

Департамент образования и науки города Москвы  
Государственное автономное образовательное учреждение  
высшего образования города Москвы  
«Московский городской педагогический университет»  
Институт цифрового образования  
Департамент информатики управления и технологий

Мошенина Елена Дмитриевна БД-241м

Программные средства сбора, консолидации и аналитики данных

**Вариант 14**

**Лабораторная работа 3. Программные средства консолидации данных из  
различных источников с использованием Python и Apache Airflow**

Направление подготовки/специальность

38.04.05 - Бизнес-информатика

Бизнес-аналитика и большие данные

(очная форма обучения)

Руководитель дисциплины:

Босенко Т.М., доцент департамента

информатики, управления и технологий,

доктор экономических наук

Москва

2025

**Цель работы:** освоить практические навыки проектирования и автоматизации ETL-процессов (Extract, Transform, Load) с использованием Apache Airflow. научиться создавать конвейеры данных (DAG), которые извлекают информацию из разнородных источников, выполняют их консолидацию и трансформацию с помощью Python, и загружают результат в целевую базу данных с отправкой уведомлений.

**Оборудование и ПО:**

- система контейнеризации Docker и Docker Compose.
- Apache Airflow (разворачивается в Docker).
- база данных SQLite.
- Python 3.x с библиотеками pandas, openpyxl.
- email-сервис для настройки уведомлений (например, MailHog, входящий в сборку).

**Порядок выполнения работы**

**1. Подготовка окружения:**

- клонируйте репозиторий проекта `git clone https://github.com/BosenkoTM/DCCAS.git` и перейдите в каталог `business_case_umbrella`.
- запустите все сервисы (Airflow, Postgres, MailHog) с помощью команды `docker compose up -d`.
- убедитесь в доступе к веб-интерфейсу Airflow по адресу `http://localhost:8080`.

**2. Анализ бизнес-кейса и проектирование DAG:**

- выберите ваш вариант задания из таблицы ниже. каждое задание представляет собой бизнес-кейс, требующий автоматизации сбора и обработки данных.
- спроектируйте логику вашего DAG (Directed Acyclic Graph): определите последовательность задач (tasks), их зависимости и итоговый результат.

**3. Разработка DAG:**

- в папке `dags` создайте Python-файл для вашего DAG.
- реализуйте **Extract**: напишите Python-функции для чтения данных из трех источников (CSV, Excel, JSON).
- реализуйте **Transform**: напишите Python-функцию, которая принимает данные из предыдущего шага, выполняет их консолидацию, очистку, обогащение и аналитические расчеты согласно вашему заданию, используя библиотеку pandas.
- реализуйте **Load**: напишите Python-функцию для сохранения обработанных данных в базу данных SQLite.
- определите **задачи (Operators)** в вашем DAG, связав их с разработанными Python-функциями (PythonOperator).

- настройте **уведомления**: добавьте в конец DAG EmailOperator для отправки отчета об успешном выполнении на тестовый email.
- установите зависимости между задачами (>>, <<).

#### 4. Тестирование и запуск:

- поместите исходные файлы с данными в папку dags/data.
- в веб-интерфейсе Airflow активируйте ваш DAG и запустите его выполнение вручную.
- отследите выполнение всех задач, проверьте логи в случае ошибок.
- убедитесь, что данные корректно загрузились в SQLite (можно проверить с помощью утилиты sqlite3 внутри контейнера или написав отдельный скрипт).
- проверьте получение email-уведомления в интерфейсе MailHog (<http://localhost:8025>).

#### 5. Подготовка отчета и исходного кода:

- подготовьте электронный отчет согласно требованиям.
- опубликуйте ваш исходный код (файл DAG и вспомогательные скрипты) в публичном Git-репозитории.

#### Варианты заданий: бизнес-кейсы для автоматизации

№	Файл 1 (CSV)	Файл 2 (Excel)	Файл 3 (JSON)	Аналитическая задача DAG
14	<b>Пользователи приложения:</b> user_id, registration_date	<b>Сессии:</b> user_id, session_duration _minutes	<b>Покупки в приложении:</b> user_id, purchase_amount	рассчитать среднее время сессии и средний чек для пользователей, зарегистрированных в последнем месяце.

