Santhosh Andavar

2818372

CIS 492(Big Data)

04/20/2025

Lab4 part 2

In the second part of the lab, we computed the cosine similarity between 10 selected documents from the State of the Union dataset using the TF-IDF vector space model.
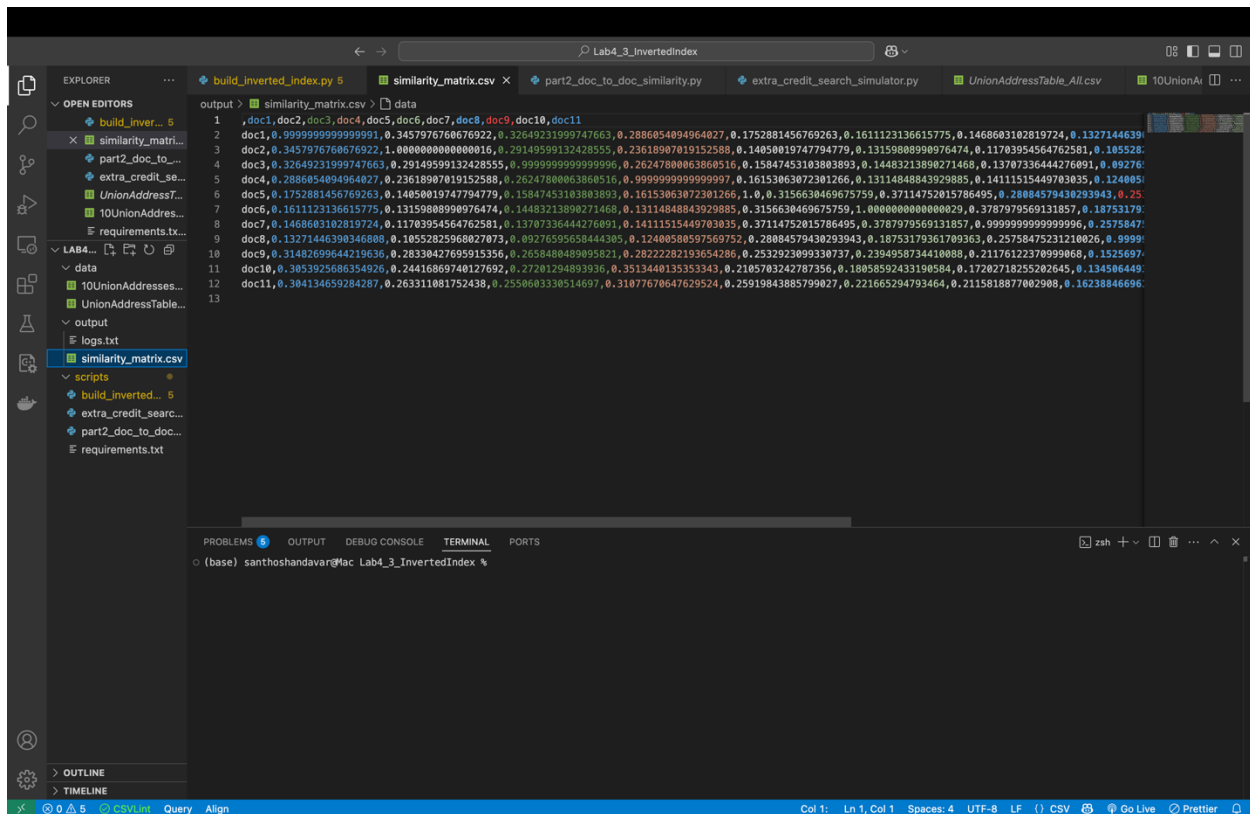
Steps:

- Used TfidfVectorizer from scikit-learn to compute document vectors
- Computed pairwise cosine similarity across documents
- Saved a 10x10 similarity matrix to CSV format

**Query:** `freedom peace economy`

<mark>**Top 5 Similar Documents:**</mark>

- doc5
- doc2
- doc6
- doc3
- doc7

Cosine similarity matrix showing pairwise document similarity across the 10 State of the Union addresses using TF-IDF vectors.

