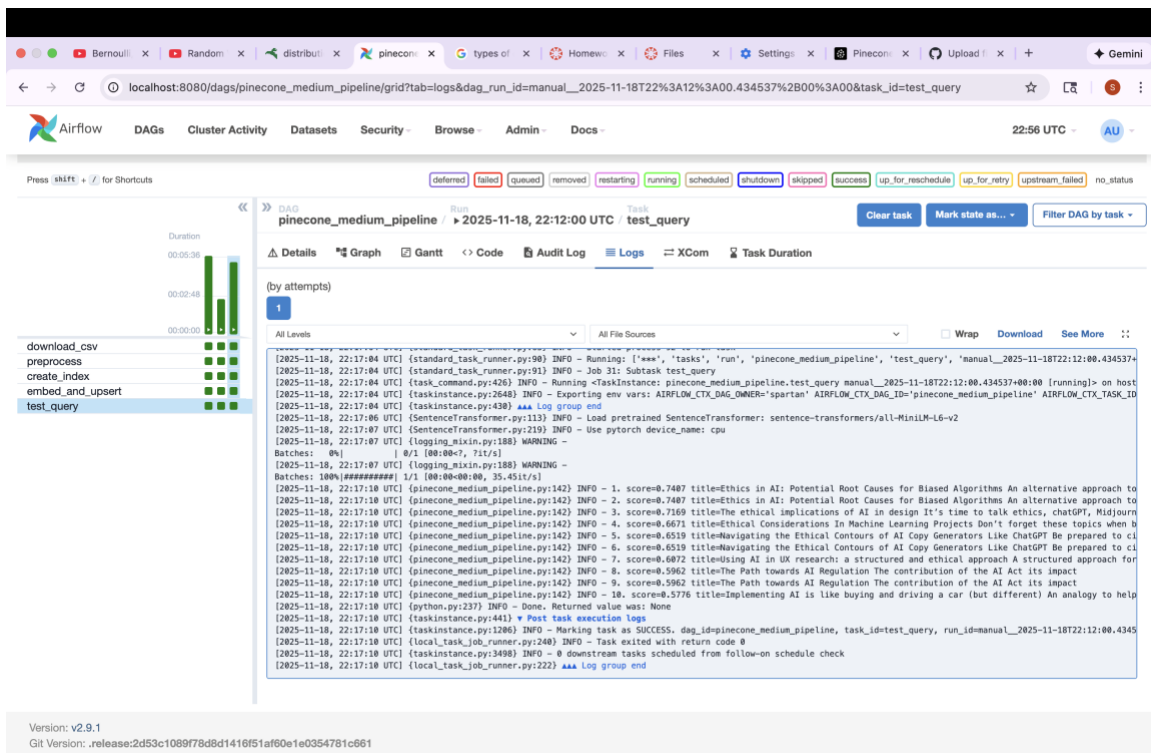
⚠ Details   🔲 Graph   📊 Gantt   <> Code   📄 Audit Log   ≡ Logs   ⇄ XCom   ⏱ Task Duration

(by attempts)
[ 1 ]

