# DA5402: Machine Learning Operations Laboratory

# Assignment 4

Nandhakishore C S (DA24M011)

March 14, 2025

## Problem Statement

We have learned about scrapping websites into a dataset of images and captions. Now, let's turn our attention to RSS feeds. RSS feeds are the XML structured data streams that news websites typically expose for easier consumption into a software application.

Python has a 'feedparser' library that can read an RSS feed either from a saved XML file or a web URL. The following code snippet allows one to read the feed directly into a Python dictionary. Here, we have used the RSS feed from The Hindu newspaper.

```
feed = feedparser.parse('https://www.thehindu.com/news/national/?service=rss')
```

The generated dictionary is organized as entries, where each entry contain the details of a specific news article streamed via the RSS feed. A typical dict structure of an entry without media content.

```
{'title': 'B-khata campaign to be inaugurated in MCC limits on March 1',
 'title_detail': {'type': 'text/plain',
  'language': None,
  'base': 'https://www.thehindu.com/news/national/?service=rss',
  'value': 'B-khata campaign to be inaugurated in MCC limits on March 1'},
 'summary': '',
 'summary_detail': {'type': 'text/html',
  'language': None,
  'base': 'https://www.thehindu.com/news/national/?service=rss',
  'value': ''},
 'links': [{'rel': 'alternate',
   'type': 'text/html',
   'href':
 'https://www.thehindu.com/news/national/karnataka/b-khata-campaign-to-be-inaugurated-in-mcc-limits-on-march-1/
 article69274675.ece'}],
 'link':
 'https://www.thehindu.com/news/national/karnataka/b-khata-campaign-to-be-inaugurated-in-mcc-limits-on-march-1/
 article69274675.ece',
 'id': 'article-69274675',
 'guidislink': False,
 'tags': [{'term': 'Karnataka', 'scheme': None, 'label': None}],
 'published': 'Fri, 28 Feb 2025 18:30:37 +0530',
 'published_parsed': time.struct_time(tm_year=2025, tm_mon=2, tm_mday=28, tm_hour=13, tm_min=0,
 tm_sec=37, tm_wday=4, tm_yday=59, tm_isdst=0)}
```

The structure of the RSS feed differs from one provider to another. You may have to review the dictionary to suite to your needs. Alright, let's get the business now. The objective is to setup an RSS reader which will automatically fetch and process the required news articles into a list of tuples. We shall push these tuples into a database table of our choice to persistent storage. Beware the the RSS feed is a dynamic stream. Depending on when you read, the data received shall be different. Nonetheless, when you read the stream to often, you may get duplicates as expected. The tuple should have the following information, which eventually get stored in the persistent database table.

1. Title (cannot be blank)

2. Publication Timestamp (to be stored as datetime)

3. Weblink to the article (no need to scrape the page; can't be blank)

4. The picture of the article. (the actual image should be downloaded; can be blank)

5. Tags (can be one or more tag terms)

6. Summary (can be blank)

   * **Necessary logs should be created for debugging and information needs.**

# Task 1 [20 points]

Create a docker container of your favorite database. Add scripts to initialize your database with the necessary credentials and create the necessary table(s) to save the data with the above-listed columns. So, when the container is up (during the service creation), the database should get created with the supplied credentials, tables are to be created. The database gets created when the docker container is created for the first time. When the container is restarted, there should be a checker that validates the existence of the database and the required tables. If the validation fails, the db  table creation should be redone again. If there is some issue with the initialization, the container should not start up. The required configuration settings from setting up the db should be via environment variables.