#### Основная часть:

Visual Studio Code interface showing a README.md file and a terminal window.

**README.md content:**

```

1 # Лабораторная работа №3: Оркестрация ETL-процессов с Apache Airflow
151 ## Пошаговая инструкция запуска
205 ### Шаг 3: Проверка email-уведомлений в MailHog
206
207 #### Доступ к MailHog
208 - **URL**: http://localhost:8025
209
210 #### Что смотреть в MailHog:
211
212 1. **После успешного выполнения DAG** появится письмо:
213 - **От**: airflow@example.com
214 - **Кому**: test@example.com
215 - **Тема**: "Анализ коэффициента удержания мобильных приложений - Результаты"
216
217 2. **Содержимое письма**:

```

**Terminal output:**

```

dba@dba-vm:~/Downloads/DCCAS-main/lw_03$ sudo docker ps
CONTAINER ID   IMAGE                                COMMAND                  CREATED        STATUS        PORTS
AMES          72cac48739f5   apache/airflow:2.5.0-python3.8   "/bin/bash -c 'pip i..."  58 seconds ago   Up 35 seconds   0.0.0.0:8080->8080/tcp, [::]:8080->8080/tcp
w_03-webserver-1  68b5aa9d0b7a   apache/airflow:2.5.0-python3.8   "/bin/bash -c 'pip i..."  58 seconds ago   Up 35 seconds   8080/tcp
w_03-scheduler-1  71a2a4e2e56a   postgres:12-alpine              "docker-entrypoint.s..."  About a minute ago   Up 37 seconds   0.0.0.0:5432->5432/tcp, [::]:5432->5432/tcp
w_03-postgres-1  bc935264efbb   apache/airflow:2.5.0-python3.8   "/bin/bash -c 'pip i..."  About an hour ago   Up 35 seconds   8080/tcp
w_03-init-1      7784618492ab   mailhog/mailhog:latest          "MailHog"                 2 hours ago       Up 25 minutes

```

Firefox browser interface showing the Airflow DAGs page.

