

DSCI-560 Lab 5

<https://github.com/herunkan/Data-Science-Professional-Practicum/tree/main/Lab5>

Team Name: The Coach

Name	USC ID
Tianqi Qiu	1716906139
Shida Yan	5725964711
Herun Kan	7222919427

Initial Setup

This project was implemented in Python on a local macOS environment with MySQL used for persistent storage. The system was designed to perform web scraping, preprocessing, database storage, clustering, and automated periodic updates.

The primary libraries used in this project include:

- requests and beautifulsoup4 for web scraping
- pandas for data processing
- scikit-learn for TF-IDF embedding and K-Means clustering
- matplotlib for visualization
- mysql-connector-python for database interaction

Data Collection and Storage

Topic Selection

The initial target subreddit was r/Fitness. However, during preliminary data exploration, it was observed that r/Fitness contained a significant number of repetitive megathreads and lower topical consistency. As a result, the project was refocused on r/workout to obtain cleaner and more coherent content. All final analysis presented in this report is based on r/workout.

Scraping Method

Data was collected by crawling old.reddit.com using the BeautifulSoup library. Pagination was handled by repeatedly following the “next” links to retrieve additional posts, avoiding reliance on single-page limits. The scraping script accepts a command-line parameter specifying the number of posts to collect, allowing flexible experimentation with dataset size.

Database Storage

All collected posts were stored in a MySQL database named lab5_workout. The primary table, reddit_posts, stores post-level data, while a supporting table, cluster_keywords, stores representative keywords for each cluster.

An upsert mechanism was implemented to ensure that repeated runs update existing rows instead of inserting duplicate entries. This design supports automation and periodic updates without corrupting the dataset.

Data Preprocessing

A structured preprocessing pipeline was implemented to prepare the raw Reddit text for clustering. The steps include:

- Removal of HTML tags
- Text normalization and lowercasing
- Removal of URLs and special characters
- Username masking (author_masked) to preserve privacy
- Keyword extraction for each message
- Generation of topic hints

Additional filtering steps were applied to improve data quality:

- Promoted or sponsored posts were removed
- Non-comment permalink entries were filtered
- Low-signal or junk rows were reduced prior to clustering

These steps ensure that the input to the embedding and clustering stages is clean, consistent, and meaningful.

Forum Analysis and Clustering

Embedding (Message Abstraction)

TF-IDF vectorization was used to convert each cleaned message into a fixed-dimensional numerical vector. This transformation represents each post in a high-dimensional feature space based on word importance. The TF-IDF representation satisfies the lab

requirement of embedding or message abstraction prior to clustering.

Clustering

K-Means clustering was applied with K = 8 clusters. The clustering pipeline produces the following outputs:

- Cluster assignment for each message
- Distance from each message to its assigned centroid
- Posts nearest to each centroid
- A two-dimensional cluster visualization generated using PCA projection

Representative Cluster Keywords

From the most recent run, the clusters produced representative keyword sets such as:

