# UITS
## UNIVERSITY OF INFORMATION TECHNOLOGY & SCIENCES

**Submitted To :** Mrinmoy Biswas Akash
Lecturer & Course Coordinator
Department of CSE, UITS.

**Submitted By :** Name: Md. Mehadi Hasan
ID : 2125051003
Dept : CSE, Batch : 50, Section : 8A

**Course Code :**   CSE426
**Course Title :**   Data Mining & Warehouse LAB

**Date of Submission :** 05-07-2025.

# CSE
### University of Information Technology and Sciences

'

# GitHub Links

**Project-01:**

**https://github.com/Mehadi4021/CSE426_Data_Mining_and_Warehouse_Lab/blob/main/Assignment01_2125051003.ipynb**

**Project-02:**

https://github.com/Mehadi4021/CSE426_Data_Mining_and_Warehouse_Lab/blob/main/Project02_Discovering_Edibility_Patterns_in_Heart_Disease_using_Association_Rule_Mining.ipynb

**Project-03:**

https://github.com/Mehadi4021/CSE426_Data_Mining_and_Warehouse_Lab/blob/main/Project03_Building_a_Domain_Specific_Search_Engine_with_Crawling_and_Link_Analysis.ipynb

# Project-01

**Project Title:** Movie Recommendation System.

## Introduction :

Recommendation systems are essential for helping users navigate the vast amount of content available on digital platforms. With so many options to choose from, users can easily feel overwhelmed, making these systems incredibly valuable. In this project, we built a movie recommendation system using the MovieLens dataset, aiming to suggest films that match users' individual preferences. Our approach centers on collaborative filtering, a method that analyzes user ratings to uncover similarities between movies. By analyzing how users rate various films, the system identifies patterns and recommends movies that align with a user's tastes, even if they haven't seen those movies before. It operates by creating a movie similarity matrix and providing personalized suggestions based on each user's unique rating history.

## Methodology :

- **Dataset Description -**

  Two datasets from the MovieLens collection were used:

  1. **Ratings.csv -**
     - Contains user ratings for movies.
     - Columns: `userId`, `movieId`, `rating`, `timestamp`.

  2. **Movies.csv -**
     - Contains movie details.
     - Columns: `movieId`, `title`, `genres`.

  Both datasets were stored in Google Drive and accessed through Google Colab for seamless processing.

- **Movie Similarity Calculation -**

  To recommend movies aligned with a user's preferences, we constructed a similarity

matrix between movies. This matrix captures how closely two movies are related, based on the ratings provided by users. When multiple users rate two movies similarly, these movies are considered similar. Although the implementation employed Pearson correlation, the underlying concept mirrors cosine similarity, comparing rating vectors to assess the degree of closeness between movies.

● **Personalized Recommendations -**

The system generates personalized movie suggestions by first identifying each user's highest-rated films. It then finds movies similar to those favorites. To maintain relevance and avoid recommending already-viewed content, the system filters out any movies the user has previously rated. From the remaining pool, it prioritizes movies with higher average ratings from the broader user base, delivering quality recommendations tailored to the user's taste.

## How It Works :

1. **Movie Similarity Matrix Computation -**
   The system determines the relationship between movies by analyzing how users have rated them. It constructs a matrix where each element represents the similarity score between two movies. A higher similarity score indicates that the two films received similar rating patterns from users, suggesting they appeal to similar audiences.

2. **Generating Recommendations -**

   ○ A user is selected from the dataset.

   ○ The user's top-rated movies are identified.

   ○ For each favorite movie, similar films are retrieved based on the similarity matrix.

   ○ Movies the user has already rated are filtered out to avoid redundancy.

   ○ From the remaining choices, the system recommends the highest-rated movies, offering personalized suggestions tailored to the user's preferences.

3. **Personalization -**

By focusing on movies that resemble a user's favorite selections, the system delivers recommendations that feel personal and meaningful. Instead of simply suggesting popular films, it customizes the experience to match individual user tastes, enhancing satisfaction and engagement.

## Codes :

```python
# Import necessary libraries
import pandas as pd
import numpy as np
```

```python
# Mount Google Drive
from google.colab import drive
drive.mount('/content/drive')
```

```python
#Load Dataset
ratings =
pd.read_csv("/content/drive/MyDrive/CSE426_Data_Mining_and_Wareh
ouse_Lab/LAB_2/03 Recommendation System 2/ratings.csv")
movies =
pd.read_csv("/content/drive/MyDrive/CSE426_Data_Mining_and_Wareh
ouse_Lab/LAB_2/03 Recommendation System 2/movies.csv")
```

```python
#ratings
#movies
```

```python
# Display raw data
ratings.head()
movies.head()
```

```python
# Step 1: Create movie-to-movie similarity matrix
pivot_table = ratings.pivot_table(index='userId',
columns='movieId', values='rating') # Changed 'ratings_data' to
'ratings'
similarity_matrix = pivot_table.corr(method='pearson')
similarity_matrix.head()



# Step 2: Movie recommendation based on a given movie
def get_similar_movies(target_movie_id, num_recommendations=5):
    if target_movie_id not in similarity_matrix:
        return "Selected movie not found in the dataset."

    similarity_scores =
similarity_matrix[target_movie_id].dropna()
    top_matches =
similarity_scores.sort_values(ascending=False)[1:num_recommendat
ions+1]

    top_movies =
movies_data[movies_data["movieId"].isin(top_matches.index)][["mo
vieId", "title"]]
    return top_movies



# Example: Recommend movies similar to movieId = 2
top_recommendations = get_similar_movies(2,
num_recommendations=5)
top_recommendations

selected_user = int(input("Enter your user ID: "))
# Use 'ratings' instead of 'ratings_data'
user_rated = ratings[ratings['userId'] == selected_user]
user_rated.head()



# Step 3: Find the highest-rated movie by the user
fav_movie = user_rated.loc[user_rated['rating'].idxmax()]
fav_movie
```

```python
# Step 4: Identify movies not rated by the user

all_movie_ids = set(movies["movieId"])
rated_by_user = set(user_rated["movieId"])
not_rated_yet = all_movie_ids - rated_by_user
unseen_movies = movies[movies["movieId"].isin(not_rated_yet)]
unseen_movies.head()



# Step 5: Recommend movies not rated by the user, sorted by
average rating
def recommend_unseen_top_movies(user_id, num_movies=5):
    user_history = ratings[ratings["userId"] == user_id]
    movies_rated = set(user_history["movieId"])
    all_movies = set(movies["movieId"])
    movies_left = all_movies - movies_rated

    candidate_movies =
movies[movies["movieId"].isin(movies_left)]

    avg_movie_scores =
ratings.groupby("movieId")["rating"].mean()

    final_recommendations = candidate_movies.merge(
        avg_movie_scores, on="movieId", how="left"
    ).sort_values(by="rating", ascending=False).head(num_movies)

    return final_recommendations[["movieId", "title", "rating"]]



user_id = int(input("Enter your user ID: "))
final_suggestions = recommend_unseen_top_movies(user_id=user_id,
num_movies=10)
print(final_suggestions)
```

