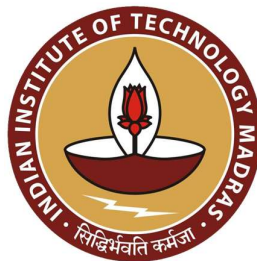


Scrapping of Google News Top Stories Using Airflow

DA5402: MLOPS

Assignment – 2 Report

Rajat Abhijit Kambale [DA24M014]



DEPARTMENT OF DATA SCIENCE AND ARTIFICIAL INTELLIGENCE

INDIAN INSTITUTE OF TECHNOLOGY MADRAS

CHENNAI – 600036

DAGS

1. Scraper DAG:

The news_scraper DAG is in charge of automatically extracting the most important news items from Google News, saving them in a PostgreSQL database, and sending out an email alert whenever new content is published. The first step involves making sure the database is configured properly using a PostgresOperator, which generates a table (news.news_articles) to hold news articles that have been scraped. The PythonOperator, which is at the heart of this DAG, uses a script called scrap.py to extract the most recent news headlines and related metadata. After that, the scraped data is momentarily saved in XComs so that it may be accessed by further tasks.

After the data is available, a different PythonOperator runs a function from dbload.py to add the news stories to the PostgreSQL database. A UNIQUE constraint on headlines and timestamps is used to prevent duplicate entries. As a trigger for the send_email_dag, the DAG then creates a status file with the quantity of recently added articles. In order to enable automated email notifications, the workflow's last task, a TriggerDagRunOperator, initiates the send_email_dag. As a scalable system for automated news monitoring, this structured process guarantees effective and dependable news aggregation, storage, and reporting.