A corresponding Dockerfile is created for virtualising the database creating process.

```
FROM postgres:15
ENV POSTGRES_USER=news_admin
ENV POSTGRES_PASSWORD=securepassword
ENV POSTGRES_DB=news_db
COPY init.sql /docker-entrypoint-initdb.d/
```

Figure 1: Dockerfile for Task1

A shell script is used to initiate the data using postgres server with required sql queries.

```
#!/bin/bash
set -e
echo "Checking database initialization..."

PGPASSWORD=$POSTGRES_PASSWORD psql -U $POSTGRES_USER -d $POSTGRES_DB -tc "SELECT 1 FROM pg_database
    WHERE datname='$POSTGRES_DB';" | grep -q 1 || {
    echo "Database $POSTGRES_DB does not exist, initializing..."
    psql -U $POSTGRES_USER -d postgres -c "CREATE DATABASE $POSTGRES_DB;"
}

TABLE_EXISTS=$(PGPASSWORD=$POSTGRES_PASSWORD psql -U $POSTGRES_USER -d $POSTGRES_DB -tc "SELECT
    to_regclass('public.news_articles');")

if [[ "$TABLE_EXISTS" != "news_articles" ]]; then
    echo "Table does not exist. Creating..."
    psql -U $POSTGRES_USER -d $POSTGRES_DB -f /docker-entrypoint-initdb.d/init.sql
else
    echo "Table exists. Skipping initialization."
fi

echo "Database is ready."
exec docker-entrypoint.sh postgres
```

Figure 2: Shell script to initialise database

A set of SQL queries are send to create a table with the necessary aatributes.

```
CREATE TABLE IF NOT EXISTS news_articles (
    id SERIAL PRIMARY KEY,
    headline TEXT NOT NULL,
    url TEXT NOT NULL UNIQUE,
    published_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
    thumbnail BYTEA,  -- Binary storage for images
    summary TEXT
);
```

Figure 3: SQL script for creating tables

# Task 2 [20 points]

Build an application (a python runnable script here) that would fetch the news articles from the configured RSS feed and push the required fields into the containerized database. The URL of the RSS feed, the dictionary paths of the required fields are to be configurable via environment variables. This way, we can configure the application to fetch from any RSS feed from source. For now, we are hanging on to The Hindu. But the expectation is to make it work for others such as, TOI, Indian Express, NDTV, India Today, News18, etc. Dockerize the application. The application should poll the RSS feed for updates every 10 mins (configurable) and fetch the feed when the data has changed. Mind that data change does not mean no-duplicates.

For creating the database in a docker, the psycopg library is used with required credentials (db_name, password, port number, etc). There is a user defined function which creates the database and checks for existing database when the docker is run again. A.env file for configuring the docker file with database credentials is created to make the process easy.

```
1  # RSS Feed Configuration
2  RSS_FEED_URL=https://www.thehindu.com/news/national/feeder/default.rss
3  FETCH_INTERVAL=600   # Time interval in seconds
4
5  # PostgreSQL Configuration
6  POSTGRES_HOST=database   # Matches Docker service name
7  POSTGRES_PORT=5432
8  POSTGRES_DB=news_db
9  POSTGRES_USER=news_admin
10 POSTGRES_PASSWORD=securepassword
```

Figure 4: Database credentials for Task1

A separate Dockerfile is created to virtualise task2 (fetching RSS separately)

```
1  FROM python:3.10
2  WORKDIR /app
3  COPY fetch_news.py requirements.txt ./
4  RUN pip install --no-cache-dir -r requirements.txt
5  CMD ["python", "fetch_news.py"]
```

Figure 5: DockerFile for task2

The project directory is structured as follows:

# Task 3 [10 points]

We got two containerized applications. Use docker-compose to create a multi-application (services) dockerization such that the containerized RSS reader app could push the data into the containerized DB. Upon docker-compose up, both services should be created and applications should be started for RSS reading and db storage. Docker-compose down should stop services and remove them. Likewise, docker-compose start and stop should have the necessary functionalities.

For multi application services, a common docker-compose file is created and both the docker container's details are read. This file initiates the docker and creates two containers and make them work in tandem by passing the information from one to other.

# Brownie [10 points]

Create a simple containerized web application to read the data from the database with suitable date filters (default is today's news). The web application should show the Title, Image  Summary from the database for every record in the table. Additionally, upon clicking the news title or image, you should open another tab with the actual news using the URL that you got from the db table. The web application should also become a part of the docker-compose environment, which will now have three services. Up, down, start and stop commands of docker-compose should have the necessary functionalities.