**Input's & Output's:**

```
In [ ]:  ratings
```

Out[ ]:

| | userId | movieId | rating | timestamp |
|---|---|---|---|---|
| **0** | 1 | 1 | 4.0 | 964982703 |
| **1** | 1 | 3 | 4.0 | 964981247 |
| **2** | 1 | 6 | 4.0 | 964982224 |
| **3** | 1 | 47 | 5.0 | 964983815 |
| **4** | 1 | 50 | 5.0 | 964982931 |
| **...** | ... | ... | ... | ... |
| **100831** | 610 | 166534 | 4.0 | 1493848402 |
| **100832** | 610 | 168248 | 5.0 | 1493850091 |
| **100833** | 610 | 168250 | 5.0 | 1494273047 |
| **100834** | 610 | 168252 | 5.0 | 1493846352 |
| **100835** | 610 | 170875 | 3.0 | 1493846415 |

100836 rows × 4 columns

```
In [ ]:  movies
```

Out[ ]:

| | movieId | title | genres |
|---|---|---|---|
| **0** | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy |
| **1** | 2 | Jumanji (1995) | Adventure\|Children\|Fantasy |
| **2** | 3 | Grumpier Old Men (1995) | Comedy\|Romance |
| **3** | 4 | Waiting to Exhale (1995) | Comedy\|Drama\|Romance |
| **4** | 5 | Father of the Bride Part II (1995) | Comedy |
| **...** | ... | ... | ... |
| **9737** | 193581 | Black Butler: Book of the Atlantic (2017) | Action\|Animation\|Comedy\|Fantasy |
| **9738** | 193583 | No Game No Life: Zero (2017) | Animation\|Comedy\|Fantasy |
| **9739** | 193585 | Flint (2017) | Drama |
| **9740** | 193587 | Bungo Stray Dogs: Dead Apple (2018) | Action\|Animation |
| **9741** | 193609 | Andrew Dice Clay: Dice Rules (1991) | Comedy |

9742 rows × 3 columns

```
In [ ]:  # Display raw data
         ratings.head()
         movies.head()
```

Out[ ]:

| | movieId | title | genres |
|---|---|---|---|
| **0** | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy |
| **1** | 2 | Jumanji (1995) | Adventure\|Children\|Fantasy |
| **2** | 3 | Grumpier Old Men (1995) | Comedy\|Romance |
| **3** | 4 | Waiting to Exhale (1995) | Comedy\|Drama\|Romance |
| **4** | 5 | Father of the Bride Part II (1995) | Comedy |

```
In [ ]:  # Step 1: Create movie-to-movie similarity matrix
         pivot_table = ratings.pivot_table(index='userId', columns='movieId', values='rating') # Changed 'ratings_data' to 'ratings
         similarity_matrix = pivot_table.corr(method='pearson')
         similarity_matrix.head()
```

Out[ ]:

| movieId | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... | 193565 | 193567 | 1935... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **movieId** | | | | | | | | | | | | | | |
| **1** | 1.000000 | 0.330978 | 0.487109 | 1.000000 | 0.310971 | 0.106465 | 0.208402 | 0.968246 | 0.095913 | -0.021409 | ... | NaN | NaN | Na |
| **2** | 0.330978 | 1.000000 | 0.419564 | NaN | 0.562791 | 0.163510 | 0.430261 | 0.415227 | 0.277350 | 0.016626 | ... | NaN | NaN | Na |
| **3** | 0.487109 | 0.419564 | 1.000000 | NaN | 0.602266 | 0.345069 | 0.554088 | 0.333333 | 0.458591 | -0.050276 | ... | NaN | NaN | Na |
| **4** | 1.000000 | NaN | NaN | 1.000000 | 0.654654 | NaN | 0.203653 | NaN | NaN | 0.870388 | ... | NaN | NaN | Na |
| **5** | 0.310971 | 0.562791 | 0.602266 | 0.654654 | 1.000000 | 0.291302 | 0.609119 | 0.555556 | 0.319173 | 0.218263 | ... | NaN | NaN | Na |