Data stored in postgres

```
news=# SELECT COUNT(*) FROM news.news_articles;
count
-----
      1171
(1 row)

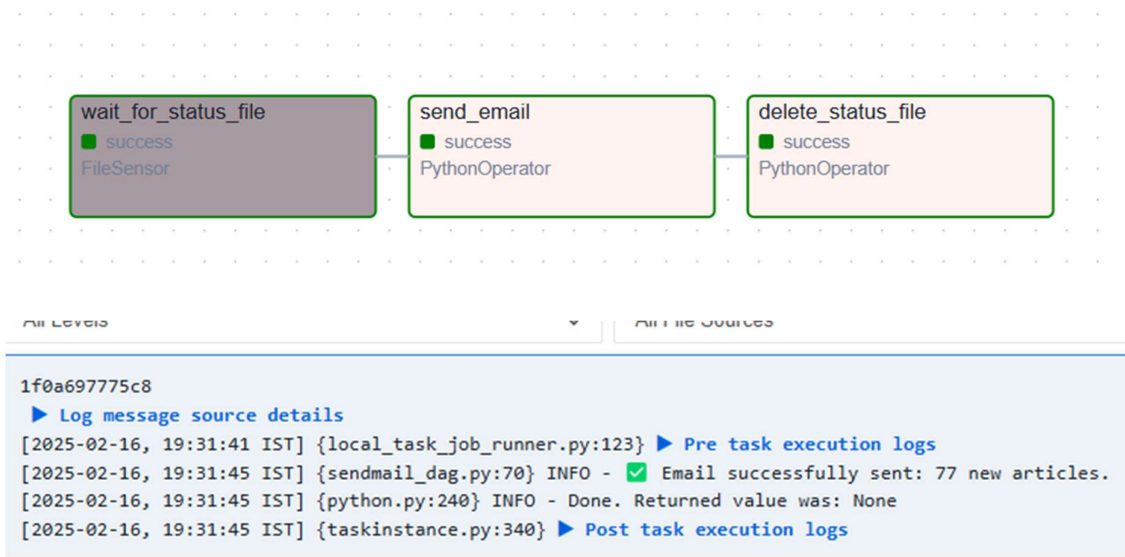
news=#
```

```
1277 | 'It was for everyone's safety': us deportees speak about handcuffing, chains and emotional struggles | UkIGriwEAABXRUQVLA4TCAEAAvD8ACELVQnm1bta30MKj53n7knP90NTH1zu/uxu
ERpIC5cEfVduZ8/+Q0g1KSaztnwFbtm3btm3btvvt2725fe9MQMazrP4KXlzzEcP+Thq5TXL4191/4enEK3C05h0Ab1DNu75QydhgHx9YbX4Dyp9FP8Qpg86812bf40yCPzm7z0IuL55Q1Q0QADAYTV/0M7e3Bhff7g58cIRE/+
PP+Gag0oht8Rm+HB3Z7jQDrX9hzDk7Cz9HV0PC9InAuiG5w64/F9v-rjFKwIQexPkBQ94/+QZ2F8mf+TCKy8Mz71/3622183H3+g2ML30f1YyMT9+3B5IHuExDgguVocF2BJyE0zHFFMzBDVbedxueUhtQHdFfX840rwc4d6pJFTIU8
FMQ94C1Ty+ZAV3vT1g26UvU5GqNMUUCSOL1XwbThuj5ADvht1WZ2YgR5uLud7HCF8Z/Leq076K/nffv5C12Qws0Y3K3K3p6TC65QcAewKOHQ3H3L0sgyGSPK8eqC9ZRA890wckq51cU8KZ3+47K+77coQUEtMFL4mjDj4H89HuxsY2DU
Nog66PjAtywBDU8QYq58Rwz1urxE6vq9NwBm1r1d6xPw6nTY7T/bseJJ+PUTV6QCUFYUj3yYSngTAZBqBx0wYfmeU11d/KX-luH3mjv8fck3pBhfCoegr3c1VEBMO0EnQdWgXtgyEb5SyUEEBbtu6bhsd1bVt9HYSDrnmly1C1CK7AmBj
rnh53ja71WggqL0BPH67wWAgUwMg5B1sAaVjJp5/HAHEIKR05rgvq4Q1HKAUYCLSR2b3B3840JmHfH9b1dyMhE8Bb09gf1T64bHdFC7zG9P892KbE7dnWJowmP82H4Z3Pc+K9Pwmdr7pyUQ336cMS1oufWu8RT1woc10vXJe
uQKqPw8dK7+4wms1cZk4rFjKdC1W4Y50x0pJ/fk07zccLghqfohyf1AtczR0wpmrnb3Dwb1mVjgBZedNukckGqD7Cz8b7ETL4T2u8M11b1LTCV9BmN1XBw6J7E9b8duyqHfQ8h5QZREL7E18nmv/4ZCQEG0MmhyA+
Ngvz7DL3DAPfUu/XokvChnPrwaOymFwJgC0xLoZED12CLEZ7h5uWwPzhCCc1EHks1NvPdWg2BAPAAc1gi1G011bV8Hv2kgHfMh1DX812pyyV6GowdHfcXb8uRuh4/zhBxrvonZ733yCk/9OjXR3HahukBZrHC9CB4VET1wJ
B8031hyfckf0dydyw0A9hCua/XEDRCK1fS31xVYGGdYExvkiK98W75zTh/P8e4k8fUpb1/cV0Uhm7jghhF6/UnH6M29M2YUcv/jekuzw8QXh1/k+ | 2025-02-16 14:00:07.922528
1278 | "Tracherous Routes, Phones Confiscated": Indians recount Deportation Horror | iVBORw0KGgoAAAANSUHElUgAAADkAAASCAyAAACNGTOFAADGELQvRogeZ4M0gLUURi
Av3f08TDLv0u8KtSStQyikk1KtPYRA9BERJ1BUy1UlmQ0KX40EVSyAc1ya1o6E26C2pKvUhmKZ1Ebp7plu5qV1fLdXowRyzzx+hcp/PBw0755798NM0zo2RDPdptVdW6CikYnUTMRBAWHz29844kgpJUV8SoqC1BQfKSkU1IupKQo5
ELRk1Ki1CVFQkgqgQU8X0M8TXHQ43T5/uQ5wMhahhVqM3UjMcnwH01LUlMprBwdd00L/GUD3s9MgSKR9t366nwbH2H0L0LuGfHrzmYhueddMMMH6s7Q718ASU1Pe293pWPLJ3u3q6GZEHFYbVrs46vm/VaZmm29s78Ngmb4VE
y+Q31n7hAGSTwC0S0/qcnaX09Hs4kbjJqehO/XkUQfS0VszTIBcY/mfP5QQ0A8FbD8kBA/ID4uej4rt18A0o2Qfmx83rzRvoQ96Akg1J9u31emaAASxgT64ru4CHzcyEndwMkLWpM4G4AXpb40XkLxbvunb1V11CSnAS75Q0g
pFeaLgus+3J0Hk4z18pqC2tq+Y3317f+fnzwfL/062+XU5WQ8ELO2Fp/XEwAIXrocU2g48Uz3t1Pqy053klHcayUg8mbzVU1p3123o/0Bz02Gw6uzi7unt1qz6VYettFYPRyQxde4H+X0xeVEB9lUddY29ut4/kQ5A/7Q10g6kyJ8
64TWag/fRk0zauh1D0n1geUfzoEYJQC8a2j01au7o4PKzAwAZctMEbkmTEhoaiqssUgh9CP8P2f700PprzrtcuY169k2v04bf3tmmWmp/RuMxLCK2fYlUR8mJPPAPACDM8zqak+7p71d788wRG7cwam1W4h3A1Cffx37oxrdX7q0
Vmc1Ejnmvz1BPwld8vVACDTGuEfw0u536/8tCjYtFdTp19+rp7qb8+8Htu6qg25ShkL3Mgeq48r4B/kv3n1kPchT5VgkqTgJUV8SoqC1BQfKSkU1IupKQo5ELR+D8kFUbX/u0hfiEko1i/AtmA9jNnsP1AAAAAE1fKSuqCC
| 2025-02-16 14:00:07.922528
```

2. Send Email DAG:

When new articles are added to a database, the `send_email_dag` is intended to automatically send out email notifications. It starts with a `FileSensor` that keeps an eye on a specified status file (`status`) all the time. It signifies the addition of new articles if the file is present. After that, a `PythonOperator` is triggered by the DAG, which creates an email notice after reading the file's number of new articles. Using login credentials that are safely kept in Airflow Variables, this email is delivered via SMTP. The quantity of newly discovered articles is dynamically reflected in the subject and text of the email. The job logs the occurrence and ends without sending an email if no new entries are discovered.

After the email is sent, the status file is deleted by another `PythonOperator` as the last step in the DAG to guarantee the system stays clean. By doing this, redundant alerts are avoided, and notifications are only triggered by new data in subsequent runs. Because the DAG is set up with error handling and recording, it is possible to monitor problems (such SMTP authentication errors). Furthermore, Airflow's retry procedures make sure that the pipeline as a whole is not disrupted by transitory failures (like SMTP outages).






News Update - 77 Articles Added Inbox x



rajatkambale02@gmail.com

to me ▼

Hello,

 77 new articles were added to the database.

Regards,
Your Airflow Pipeline