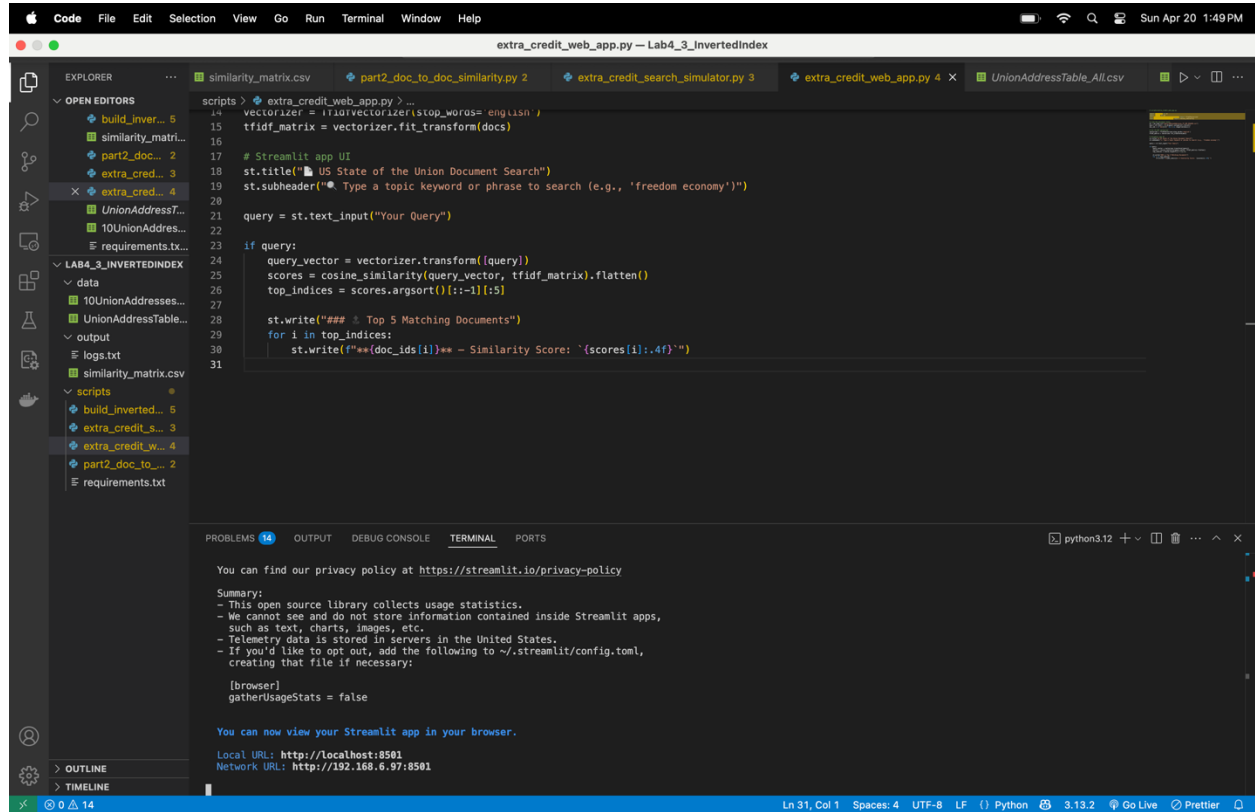## Extra Credit: Web-Based Search Engine with User Input

A Streamlit web application was developed to simulate a real-time search engine. Users can enter a query, and the system returns the top 5 most relevant documents using TF-IDF and cosine similarity.