5 rows × 9724 columns

```
In [ ]:   # Step 2: Movie recommendation based on a given movie
          def get_similar_movies(target_movie_id, num_recommendations=5):
              if target_movie_id not in similarity_matrix:
                  return "Selected movie not found in the dataset."

              similarity_scores = similarity_matrix[target_movie_id].dropna()
              top_matches = similarity_scores.sort_values(ascending=False)[1:num_recommendations+1]

              top_movies = movies_data[movies_data["movieId"].isin(top_matches.index)][["movieId", "title"]]
              return top_movies
```

```
In [ ]:   # Example: Recommend movies similar to movieId = 2
          top_recommendations = get_similar_movies(2, num_recommendations=5)
          top_recommendations
```

Out[ ]:

| | movieId | title |
|---|---|---|
| 2311 | 3063 | Poison Ivy (1992) |
| 2825 | 3774 | House Party 2 (1991) |
| 2826 | 3783 | Croupier (1998) |
| 3583 | 4912 | Funny Girl (1968) |
| 3856 | 5420 | Windtalkers (2002) |

```
In [ ]:   selected_user = int(input("Enter your user ID: "))
          # Use 'ratings' instead of 'ratings_data'
          user_rated = ratings[ratings['userId'] == selected_user]
          user_rated.head()
```

Enter your user ID: 03

Out[ ]:

| | userId | movieId | rating | timestamp |
|---|---|---|---|---|
| 261 | 3 | 31 | 0.5 | 1306463578 |
| 262 | 3 | 527 | 0.5 | 1306464275 |
| 263 | 3 | 647 | 0.5 | 1306463619 |
| 264 | 3 | 688 | 0.5 | 1306464228 |
| 265 | 3 | 720 | 0.5 | 1306463595 |

```
In [ ]:   # Step 3: Find the highest-rated movie by the user
          fav_movie = user_rated.loc[user_rated['rating'].idxmax()]
          fav_movie
```

Out[ ]:

| | 266 |
|---|---|
| userId | 3.000000e+00 |
| movieId | 8.490000e+02 |
| rating | 5.000000e+00 |
| timestamp | 1.306464e+09 |

dtype: float64

```
In [ ]:   # Step 4: Identify movies not rated by the user

          all_movie_ids = set(movies["movieId"])
          rated_by_user = set(user_rated["movieId"])
          not_rated_yet = all_movie_ids - rated_by_user
          unseen_movies = movies[movies["movieId"].isin(not_rated_yet)]
          unseen_movies.head()
```

Out[ ]:

| | movieId | title | genres |
|---|---|---|---|
| **0** | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy |
| **1** | 2 | Jumanji (1995) | Adventure\|Children\|Fantasy |
| **2** | 3 | Grumpier Old Men (1995) | Comedy\|Romance |
| **3** | 4 | Waiting to Exhale (1995) | Comedy\|Drama\|Romance |
| **4** | 5 | Father of the Bride Part II (1995) | Comedy |

```
In [ ]:   # Step 5: Recommend movies not rated by the user, sorted by average rating
          def recommend_unseen_top_movies(user_id, num_movies=5):
              user_history = ratings[ratings["userId"] == user_id]
              movies_rated = set(user_history["movieId"])
              all_movies = set(movies["movieId"])
              movies_left = all_movies - movies_rated

              candidate_movies = movies[movies["movieId"].isin(movies_left)]

              avg_movie_scores = ratings.groupby("movieId")["rating"].mean()

              final_recommendations = candidate_movies.merge(
                  avg_movie_scores, on="movieId", how="left"
              ).sort_values(by="rating", ascending=False).head(num_movies)

              return final_recommendations[["movieId", "title", "rating"]]
```

```
In [ ]:   user_id = int(input("Enter your user ID: "))
          final_suggestions = recommend_unseen_top_movies(user_id=user_id, num_movies=10)
          print(final_suggestions)
```

```
Enter your user ID: 03
        movieId                                              title  rating
2855       3851                    I'm the One That I Want (2000)     5.0
2921       3951                          Two Family House (2000)     5.0
9672     187717               Won't You Be My Neighbor? (2018)     5.0
2914       3942                Sorority House Massacre II (1990)     5.0
9649     184245                        De platte jungle (1978)     5.0
2456       3303  Black Tar Heroin: The Dark End of the Street (...     5.0
3084       4180                        Reform School Girls (1986)     5.0
9549     175431                     Bobik Visiting Barbos (1977)     5.0
9547     175397          In the blue sea, in the white foam. (1984)     5.0
9546     175387    On the Trail of the Bremen Town Musicians (1973)     5.0
```

## Results:

- The system successfully generated a movie similarity matrix by utilizing collaborative filtering on user ratings, enabling it to discover films with strong patterns of shared viewer preferences.

- Personalized recommendations were created by pinpointing each user's top-rated movies and proposing similar films they have not yet rated. This approach guarantees that the suggestions are relevant and closely match each user's individual preferences.

- Example:
  For **User ID 03**, the system generated the following top 10 movie recommendations that the user had not previously rated. These suggestions are based on movies that are highly rated by users with similar preferences:

```
Enter your user ID: 03
      movieId                                              title  rating
2855     3851                        I'm the One That I Want (2000)     5.0
2921     3951                               Two Family House (2000)     5.0
9672   187717                      Won't You Be My Neighbor? (2018)     5.0
2914     3942                     Sorority House Massacre II (1990)     5.0
9649   184245                             De platte jungle (1978)     5.0
2456     3303  Black Tar Heroin: The Dark End of the Street (...     5.0
3084     4180                            Reform School Girls (1986)     5.0
9549   175431                          Bobik Visiting Barbos (1977)     5.0
9547   175397          In the blue sea, in the white foam. (1984)     5.0
9546   175387  On the Trail of the Bremen Town Musicians (1973)     5.0
```

  These recommendations demonstrate the system's ability to align movie suggestions with a user's likely interests, even for less mainstream or internationally produced films.