All Levels                    All File Sources                    ☐ Wrap   Download   See More

```
0969793bd6e0
*** Found local files:
***   * /opt/airflow/logs/dag_id=pinecone_medium_pipeline/run_id=manual__2025-11-18T22:12:00.434537+00:00/task_id=preprocess/attempt=1.log
[2025-11-18, 22:12:02 UTC] {local_task_job_runner.py:120} ▼ Pre task execution logs
[2025-11-18, 22:12:02 UTC] {taskinstance.py:2076} INFO - Dependencies all met for dep_context=non-requeueable deps ti=<TaskInstance: pinecone_medium_pipeline.preprocess manua
[2025-11-18, 22:12:02 UTC] {taskinstance.py:2076} INFO - Dependencies all met for dep_context=requeueable deps ti=<TaskInstance: pinecone_medium_pipeline.preprocess manual__2
[2025-11-18, 22:12:02 UTC] {taskinstance.py:2306} INFO - Starting attempt 1 of 2
[2025-11-18, 22:12:02 UTC] {taskinstance.py:2330} INFO - Executing <Task(PythonOperator): preprocess> on 2025-11-18 22:12:00.434537+00:00
[2025-11-18, 22:12:02 UTC] {warnings.py:110} WARNING - /home/***/.local/lib/python3.12/site-packages/***/task/task_runner/standard_task_runner.py:61: DeprecationWarning: This
  pid = os.fork()
[2025-11-18, 22:12:02 UTC] {standard_task_runner.py:63} INFO - Started process 211 to run task
[2025-11-18, 22:12:02 UTC] {standard_task_runner.py:90} INFO - Running: ['***', 'tasks', 'run', 'pinecone_medium_pipeline', 'preprocess', 'manual__2025-11-18T22:12:00.434537+
[2025-11-18, 22:12:02 UTC] {standard_task_runner.py:91} INFO - Job 21: Subtask preprocess
[2025-11-18, 22:12:02 UTC] {task_command.py:426} INFO - Running <TaskInstance: pinecone_medium_pipeline.preprocess manual__2025-11-18T22:12:00.434537+00:00 [running]> on host
[2025-11-18, 22:12:02 UTC] {taskinstance.py:2648} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER='spartan' AIRFLOW_CTX_DAG_ID='pinecone_medium_pipeline' AIRFLOW_CTX_TASK_ID
[2025-11-18, 22:12:02 UTC] {taskinstance.py:430} ▲▲▲ Log group end
[2025-11-18, 22:12:02 UTC] {pinecone_medium_pipeline.py:50} INFO - Processed rows: 2498 -> /opt/***/data/medium_processed.parquet
[2025-11-18, 22:12:02 UTC] {python.py:237} INFO - Done. Returned value was: None
[2025-11-18, 22:12:02 UTC] {taskinstance.py:441} ▼ Post task execution logs
[2025-11-18, 22:12:02 UTC] {taskinstance.py:1206} INFO - Marking task as SUCCESS. dag_id=pinecone_medium_pipeline, task_id=preprocess, run_id=manual__2025-11-18T22:12:00.4345
[2025-11-18, 22:12:02 UTC] {local_task_job_runner.py:240} INFO - Task exited with return code 0
[2025-11-18, 22:12:02 UTC] {taskinstance.py:3498} INFO - 1 downstream tasks scheduled from follow-on schedule check
[2025-11-18, 22:12:02 UTC] {local_task_job_runner.py:222} ▲▲▲ Log group end
```

Version: v2.9.1
Git Version: .release:2d53c1089f78d8d1416f51af60e1e0354781c661

## 5. Conclusion

This project demonstrated the integration of Apache Airflow with Pinecone for managing and automating semantic search workflows.
By using Docker, the environment remained fully isolated and reproducible. Airflow ensured that each stage of the data pipeline
executed in the correct order and was fully traceable through logs and UI visualization.

Key learnings from this assignment include:
- Managing Python dependencies in Airflow via Dockerfiles.
- Handling external API integrations (Pinecone) using Airflow Variables.
- Debugging distributed task failures in Airflow.
- Designing a robust, idempotent data ingestion and search workflow.

The final output proved that the pipeline could autonomously process text data, generate embeddings, store them in Pinecone,
and perform vector-based semantic search efficiently. The workflow's modular design allows it to be extended for larger datasets
or other vector databases in the future.
**(Pine cone verification screenshot on next page)**

Pinecone / sjsu ∨ / Default ∨ / Database

Docs    Settings    Get help    SSP

- Get started
- Database
  - Indexes (1)
  - Backups
- Assistant
- Inference
- API keys
- Manage

| Namespace | Operation | ID | Top K |
|---|---|---|---|
| __default__ | Search by ID ∨ | 1 | 10 |

+ Filter    + Rerank    Search

**Search:** 10 results (top_k=10)    1 RUs ⓘ

**1**
ID: 1
**title:** "Not All Rainbows and Sunshine: The Darker Side of ChatGPT Part 1: The Risks and Ethical Issues..."
SCORE
1.0009

**2**
ID: 1634
**title:** "Not All Rainbows and Sunshine: The Darker Side of ChatGPT Part 1: The Risks and Ethical Issues..."
SCORE
1.0009

**3**
ID: 1666
**title:** "ChatGPT—Handle With Care Behind the Hype—Understanding what ChatGPT can do and what it cannot do..."
SCORE
0.5885

STARTER USAGE ⓘ

Storage ⓘ    0.0041GB / 2GB

RUs ⓘ    0 / 1M

WUs ⓘ    19K / 2M

Upgrade now