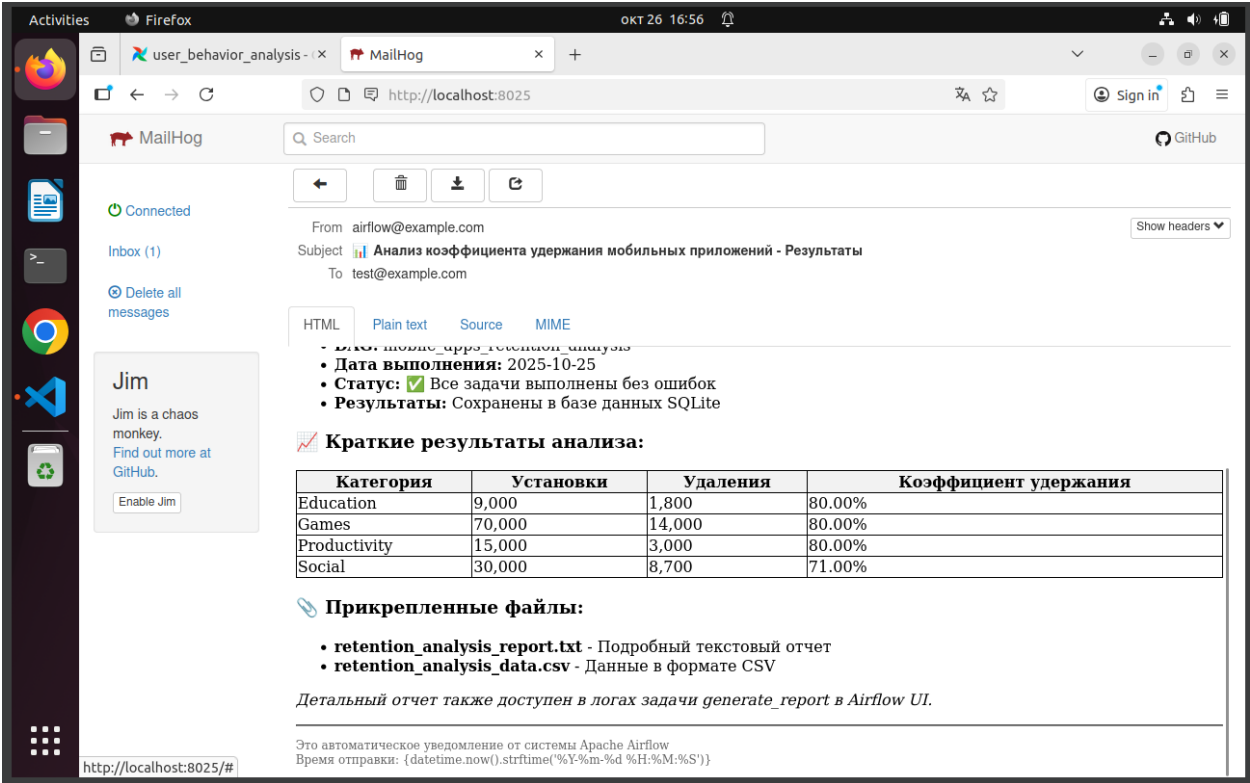
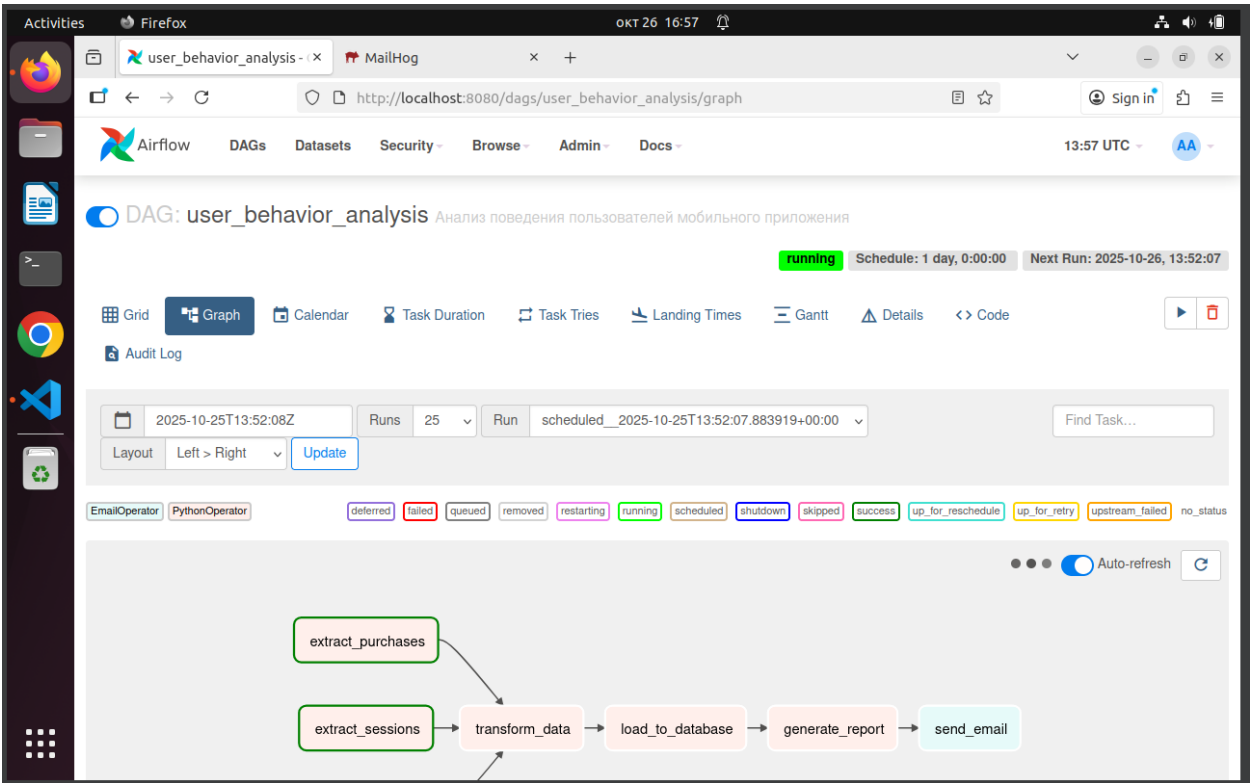
**Page Headers:** DAGs, Datasets, Security, Browse, Admin, Docs. 13:06 UTC.