A simple RESTAPI based web application is created in python with a standard HTML, CSS layer.

```python
# Load environment variables
RSS_FEED_URL = os.getenv("RSS_FEED_URL", "https://www.thehindu.com/news/national/feeder/default.rss")
FETCH_INTERVAL = int(os.getenv("FETCH_INTERVAL", 600))  # Default: 10 minutes
DB_CONFIG = {
    "dbname": os.getenv("POSTGRES_DB", "news_db"),
    "user": os.getenv("POSTGRES_USER", "news_admin"),
    "password": os.getenv("POSTGRES_PASSWORD", "securepassword"),
    "host": os.getenv("POSTGRES_HOST", "database"),  # Updated to match service name
    "port": os.getenv("POSTGRES_PORT", "5432"),
}

LAST_FEED_CONTENT = None  # Stores last feed content for change detection

def connect_db():
    """Connect to PostgreSQL database."""
    return psycopg2.connect(**DB_CONFIG)

def fetch_rss_feed():
    """Fetch and parse the RSS feed."""
    global LAST_FEED_CONTENT
    print(f"Fetching RSS feed from {RSS_FEED_URL}...")
    feed = feedparser.parse(RSS_FEED_URL)

    if not feed.entries:
        print("No new entries found.")
        return []

    feed_content = [(entry.title, entry.link) for entry in feed.entries]

    # Check if content has changed
    if feed_content == LAST_FEED_CONTENT:
        print("Feed has not changed since the last fetch. Skipping...")
        return []

    LAST_FEED_CONTENT = feed_content  # Update last fetched content
    return feed.entries

def process_and_store_news(entries):
    """Process and store news articles in PostgreSQL."""
    conn = connect_db()
    cur = conn.cursor()

    for entry in entries:
        headline = entry.title
        url = entry.link
        summary = entry.get("summary", "")

        # Fetch the first image (if available)
        thumbnail_data = None
        if "media_content" in entry:
            image_url = entry.media_content[0]['url']
            try:
                response = requests.get(image_url)
                if response.status_code == 200:
                    img = Image.open(BytesIO(response.content))
                    img = img.convert("RGB")
                    img_io = BytesIO()
                    img.save(img_io, format="JPEG")
                    thumbnail_data = img_io.getvalue()
            except Exception as e:
                print(f"Could not fetch image: {e}")

        try:
            cur.execute("""
                INSERT INTO news_articles (headline, url, summary, thumbnail)
                VALUES (%s, %s, %s, %s)
                ON CONFLICT (url) DO NOTHING;
            """, (headline, url, summary, psycopg2.Binary(thumbnail_data) if thumbnail_data else None)
    )

            conn.commit()
            print(f"Inserted: {headline}")
        except psycopg2.IntegrityError:
            conn.rollback()
            print(f"Skipping duplicate: {headline}")

    cur.close()
    conn.close()
```