Cluster 0: workout, routine, good, beginner, body, day, split, after

Cluster 1: can, should, anyone, workouts, only, does, training, leg

Cluster 2: gym, how, going, got, anyone, else, should, most

Cluster 3: how, guys, workout, why, exercise, think, actually, body

Cluster 4: day, split, how, can, routine, workout, help, any

Cluster 5: advice, need, help, body, gym, weight, stuck, split

Cluster 6: out, working, how, work, after, lost, help, back

Cluster 7: how, routine, weight, muscle, start, can, day, best

```

- Cluster 5: advice, need, help, body, gym, weight, stuck, split
- Cluster 6: out, working, how, work, after, lost, help, back
- Cluster 7: how, routine, weight, muscle, start, can, day, best
Update cycle completed.
Next update in 300 seconds.
[2026-02-14 17:24:16] Fetching/processing/updating DB...
[1/5] Initializing database...
[2/5] Fetching posts with BeautifulSoup scraper...
Fetched 1318 posts from r/workout.
[3/5] Preprocessing posts...
Upserter 1318 posts.
[4/5] Clustering messages...
[5/5] Done.
Saved clustered data: outputs/clustered_posts.csv
Saved centroid-near posts: outputs/posts_near_centroid.csv
Saved cluster plot: outputs/cluster_scatter.png

Top keywords per cluster:
- Cluster 0: workout, routine, good, beginner, body, day, split, after
- Cluster 1: can, should, anyone, workouts, only, does, training, leg
- Cluster 2: gym, how, going, got, anyone, else, should, most
- Cluster 3: guys, workout, why, exercise, think, actually, body
- Cluster 4: day, split, how, can, routine, working, help, any
- Cluster 5: advice, need, help, body, gym, weight, stuck, split
- Cluster 6: out, working, how, work, after, lost, help, back
- Cluster 7: how, routine, weight, muscle, start, can, day, best
Update cycle completed.
Next update in 300 seconds.
[2026-02-14 17:32:26] Fetching/processing/updating DB...
[1/5] Initializing database...
[2/5] Fetching posts with BeautifulSoup scraper...
Fetched 1318 posts from r/workout.
[3/5] Preprocessing posts...
Upserter 1318 posts.
[4/5] Clustering messages...
[5/5] Done.
Saved clustered data: outputs/clustered_posts.csv
Saved centroid-near posts: outputs/posts_near_centroid.csv
Saved cluster plot: outputs/cluster_scatter.png

Top keywords per cluster:
- Cluster 0: workout, routine, good, beginner, body, day, split, after
- Cluster 1: can, should, anyone, workouts, only, does, training, leg
- Cluster 2: gym, how, going, got, anyone, else, should, most
- Cluster 3: guys, workout, why, exercise, think, actually, body
- Cluster 4: day, split, how, can, routine, working, help, any
- Cluster 5: advice, need, help, body, gym, weight, stuck, split
- Cluster 6: out, working, how, work, after, lost, help, back
- Cluster 7: how, routine, weight, muscle, start, can, day, best
Update cycle completed.
Next update in 300 seconds.

```

These clusters reflect distinct workout-related themes.

Automation

A real-time periodic update system was implemented in `automation.py`. The script accepts an interval in minutes and executes the following cycle repeatedly:

1. Fetch new posts
2. Preprocess the text
3. Update the database
4. Recompute clustering
5. Wait for the specified interval

During the waiting period, the system accepts free-text queries from the terminal. Each query is embedded and mapped to the nearest cluster. The system then retrieves and displays the most

relevant posts from that cluster. Additionally, a query-specific visualization is generated and saved as a PNG file.

This automation component demonstrates a practical near real-time analysis workflow.

Results and Artifacts

Runtime Summary

In the latest run:

- 500 posts were requested
- 1,277 rows were stored and processed
- 8 clusters were generated

```
Next update in 300 seconds.
back day

Closest cluster: 1
Cluster keywords: can, should, anyone, workouts, only, does, training, leg, split, weight
Top related messages:
- (0) lower back pain is ruining leg day
  https://old.reddit.com/r/workout/comments/1rl2dw/lower_back_pain_is_ruining_leg_day/
- (0) What are your back day exercises?
  https://old.reddit.com/r/workout/comments/1qjygf/what_are_your_back_day_exercises/
- (0) Slight Back injury, again
  https://old.reddit.com/r/workout/comments/1r4f3d/slight_back_injury_again/
- (0) Gymming @day 1
  https://old.reddit.com/r/workout/comments/1qhzb5/gymming_day_1/
- (0) Lower back soreness
  https://old.reddit.com/r/workout/comments/1lqu0k/lower_back_soreness/
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:67: UserWarning: Glyph 128075 (\N{WAVING HAND SIGN}) missing from font(s) DejaVu Sans.
  plt.tight_layout()
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:67: UserWarning: Glyph 128064 (\N{EYES}) missing from font(s) DejaVu Sans.
  plt.tight_layout()
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:68: UserWarning: Glyph 128075 (\N{WAVING HAND SIGN}) missing from font(s) DejaVu Sans.
  plt.savefig(output_path, dpi=160)
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:68: UserWarning: Glyph 128064 (\N{EYES}) missing from font(s) DejaVu Sans.
  plt.savefig(output_path, dpi=160)
Saved cluster view: outputs/query_cluster_1.png

bench press

Closest cluster: 1
Cluster keywords: can, should, anyone, workouts, only, does, training, leg, split, weight
Top related messages:
- (0) Bench press
  https://old.reddit.com/r/workout/comments/1qxv14/bench_press/
- (0) Dumbbell bench press
  https://old.reddit.com/r/workout/comments/1qxw7d/dumbbell_bench_press/
- (0) Overcoming fear of barbell bench press?
  https://old.reddit.com/r/workout/comments/1r0mfp/overcoming_fear_of_barbell_bench_press/ what does that mean exactly?
- (0) If a workout plan says "bench press into incline bench press" what does that mean exactly?
  https://old.reddit.com/r/workout/comments/1rl0map/if_a_workout_plan_says_bench_press_into_incline/
- (0) Why use a bench without safeties for bench press if choice available?
  https://old.reddit.com/r/workout/comments/1lqatul/why_use_a_bench_without_safeties_for_bench_press/
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:67: UserWarning: Glyph 128075 (\N{WAVING HAND SIGN}) missing from font(s) DejaVu Sans.
  plt.tight_layout()
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:67: UserWarning: Glyph 128064 (\N{EYES}) missing from font(s) DejaVu Sans.
  plt.tight_layout()
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:68: UserWarning: Glyph 128075 (\N{WAVING HAND SIGN}) missing from font(s) DejaVu Sans.
  plt.savefig(output_path, dpi=160)
/Users/herunkan/GITFolder/DSCL 560/Lab5/automation.py:68: UserWarning: Glyph 128064 (\N{EYES}) missing from font(s) DejaVu Sans.
  plt.savefig(output_path, dpi=160)
Saved cluster view: outputs/query_cluster_1.png
```

Challenges and Decisions

The BeautifulSoup scraping approach was chosen to avoid Reddit API registration requirements.

The subreddit r/Fitness was initially selected but later replaced with r/workout due to excessive repetitive threads naming

Selecting the number of clusters ($K = 8$) required experimentation to balance interpretability and granularity.

Team Meeting Notes

Day 1 – 2/9/2026

- Reviewed the Lab 5 assignment requirements and clarified deliverables (scraping, preprocessing, clustering, automation, and team notes)
- Discussed scraping options (PRAW vs BeautifulSoup) and agreed to begin with PRAW as the initial approach
- Divided tasks so each team member could focus on scraper setup, database design, and clustering research

Day 2 – 2/10/2026

- Shared progress on Reddit scraping setup and encountered API approval/access limitations with the PRAW route

- Discussed assignment-allowed alternatives and agreed to switch to the BeautifulSoup-based scraping method
- Finalized use of MySQL for storage to align with lab instructions
- Outlined preprocessing requirements including cleaning text, masking usernames, and keyword extraction

Day 3 – 2/11/2026

- Implemented and reviewed the BeautifulSoup scraper with pagination to handle larger post requests
- Built MySQL schema and upsert logic for post storage and cluster metadata
- Integrated preprocessing pipeline and started initial clustering experiments using TF-IDF + KMeans
- Validated that representative posts could be retrieved near each cluster centroid

Day 4 – 2/12/2026

- Developed automation workflow to run scraping, preprocessing, storage, and clustering at fixed intervals
- Added interactive query mode to map user input to the closest cluster and display related messages
- Reviewed output quality and identified irrelevant/promoted/noisy posts affecting cluster quality
- Applied filtering and query-ranking improvements to return more relevant subreddit-specific results

Day 5 – 2/13/2026

- Completed end-to-end testing on r/workout and confirmed successful generation of cluster outputs and query visuals

- Captured terminal evidence/screenshots and verified output artifacts for submission
- Assigned final responsibilities for report formatting, meeting notes export, and repository cleanup