**Tools & Technologies :**

- Python 3
- Google Colab
- Pandas, NumPy (for data handling)
- Pearson Correlation (for similarity)

# Future Improvements :

- Integrate genre information to refine recommendations based on users' favorite movie categories.
- Implement user-user collaborative filtering to suggest movies liked by users with similar tastes.
- Combine collaborative filtering with content-based methods for more accurate suggestions.
- Explore advanced models like Neural Collaborative Filtering (NCF) or Autoencoders for more intelligent and dynamic recommendations.

# Conclusion :

This project successfully developed a personalized and functional movie recommendation system using the MovieLens dataset. By implementing collaborative filtering techniques and creating a similarity matrix between movies, the system provides movie recommendations tailored to individual user tastes.

Recommendations are made by analyzing a user's highest-rated movies and suggesting similar, unseen titles that also boast strong average ratings. This method improves the user experience by reducing the effort required to find suitable content and increasing the relevance of recommendations.

**Project-02**

**Project Title:** Discovering Edibility Patterns in Heart Disease using Association Rule Mining.

## Introduction:

This project on Association Rule Mining was conducted on a Heart Disease dataset to uncover frequently co-occurring health conditions and symptoms. By employing the Apriori algorithm, the study identifies frequent itemsets and extracts meaningful association rules. These findings reveal key contributing factors to heart disease, providing valuable insights for early detection, risk assessment, and preventive healthcare strategies.

## Objective:

1. To preprocess the Heart Disease dataset to prepare it for rule mining.

2. To apply the Apriori algorithm with a minimum support threshold of 0.3 for generating frequent itemsets.

3. To generate the top 10 association rules, ranked by confidence and lift, ensuring a minimum confidence of 0.7.

4. To analyze and interpret at least one of the discovered rules, explaining its significance in understanding the risk factors associated with heart disease.

## Methodology:

**Data Loading and Exploration**

- The heart disease dataset was imported using `pandas`.

- Initial exploration was done to understand the structure, attribute types, and missing values.

**Data Preprocessing**

- All categorical variables were one-hot encoded using `pd.get_dummies()` to convert them into a binary format suitable for Association Rule Mining.

- The resulting dataframe consisted of 0s and 1s, indicating the presence or absence of specific attribute values.

**Frequent Itemset Generation**

- The Apriori algorithm from `mlxtend.frequent_patterns` was applied with a minimum support of 0.3.

- This step identified item combinations that frequently appear together in the dataset.

**Association Rule Mining**

- Using the generated frequent itemsets, association rules were extracted with:

  - Confidence ≥ 0.7

  - Rules were ranked by Lift and Confidence

- The top 10 rules were selected for analysis.

**Rule Interpretation**

- One of the high-confidence rules was selected and analyzed in detail to explain the relationship between medical features and heart disease risk.

## Codes:

```python
import numpy as np
import pandas as pd
from mlxtend.frequent_patterns import apriori, association_rules

from google.colab import drive
drive.mount('/content/drive')

df =
pd.read_csv('/content/drive/MyDrive/CSE426_Data_Mining_and_Warehouse_Lab/Final_
Project/heart_disease.csv')
df



import pandas as pd

df_encoded = pd.get_dummies(df, columns=df.columns)
df_encoded

# Frequent itemset generation

frequent_itemsets = apriori(df_encoded, min_support=0.3, use_colnames=True)
print(frequent_itemsets)



# Association rule generation
num_itemsets = len(frequent_itemsets)
rules = association_rules(frequent_itemsets, metric="confidence",
min_threshold=0.7, num_itemsets = num_itemsets)
rules



rules = rules.sort_values(by=['confidence', 'lift'], ascending=False)
rules

# Step 3: Sort rules by lift in descending order and select top 10
top_10_rules = rules.sort_values(by=["confidence", "lift"],
ascending=False).head(10)

# Step 4: Display the top 10 rules
print(top_10_rules[['confidence', 'lift']])
```

# Input's & Output's:

```
In [ ]:  df = pd.read_csv('/content/drive/MyDrive/CSE426_Data_Mining_and_Warehouse_Lab/Final_Project/heart_disease.csv')
         df
```

Out[ ]:

| | age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | target |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | high | 1 | 3 | high | medium | 1 | 0 | medium | 0 | high | 0 | 0 | 1 | 1 |
| 1 | low | 1 | 2 | medium | medium | 0 | 1 | high | 0 | high | 0 | 0 | 2 | 1 |
| 2 | low | 0 | 1 | medium | low | 0 | 0 | high | 0 | high | 2 | 0 | 2 | 1 |
| 3 | medium | 1 | 1 | low | medium | 0 | 1 | high | 0 | medium | 2 | 0 | 2 | 1 |
| 4 | medium | 0 | 0 | low | high | 0 | 1 | high | 1 | medium | 2 | 0 | 2 | 1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 298 | medium | 0 | 0 | high | medium | 0 | 1 | low | 1 | medium | 1 | 0 | 3 | 0 |
| 299 | low | 1 | 3 | low | high | 0 | 1 | low | 0 | medium | 1 | 0 | 3 | 0 |
| 300 | high | 1 | 0 | high | low | 1 | 1 | low | 0 | high | 1 | 2 | 3 | 0 |
| 301 | medium | 1 | 0 | medium | low | 0 | 1 | low | 1 | medium | 1 | 1 | 3 | 0 |
| 302 | medium | 0 | 1 | medium | medium | 0 | 0 | high | 0 | low | 1 | 1 | 2 | 0 |