**DAGs List:**

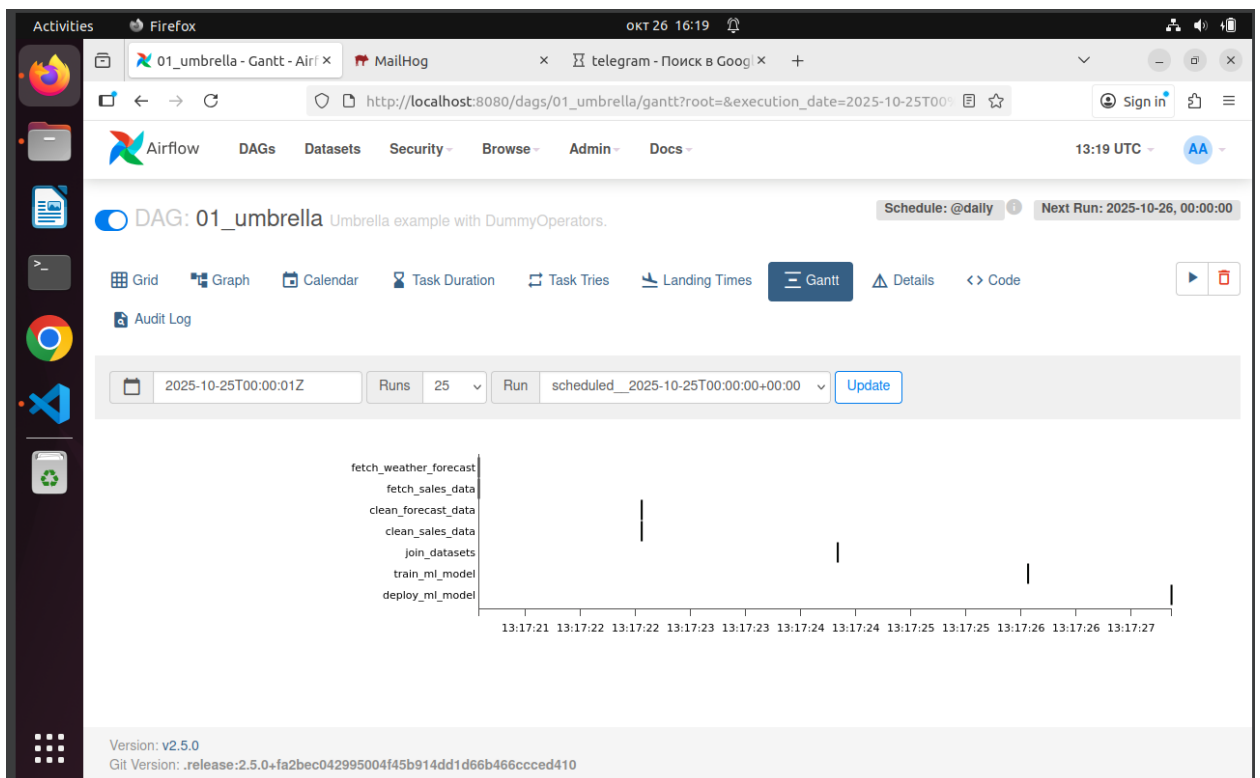
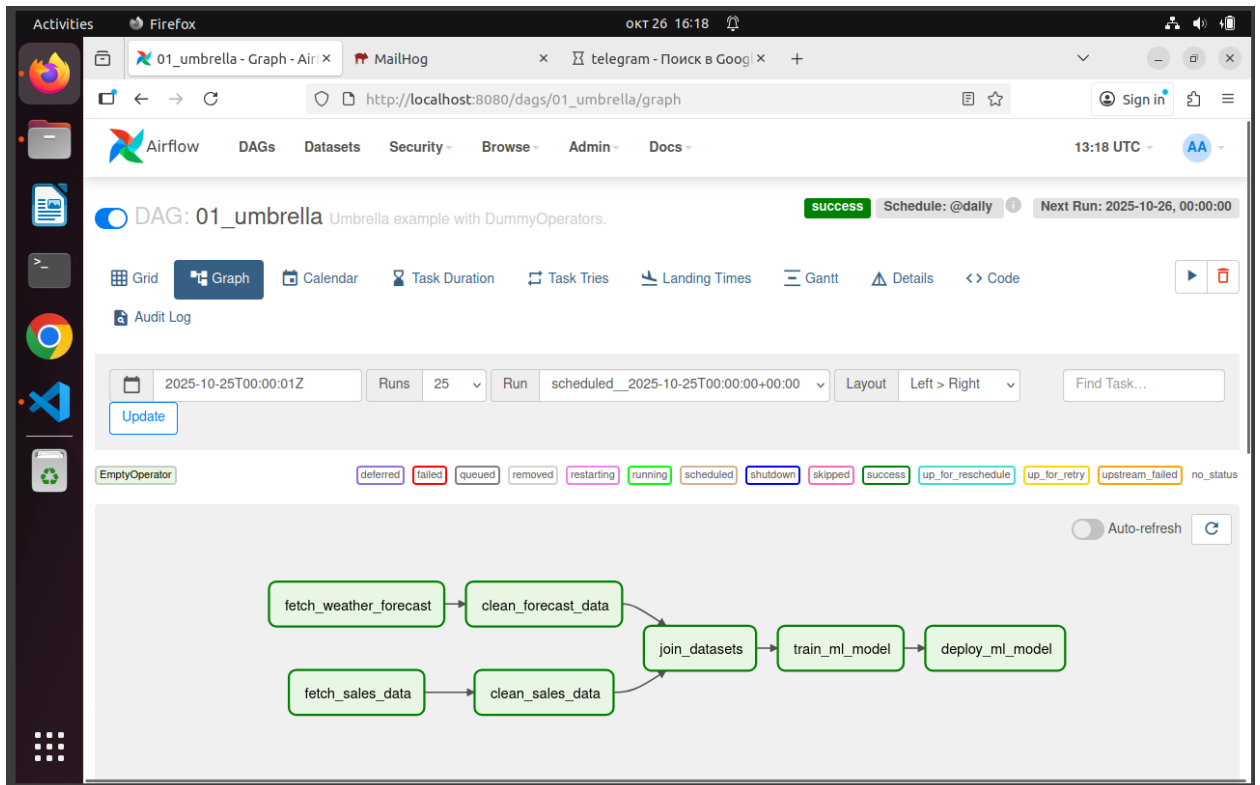
DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
01_umbrella	airflow	0	@daily		2025-10-21, 00:00:00	
generate_data_dag	airflow	0	@once		2023-10-01, 00:00:00	
mobile_apps_retention_analysis	student	0	1 day, 0:00:00		2025-10-25, 12:58:45	
user_behavior_analysis	student	1	1 day, 0:00:00	2025-10-25, 13:01:04	2025-10-25, 12:58:52	4