Figure 6: Code for creating database and fetching RSS for task2
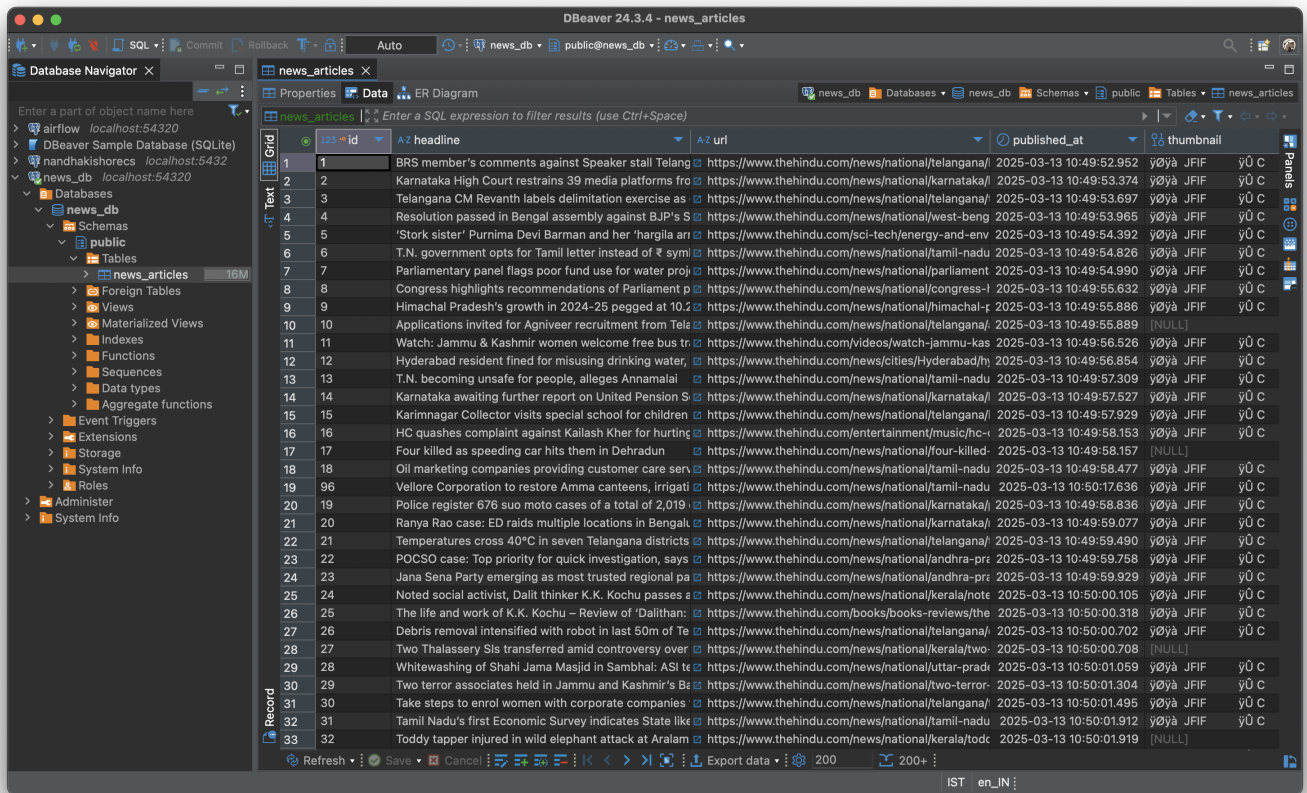