303 rows × 14 columns

```
In [ ]:  import pandas as pd

         df_encoded = pd.get_dummies(df, columns=df.columns)
         df_encoded
```

Out[ ]:

| | age_high | age_low | age_medium | sex_0 | sex_1 | cp_0 | cp_1 | cp_2 | cp_3 | trestbps_high | ... | ca_1 | ca_2 | ca_3 | ca_4 | thal_0 | thal_1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | True | False | False | False | True | False | False | False | True | True | ... | False | False | False | False | False | True |
| 1 | False | True | False | False | True | False | False | True | False | False | ... | False | False | False | False | False | False |
| 2 | False | True | False | True | False | False | True | False | False | False | ... | False | False | False | False | False | False |
| 3 | False | False | True | False | True | False | True | False | False | False | ... | False | False | False | False | False | False |
| 4 | False | False | True | True | False | True | False | False | False | False | ... | False | False | False | False | False | False |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 298 | False | False | True | True | False | True | False | False | False | True | ... | False | False | False | False | False | False |
| 299 | False | True | False | False | True | False | False | False | True | False | ... | False | False | False | False | False | False |
| 300 | True | False | False | False | True | True | False | False | False | True | ... | False | True | False | False | False | False |
| 301 | False | False | True | False | True | True | False | False | False | False | ... | True | False | False | False | False | False |
| 302 | False | False | True | True | False | False | True | False | False | False | ... | True | False | False | False | False | False |

303 rows × 42 columns

```
# Frequent itemset generation

frequent_itemsets = apriori(df_encoded, min_support=0.3, use_colnames=True)
print(frequent_itemsets)
```

```
      support                       itemsets
0    0.343234                     (age_high)
1    0.313531                      (age_low)
2    0.343234                   (age_medium)
3    0.316832                        (sex_0)
4    0.683168                        (sex_1)
..        ...                            ...
88   0.376238      (exang_0, thal_2, target_1)
89   0.336634          (ca_0, thal_2, target_1)
90   0.323432   (ca_0, exang_0, fbs_0, target_1)
91   0.330033   (exang_0, thal_2, fbs_0, target_1)
92   0.303630      (ca_0, thal_2, fbs_0, target_1)

[93 rows x 2 columns]
```

```
# Association rule generation
num_itemsets = len(frequent_itemsets)
rules = association_rules(frequent_itemsets, metric="confidence", min_threshold=0.7, num_itemsets = num_itemsets)
rules
```

| | antecedents | consequents | antecedent support | consequent support | support | confidence | lift | representativity | leverage | conviction | zhangs_r |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | (cp_0) | (sex_1) | 0.471947 | 0.683168 | 0.343234 | 0.727273 | 1.064559 | 1.0 | 0.020815 | 1.161716 | 0.1 |
| 1 | (sex_1) | (fbs_0) | 0.683168 | 0.851485 | 0.574257 | 0.840580 | 0.987192 | 1.0 | -0.007450 | 0.931593 | -0.0 |
| 2 | (restecg_0) | (sex_1) | 0.485149 | 0.683168 | 0.339934 | 0.700680 | 1.025633 | 1.0 | 0.008496 | 1.058506 | 0.0 |
| 3 | (thal_3) | (sex_1) | 0.386139 | 0.683168 | 0.336634 | 0.871795 | 1.276106 | 1.0 | 0.072836 | 2.471287 | 0.3 |
| 4 | (target_0) | (sex_1) | 0.455446 | 0.683168 | 0.376238 | 0.826087 | 1.209200 | 1.0 | 0.065092 | 1.821782 | 0.3 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 94 | (ca_0, fbs_0, target_1) | (thal_2) | 0.379538 | 0.547855 | 0.303630 | 0.800000 | 1.460241 | 1.0 | 0.095699 | 2.260726 | 0.5 |
| 95 | (target_1, fbs_0, thal_2) | (ca_0) | 0.376238 | 0.577558 | 0.303630 | 0.807018 | 1.397293 | 1.0 | 0.086331 | 2.189019 | 0.4 |
| 96 | (ca_0, thal_2) | (fbs_0, target_1) | 0.376238 | 0.468647 | 0.303630 | 0.807018 | 1.722016 | 1.0 | 0.127308 | 2.753375 | 0.6 |
| 97 | (ca_0, target_1) | (fbs_0, thal_2) | 0.429043 | 0.481848 | 0.303630 | 0.707692 | 1.468704 | 1.0 | 0.096897 | 1.772625 | 0.5 |
| 98 | (target_1, thal_2) | (ca_0, fbs_0) | 0.429043 | 0.511551 | 0.303630 | 0.707692 | 1.383424 | 1.0 | 0.084153 | 1.671009 | 0.4 |