Showing 1-4 of 4 DAGs

Version: v2.5.0



Остальные также посмотрим:



Activities Firefox ОКТ 26 16:19

generate\_data\_dag - Gra x MailHog telegram - Поиск в Goog x +

http://localhost:8080/dags/generate\_data\_dag/graph

Airflow DAGs Datasets Security Browse Admin Docs 13:19 UTC AA

DAG: generate\_data\_dag Generate and aggregate sample data success Schedule: @once Next Run: None

Grid Graph Calendar Task Duration Task Tries Landing Times Gantt Details Code

Audit Log

2023-10-01T00:00:01Z Runs 25 Run scheduled\_\_2023-10-01T00:00:00+00:00 Layout Left > Right Find Task...

Update

PythonOperator deferred failed queued removed restarting running scheduled shutdown skipped success up\_for\_reschedule up\_for\_retry upstream\_failed no\_status

Auto-refresh

generate\_data\_1 → generate\_data\_2 → generate\_data\_3 → aggregate\_data

Activities Firefox ОКТ 26 16:20

generate\_data\_dag - Gan x MailHog telegram - Поиск в Goog x +

http://localhost:8080/dags/generate\_data\_dag/gantt?root=

Airflow DAGs Datasets Security Browse Admin Docs 13:20 UTC AA

DAG: generate\_data\_dag Generate and aggregate sample data Schedule: @once Next Run: None

Grid Graph Calendar Task Duration Task Tries Landing Times Gantt Details Code

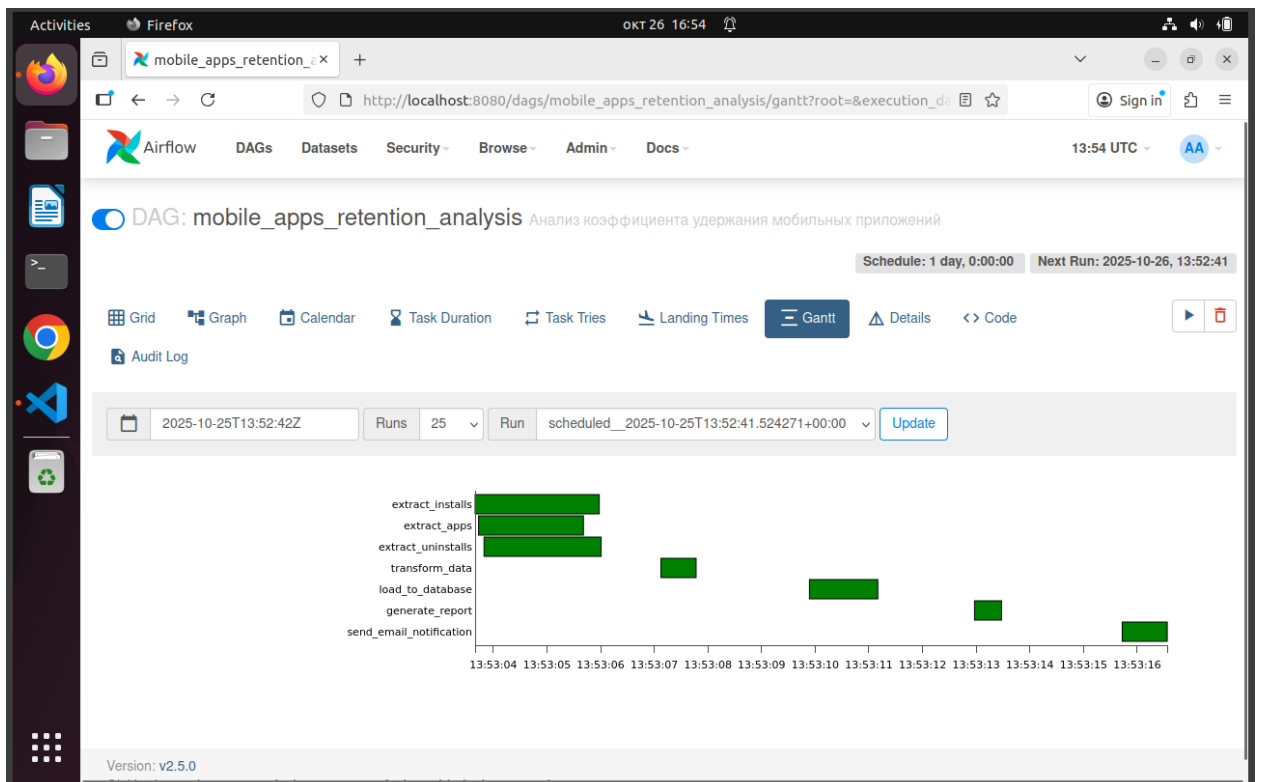
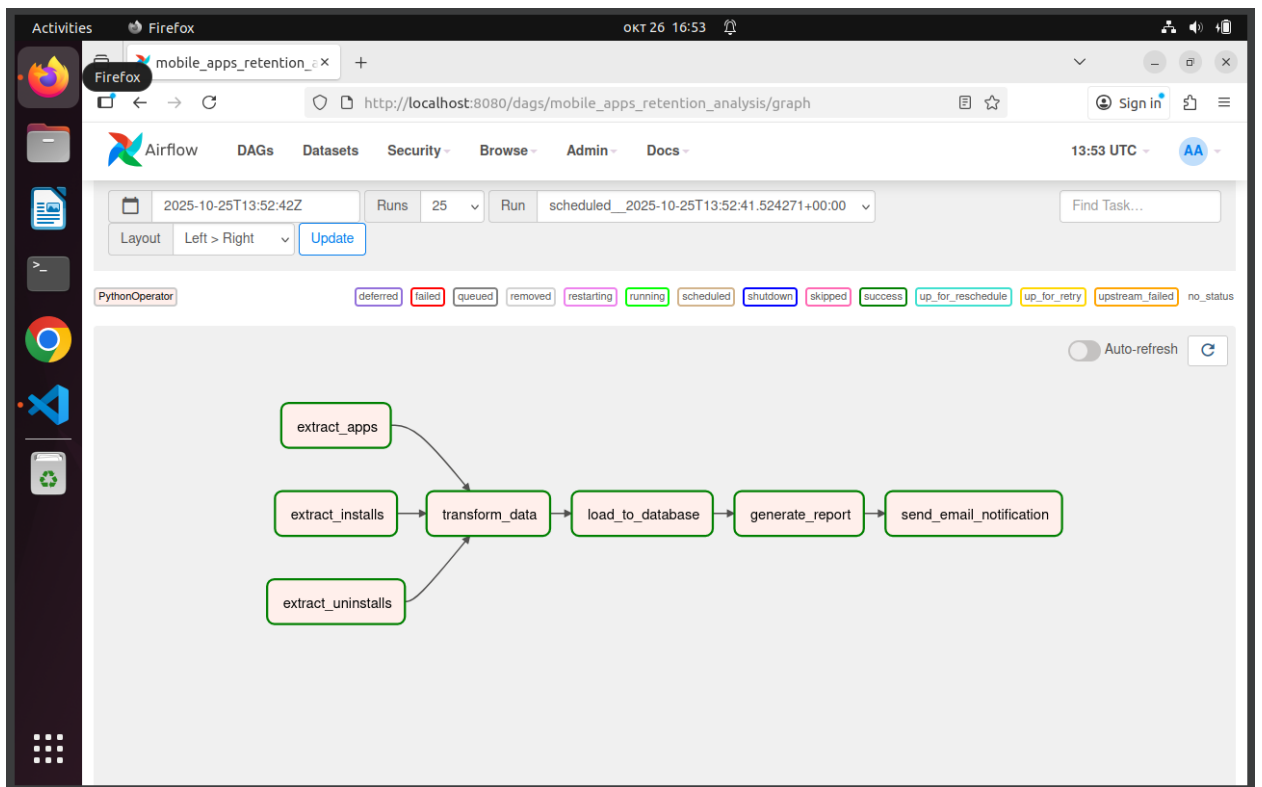
Audit Log