- Input box captures keywords like "freedom peace economy"
- Matches query vector to document vectors
- Displays most similar speeches and scores

## Queries Used for Screenshots:

1. `freedom rights democracy`
2. `economy inflation growth`
3. `war peace treaty`
4. `education science innovation`
5. `tax revenue budget`

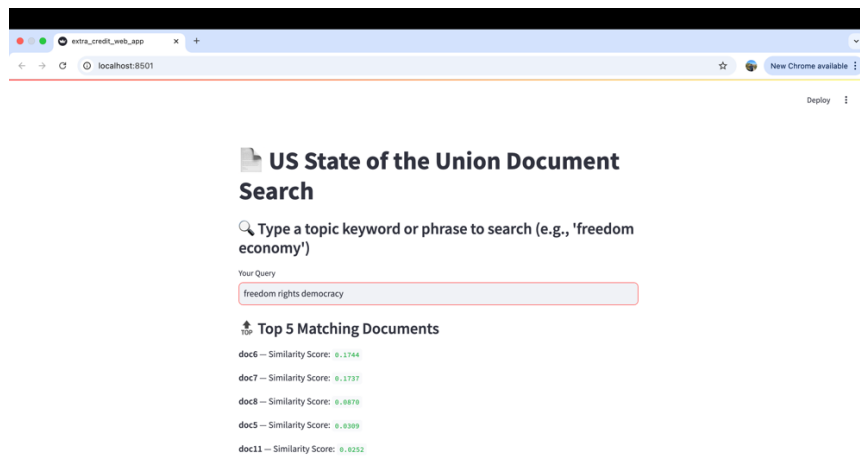Real-time Streamlit app matching user queries against speech documents using cosine similarity.



"Freedom rights democracy" — shows documents discussing civil liberties and governance.
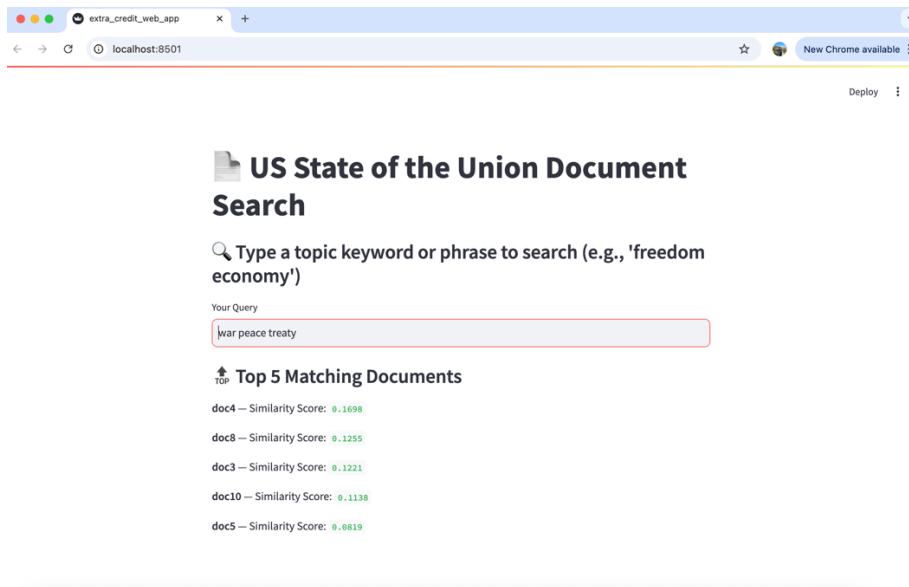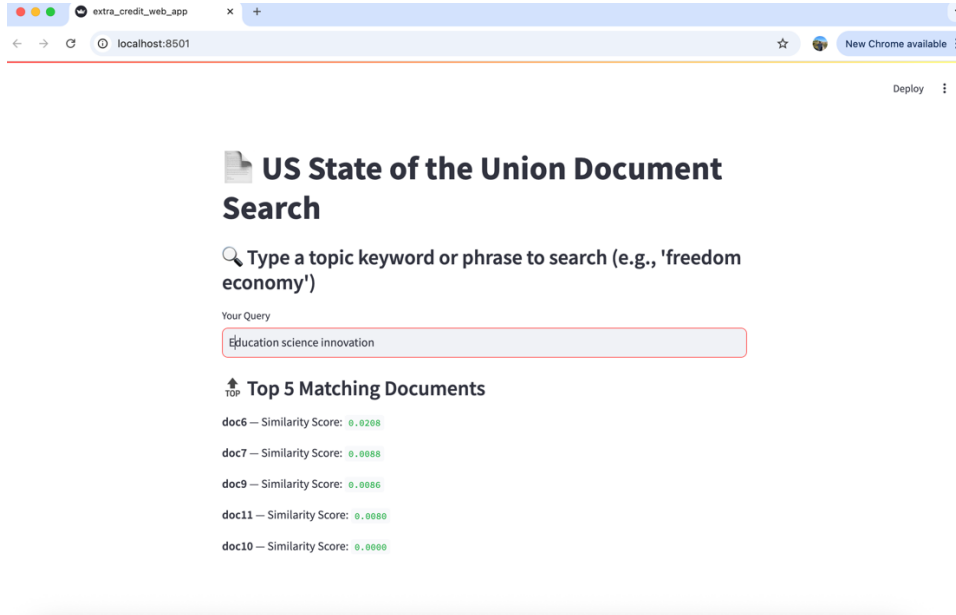
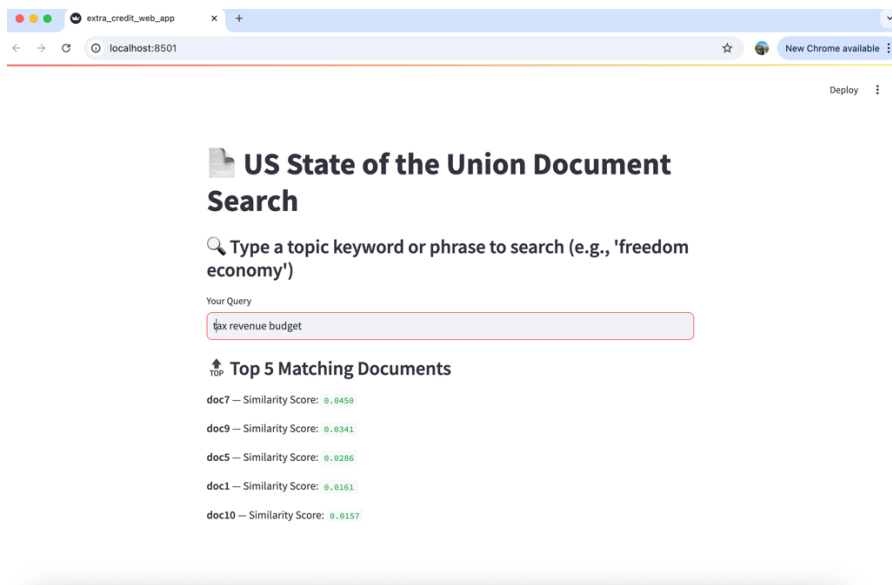"Economy inflation growth" — returns economic-focused speeches.