99 rows × 14 columns

```python
rules = rules.sort_values(by=['confidence', 'lift'], ascending=False)
rules
```

| | antecedents | consequents | antecedent support | consequent support | support | confidence | lift | representativity | leverage | conviction | zhangs_ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 8 | (trestbps_low) | (fbs_0) | 0.333333 | 0.851485 | 0.313531 | 0.940594 | 1.104651 | 1.0 | 0.029703 | 2.500000 | 0.1 |
| 54 | (ca_0, thal_2) | (fbs_0) | 0.376238 | 0.851485 | 0.343234 | 0.912281 | 1.071399 | 1.0 | 0.022874 | 1.693069 | 0.1 |
| 93 | (ca_0, target_1, thal_2) | (fbs_0) | 0.336634 | 0.851485 | 0.303630 | 0.901961 | 1.059280 | 1.0 | 0.016992 | 1.514851 | 0.0 |
| 75 | (ca_0, thal_2) | (target_1) | 0.376238 | 0.544554 | 0.336634 | 0.894737 | 1.643062 | 1.0 | 0.131752 | 4.326733 | 0.6 |
| 44 | (ca_0, exang_0) | (fbs_0) | 0.432343 | 0.851485 | 0.386139 | 0.893130 | 1.048908 | 1.0 | 0.018005 | 1.389675 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 40 | (restecg_1, fbs_0) | (exang_0) | 0.438944 | 0.673267 | 0.310231 | 0.706767 | 1.049757 | 1.0 | 0.014704 | 1.114242 | 0.0 |
| 91 | (fbs_0, target_1) | (exang_0, thal_2) | 0.468647 | 0.445545 | 0.330033 | 0.704225 | 1.580595 | 1.0 | 0.121230 | 1.874587 | 0.6 |
| 89 | (exang_0, target_1) | (fbs_0, thal_2) | 0.468647 | 0.481848 | 0.330033 | 0.704225 | 1.461509 | 1.0 | 0.104216 | 1.751847 | 0.5 |
| 43 | (slope_2) | (exang_0, fbs_0) | 0.468647 | 0.577558 | 0.330033 | 0.704225 | 1.219316 | 1.0 | 0.059362 | 1.428257 | 0.3 |
| 2 | (restecg_0) | (sex_1) | 0.485149 | 0.683168 | 0.339934 | 0.700680 | 1.025633 | 1.0 | 0.008496 | 1.058506 | 0.0 |

99 rows × 14 columns

```python
# Step 3: Sort rules by lift in descending order and select top 10
top_10_rules = rules.sort_values(by=["confidence", "lift"], ascending=False).head(10)

# Step 4: Display the top 10 rules
print(top_10_rules[['confidence', 'lift']])
```

```
    confidence      lift
8     0.940594  1.104651
54    0.912281  1.071399
93    0.901961  1.059280
75    0.894737  1.643062
44    0.893130  1.048908
63    0.892157  1.325115
15    0.885714  1.040199
57    0.884615  1.038909
92    0.884615  1.624476
16    0.879518  1.032922
```

## Results:

The Heart Disease dataset was successfully preprocessed and transformed using one-hot encoding to prepare it for Association Rule Mining. The **Apriori algorithm** was applied with a **minimum support threshold of 0.3**, resulting in several frequent itemsets.

Subsequently, **association rules** were generated using a **confidence threshold of 0.7**. The rules were then sorted based on **confidence** and **lift**, and the **top 10 rules** were extracted.

Among these, the most significant rules indicated strong correlations between specific medical attributes. For example, attributes like:

- **Chest pain type (cp_0 or cp_1)**,

- **Exercise-induced angina (exang_0)**,

- **ST depression (oldpeak_0.0)**,

were frequently associated with the absence or presence of heart disease.

These results reveal clear patterns in how combinations of symptoms and test results correlate with heart disease diagnosis. Such insights can help medical professionals identify at-risk patients earlier and improve preventive healthcare decisions.

## Future Improvements:

1. Use a larger and more diverse dataset to improve the reliability of results.

2. Apply feature selection or dimensionality reduction to focus on key attributes.

3. Compare Apriori with other algorithms like FP-Growth for better performance.

4. Develop interactive visualizations to make rules easier to interpret.

## Conclusion:

This project successfully applied Association Rule Mining on a heart disease dataset to uncover hidden patterns and relationships among medical attributes. By using the Apriori algorithm, the analysis revealed significant rules with high confidence and lift, identifying key factors commonly associated with heart disease. These insights can support early diagnosis, risk assessment, and preventive healthcare decisions. Overall, the study demonstrates how data mining techniques can be effectively used to extract valuable knowledge from medical datasets.

**Project-03**

**Project Title:** Building a Domain-Specific Search Engine with Crawling and Link Analysis

## Introduction:

In this project, I developed a domain-specific search engine focused on cricket by utilizing multiple authoritative and relevant sources. To collect the data, I used web crawling techniques on several cricket news websites including ESPN Cricinfo, Cricbuzz, BBC Sport Cricket, ICC official site, Hindustan Times Cricket, Indian Express Sports, and others. These links provided a rich set of cricket-related content, which allowed me to perform focused crawling and link analysis. This approach ensured that the search engine retrieved high-quality, topic-specific information related to cricket matches, players, records, and news updates.

## Objective:

1. To build a domain-specific search engine that focuses exclusively on cricket-related content from selected, reliable news sources.

2. To implement focused web crawling techniques for collecting relevant data from cricket news websites while avoiding unrelated content.

3. To perform link analysis (PageRank) in order to rank and evaluate the importance of web pages within the cricket domain.

4. To enhance information retrieval accuracy by indexing only cricket-specific pages and ensuring more relevant search results for users interested in cricket.

## Methodology:

**Selection of Domain and Seed URLs:** Identified the specific domain (cricket) and selected relevant seed links from trusted cricket news websites like ESPN Cricinfo, Cricbuzz, ICC, etc.

**Focused Web Crawling:** Developed a crawler to fetch only domain-specific (cricket-related) content by filtering out irrelevant pages based on keywords or content analysis.