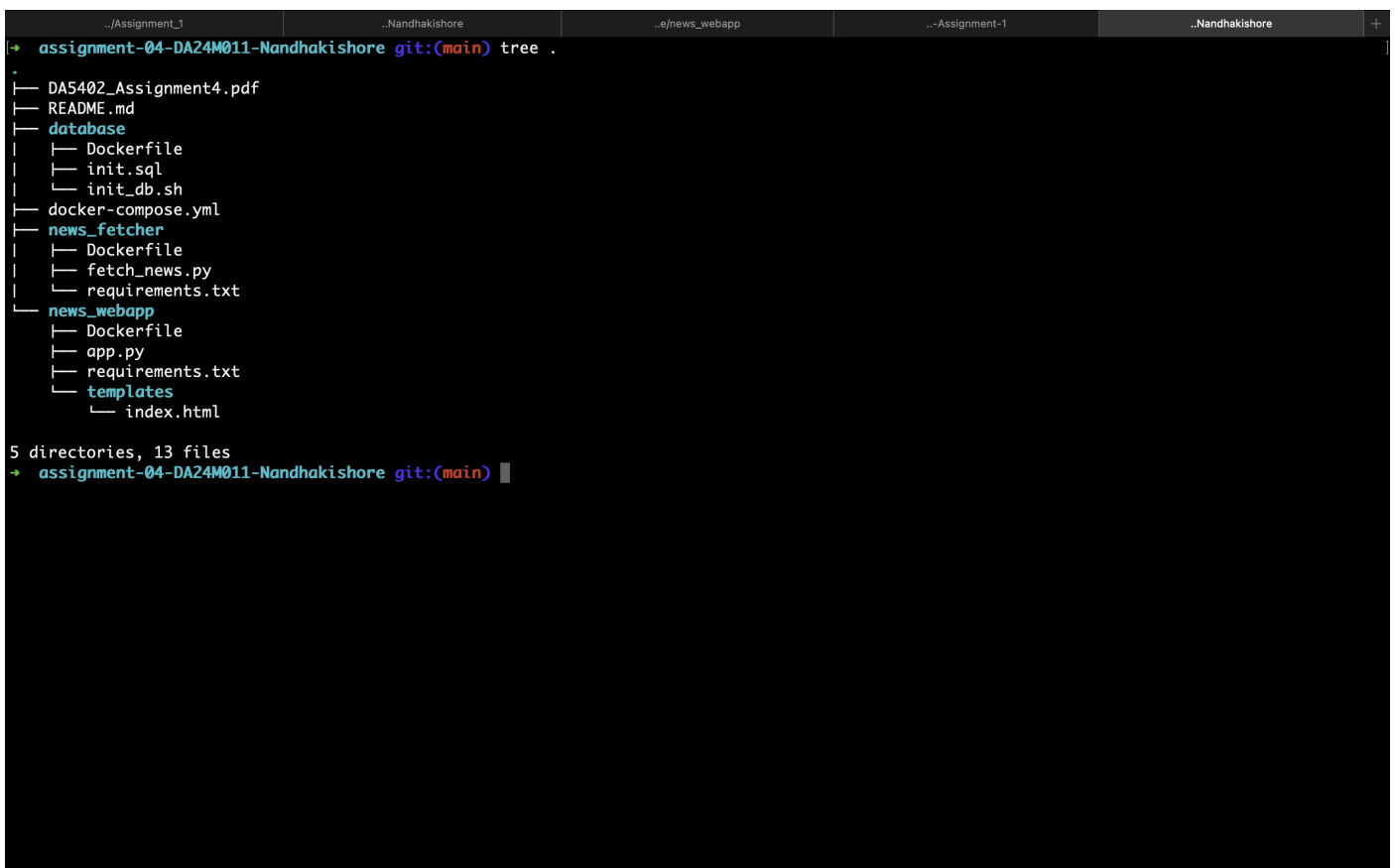Figure 7: Database viewed with Dbeaver