*"War peace treaty" — matches documents addressing defense and foreign relations.



"Education science innovation" — highlights educational and innovation content

# US State of the Union Document Search

🔍 Type a topic keyword or phrase to search (e.g., 'freedom economy')

Your Query

Education science innovation

### 🔝 Top 5 Matching Documents

**doc6** — Similarity Score: 0.0208

**doc7** — Similarity Score: 0.0088

**doc9** — Similarity Score: 0.0086

**doc11** — Similarity Score: 0.0080

**doc10** — Similarity Score: 0.0000

"Tax revenue budget" — captures fiscal discussions in presidential addresses.



# US State of the Union Document Search

🔍 Type a topic keyword or phrase to search (e.g., 'freedom economy')

Your Query

tax revenue budget

### 🔝 Top 5 Matching Documents

**doc7** — Similarity Score: 0.0450

**doc9** — Similarity Score: 0.0341

**doc5** — Similarity Score: 0.0286

**doc1** — Similarity Score: 0.0161

**doc10** — Similarity Score: 0.0157

**Conclusion:**

This lab successfully demonstrated building an inverted index and computing document similarity using TF-IDF and cosine similarity. The integration of NLP techniques and a web-based query system allowed us to explore real-time content-based search with extra credit features, making this a comprehensive document search pipeline.