**Data Extraction and Preprocessing:** Extracted useful information such as article titles, summaries, links, and metadata from the crawled web pages, and cleaned the data for consistency.

**Link Graph Construction:** Built a directed graph representing the hyperlinks between crawled web pages to capture the link structure within the domain.

**PageRank Implementation:** Applied the PageRank algorithm on the constructed link graph to evaluate and rank the importance of each web page.

**Search Engine Indexing:** Indexed the processed data to enable efficient querying, allowing users to retrieve the most relevant and high-ranked cricket-related pages.

**Result Evaluation:** Evaluated the performance of the search engine based on relevance, precision, and the effectiveness of the PageRank results within the cricket domain.

## Codes:

```python
import requests
from bs4 import BeautifulSoup


import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords


STOPWORDS = stopwords.words('english')
print(STOPWORDS)
```

```python
custom_STOPWORDS = []  # Add your own stopwords here

STOPWORDS.extend(custom_STOPWORDS)


from collections import defaultdict


# Inverted index: word -> set of URLs

inverted_index = defaultdict(set)

url_list = set()


# This dictionary will be used to build the connection between links

web_connection = {'source':[], 'target':[]}


import re


# This function will clean the content of web page in order to build the
# inverted index.

def clean_and_tokenize(text):

    text = re.sub(r'[^a-zA-Z0-9\s]', '', text.lower())  # Remove punctuation
and lowercase

    tokens = text.split()

    return [t for t in tokens if t not in STOPWORDS and len(t) > 1]
```

```python
from urllib.parse import urljoin, urlparse


# The crawl function has 5 parameters

# url = The url to crawl

# base_domain = the base domain of the url. During crawling, the crawler will
ignore links from other domains


def crawl(url, base_domain, visited, visit_limit, limit):

    if limit==0 or len(visited)==visit_limit:

        return


    try:

        response = requests.get(url, timeout=5)

        if response.status_code != 200:

            return

    except requests.RequestException:

        return


    visited.add(url)

    print("-"*(10-limit), end=" ")

    print(f"Crawled: {url}")


    soup = BeautifulSoup(response.text, 'html.parser')

    text = soup.get_text(separator=' ', strip=True)

    words = clean_and_tokenize(text)
```

```python
    for word in words:

        inverted_index[word].add(url)

        url_list.add(url)


    # Recursively follow links

    for tag in soup.find_all('a', href=True):

        link = urljoin(url, tag['href'])

        parsed = urlparse(link)


        # Store external links as connection

        web_connection['source'].append(url)

        web_connection['target'].append(link)


        if parsed.netloc == base_domain and link not in visited:

            crawl(link, base_domain, visited, visit_limit, limit-1)




def crawl_roots(root_urls, max_per_root=2, visit_limit=50):

    for root in root_urls:

        print(f"\nStarting crawl from: {root}")

        domain = urlparse(root).netloc

        visited = set()

        crawl(root, domain, visited, visit_limit, max_per_root)
```

```python
seed_urls = [

    'https://www.mykhel.com/cricket/ban-vs-zim-shadman-islam-shines-with-gritty-10
0-anchors-bangladesh-to-commanding-lead-on-day-2-in-358519.html',

    'https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangla
desh-under-pressure-dion-ebrahim',

    'https://gulfnews.com/sport/cricket-prodigy-vaibhav-suryavanshi-smashes-ipl-re
cord-at-14-wins-hearts-of-legends-and-bollywood-icons-1.500109736',

    'https://www.espncricinfo.com/cricket-news',

    'https://sports.ndtv.com/cricket/news',

    'https://www.hindustantimes.com/cricket',

    'https://www.bbc.com/sport/cricket',

    'https://www.icc-cricket.com/news',

    'https://indianexpress.com/section/sports/cricket/',

    'https://www.news18.com/cricket/',

    'https://www.cricket.com.au/news'

]


crawl_roots(seed_urls, max_per_root=10)



# Print first 20 connections


for source, target in list(zip(web_connection['source'],
web_connection['target']))[:20]:

    print(f"{source} -> {target}")
```

```python
import networkx as nx


web_graph = nx.DiGraph()

for i in range(len(web_connection['source'])):

    web_graph.add_edge(web_connection["source"][i],
web_connection["target"][i])




len(web_graph.nodes)


pagerank_scores = nx.pagerank(web_graph, alpha=0.85, max_iter=100, tol=1e-6)

print("\nPageRank Scores:", pagerank_scores)




def search_engine(query, index, scores):

    query_terms = query.lower().split()

    results = set()

    for term in query_terms:

        if term in index:

            if not results:

                results = set(index[term])

            else:

                results = results.intersection(index[term])   # Find common
websites


    # Sort results based on score

    ranked_results = []
```

```python
    for website in results:

        if website in scores:

            ranked_results.append((website, scores[website]))

    ranked_results.sort(key=lambda x: x[1], reverse=True)



    return ranked_results




# Query and display results

query = "Virat Kohli"

print(f"\nSearch Results for '{query}' using PageRank:")

results = search_engine(query, inverted_index, pagerank_scores)


for page, score in results:

    print(f"{page}: ({score})")
```

# Input's & Output's:

```
Stopwords are used when building the inverted index. The inverted index will ignore stopwords.

import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords

STOPWORDS = stopwords.words('english')
print(STOPWORDS)
```

```
['a', 'about', 'above', 'after', 'again', 'against', 'ain', 'all', 'am', 'an', 'and', 'any', 'are', 'aren', "aren't", 'as', 'at', 'be', 'because', 'been', 'before', 'being', 'below', 'between', 'both', 'but', 'by', 'can', 'couldn', "
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
```