Figure 8: Directory structure for assignment 4

```yaml
1  services:
2    database:
3      build: ./database
4      container_name: news_postgres
5      restart: unless-stopped
6      environment:
7        POSTGRES_USER: news_admin
8        POSTGRES_PASSWORD: securepassword
9        POSTGRES_DB: news_db
10     ports:
11       - "54320:5432"
12     volumes:
13       - pgdata:/var/lib/postgresql/data
14
15   news_fetcher:
16     build: ./news_fetcher
17     container_name: news_fetcher
18     restart: unless-stopped
19     env_file:
20       - ./news_fetcher/.env
21     depends_on:
22       - database
23
24   news_webapp:
25     build: ./news_webapp
26     container_name: news_webapp
27     restart: unless-stopped
28     ports:
29       - "8000:8000"
30     env_file:
31       - ./news_webapp/.env
32     depends_on:
33       - database
34
35 volumes:
36   pgdata:
```

Figure 9: common docker-compose.yaml file

```
1  from fastapi import FastAPI, Request, Query
2  from fastapi.templating import Jinja2Templates
3  import psycopg2
4  import os
5  import base64
6  from datetime import date
7  from dotenv import load_dotenv
8
9  # Load environment variables
10 load_dotenv()
11
12 app = FastAPI()
13 templates = Jinja2Templates(directory="templates")
14
15 # Database connection
16 def get_db_connection():
17     return psycopg2.connect(
18         dbname=os.getenv("POSTGRES_DB"),
19         user=os.getenv("POSTGRES_USER"),
20         password=os.getenv("POSTGRES_PASSWORD"),
21         host=os.getenv("POSTGRES_HOST"),
22         port=os.getenv("POSTGRES_PORT")
23     )
24
25 @app.get("/")
26 def home(request: Request):
27     conn = get_db_connection()
28     cur = conn.cursor()
29     cur.execute("""
30         SELECT headline, url, summary, thumbnail
31         FROM news_articles
32         WHERE date(published_at) = CURRENT_DATE
33         ORDER BY published_at DESC
34     """)
35     articles = cur.fetchall()
36     conn.close()
37
38     # Format articles properly
39     formatted_articles = []
40     for article in articles:
41         headline, url, summary, thumbnail = article
42         if isinstance(thumbnail, memoryview):  # Convert binary images to Base64
43             thumbnail = f"data:image/png;base64,{base64.b64encode(thumbnail.tobytes()).decode()}"
44         formatted_articles.append((headline, url, summary, thumbnail))
45
46     # Pass your name and roll number to the template
47     return templates.TemplateResponse("index.html", {
48         "request": request,
49         "articles": formatted_articles,
50         "name": "Nandhakishore CS",  # Replace with your actual name
51         "roll_number": "DA24M011"    # Replace with your actual roll number
52     })
```
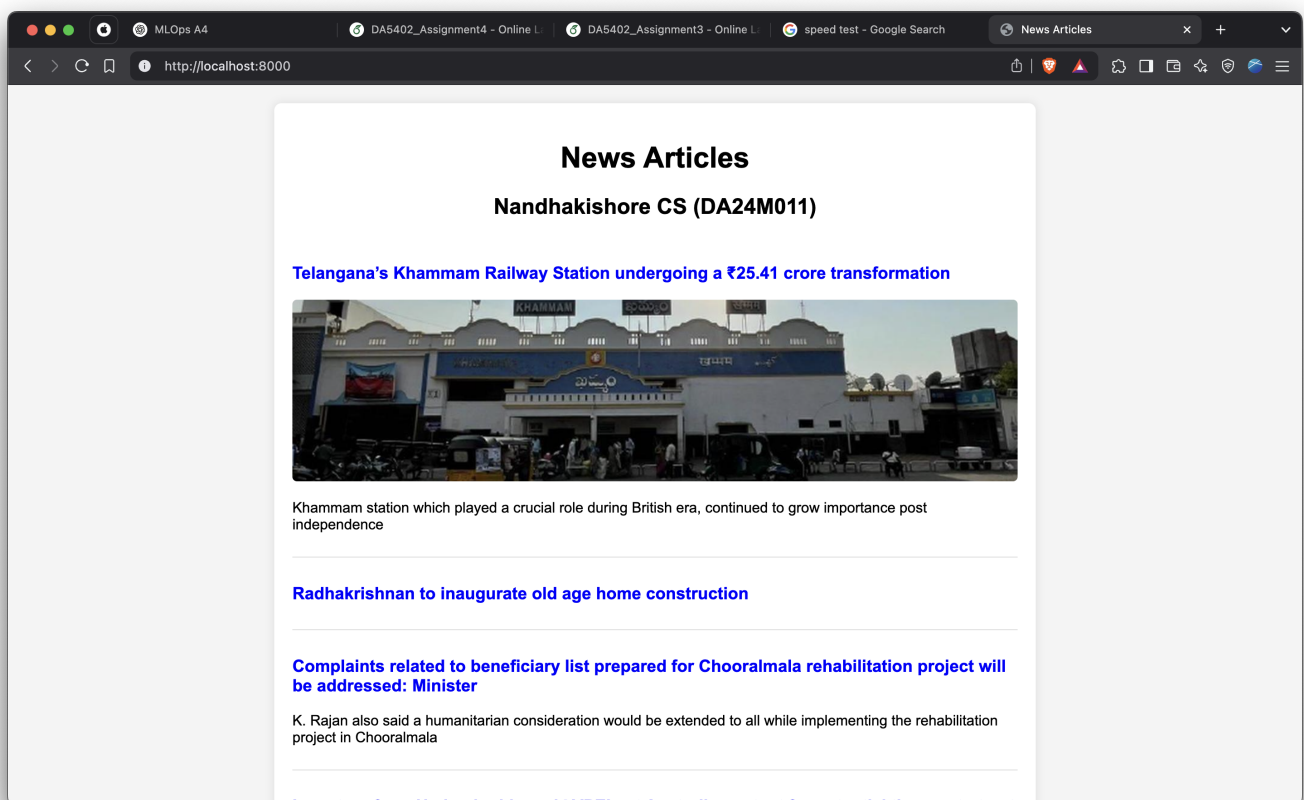
Figure 10: RESTAPI based webapp

Figure 11: Screenshot of the webapp