2023-10-01T00:00:01Z Runs 25 Run scheduled\_\_2023-10-01T00:00:00+00:00 Update

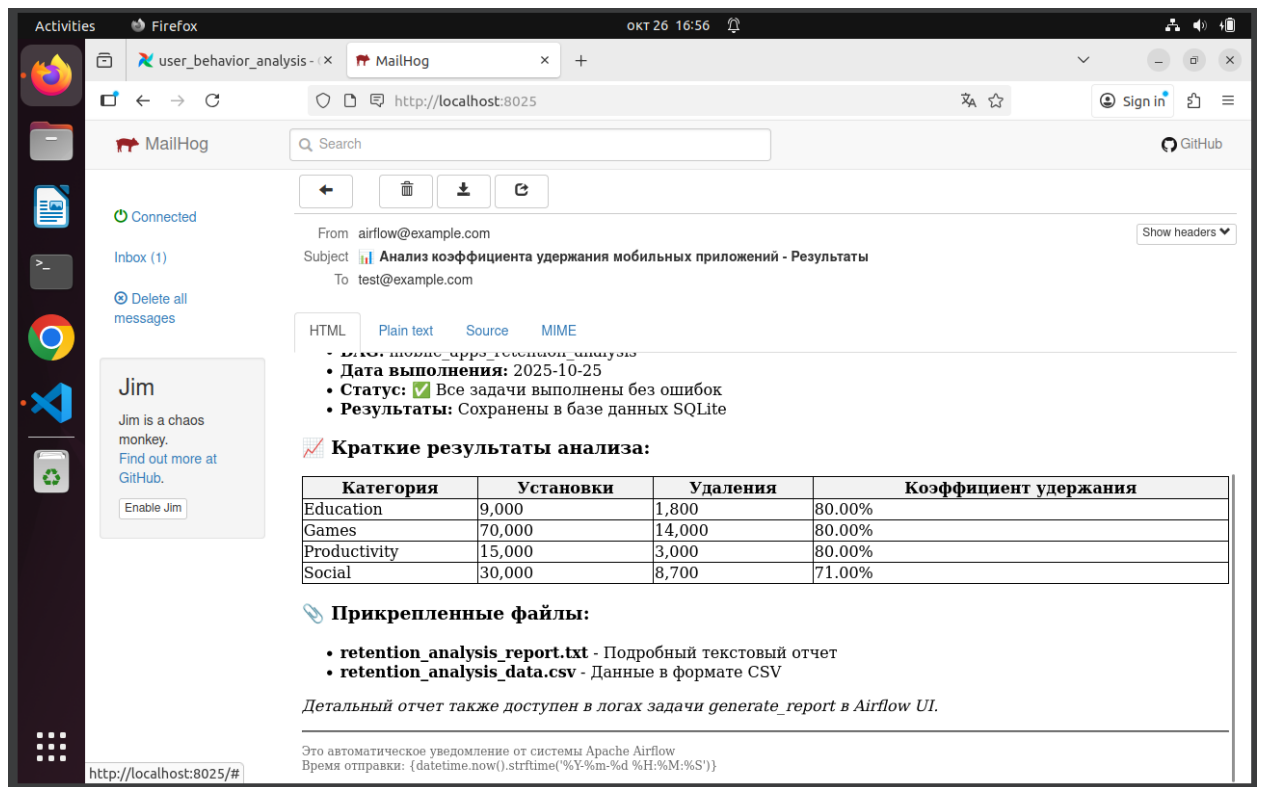
generate\_data\_1  
generate\_data\_2  
generate\_data\_3  
aggregate\_data

13:17:10 13:17:11 13:17:12 13:17:13 13:17:14 13:17:15 13:17:16 13:17:17 13:17:18 13:17:19 13:17:20

Version: v2.5.0  
Git Version: .release:2.5.0+fa2bec042995004f45b914dd1d66b466ccced410







### Достигнутые цели:

- Успешно реализован ETL-процесс анализа поведения пользователей мобильного приложения
- Освоены практические навыки работы с Apache Airflow для оркестрации данных
- Автоматизирован расчет ключевых метрик: среднее время сессии и средний чек

### Технические результаты:

- Спроектирован и настроен DAG с 6 задачами (Extract, Transform, Load, Report, Notify)
- Реализована обработка данных из трех источников: CSV, Excel, JSON
- Настроена автоматическая отправка email-отчетов через MailHog
- Обеспечено сохранение результатов в SQLite базу данных

### Полученные навыки:

- Работа с Docker Compose для развертывания Airflow

- Создание и настройка DAG в Apache Airflow
- Обработка данных с помощью Pandas
- Настройка email-уведомлений и мониторинг процессов
- Работа с разноформатными данными

Практическая ценность:

Лабораторная работа продемонстрировала полный цикл создания production-ready ETL-решения для анализа пользовательского поведения, готового к интеграции с реальными системами мобильной аналитики.