+ Code    + Text

```
seed_urls = [
    'https://www.mykhel.com/cricket/ban-vs-zim-shadman-islam-shines-with-gritty-100-anchors-bangladesh-to-commanding-lead-on-day-2-in-358519.html',
    'https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim',
    'https://gulfnews.com/sport/cricket-prodigy-vaibhav-suryavanshi-smashes-ipl-record-at-14-wins-hearts-of-legends-and-bollywood-icons-1.500109736',
    'https://www.espncricinfo.com/cricket-news',
    'https://sports.ndtv.com/cricket/news',
    'https://www.hindustantimes.com/cricket',
    'https://www.bbc.com/sport/cricket',
    'https://www.icc-cricket.com/news',
    'https://indianexpress.com/section/sports/cricket/',
    'https://www.news18.com/cricket/',
    'https://www.cricket.com.au/news'
]

crawl_roots(seed_urls, max_per_root=10)
```

```
Starting crawl from: https://www.mykhel.com/cricket/ban-vs-zim-shadman-islam-shines-with-gritty-100-anchors-bangladesh-to-commanding-lead-on-day-2-in-358519.html

Starting crawl from: https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim
 Crawled: https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim
 - Crawled: https://www.cricbuzz.com/
 -- Crawled: https://www.cricbuzz.com/cricket-match/live-scores
 --- Crawled: https://www.cricbuzz.com/cricket-schedule/upcoming-series/international
 ---- Crawled: https://www.cricbuzz.com/cricket-scorecard-archives
 ----- Crawled: https://www.cricbuzz.com/cricket-news
 ------ Crawled: https://www.cricbuzz.com/cricket-news/editorial/cb-plus
 ------- Crawled: https://www.cricbuzz.com/cricket-news/latest-news
 -------- Crawled: https://www.cricbuzz.com/cricket-news/info/
 --------- Crawled: https://www.cricbuzz.com/cricket-news/editorial/spotlight
 --------- Crawled: https://www.cricbuzz.com/cricket-news/editorial/editorial-list
 --------- Crawled: https://www.cricbuzz.com/cricket-news/editorial/specials
 --------- Crawled: https://www.cricbuzz.com/cricket-news/editorial/stats-analysis
 --------- Crawled: https://www.cricbuzz.com/cricket-news/editorial/interviews
 --------- Crawled: https://www.cricbuzz.com/cricket-news/editorial/live-blogs
 --------- Crawled: https://www.cricbuzz.com/cricket-news/experts/harsha-bhogle/170
 --------- Crawled: https://www.cricbuzz.com/cricket-schedule/series/all
 --------- Crawled: https://www.cricbuzz.com/cricket-series/9237/indian-premier-league-2025
 --------- Crawled: https://www.cricbuzz.com/cricket-series/8796/west-indies-tour-of-england-2025
```

```
# print inverted index
print("\nSample inverted index (first 20 words):")
for word in list(inverted_index.keys())[:20]:
    print(f"{word}: {list(inverted_index[word])}")
```

```
Sample inverted index (first 20 words):
confident: ['https://www.cricbuzz.com/cricket-news/latest-news', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/', 'https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-und
put: ['https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://indianexpress.com/section/sports/', 'https://www.bbc.com/news/uk', 'https://www.bbc.com/news/northern_ireland/northern_ireland_politics', 'https://www.cricket.com.au
bangladesh: ['https://www.news18.com/cricket/videos/', 'https://www.icc-cricket.com/news/team/1', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.cricbuzz.com/cricket-series/8881/icc-world-test-champion
pressure: ['https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.bbc.com/news/bbcverify', 'https://indianexpress.com/section/opinion/', 'https://www.cricket.com.au/players/CA:655/alyssa-healy', 'https://www.cricbuzz.com/
dion: ['https://www.cricbuzz.com/cricket-news/latest-news', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/', 'https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pr
ebrahim: ['https://www.cricbuzz.com/cricket-news/latest-news', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/', 'https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under
cricbuzzcom: ['https://www.cricbuzz.com/cricket-team/denmark/185', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.cricbuzz.com/cricket-series/8881/i
live: ['https://www.news18.com/cricket/videos/', 'https://www.cricket.com.au/videos/4255389/watch-highlights-every-wicket-peter-siddle-final-season-2024-25-summer-victoria-sheffield-shield-one-day-cup-retirement?tags=9022', 'https://
scores: ['https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.cricbuzz.com/cricket-series/8881/icc-world-test-championship-final-2025', 'https://www.cricket.com.au/our-partners', 'https://www.cricket.com.au/
schedule: ['https://www.news18.com/cricket/videos/', 'https://www.icc-cricket.com/news/team/1', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.news18.com/cricketnext/live-score/team-squads/delhi-capita
archives: ['https://www.cricbuzz.com/cricket-team/denmark/185', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.cricbuzz.com/cricket-series/8881/icc-
news: ['https://www.icc-cricket.com/news/team/1', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.news18.com/cricketnext/live-score/team-squads/delhi-capitals-vs-royal-challengers-bengaluru-ddbc04272025
stories: ['https://www.news18.com/cricket/videos/', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.news18.com/cricketnext/live-score/team-squads/delhi-capitals-vs-royal-challengers-bengaluru-ddbc042728
premium: ['https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://indianexpress.com/', 'https://www.cricbuzz.com/cricket-series/8881/icc-world-test-championship-final-2025', 'https://indianexpress.com/section/opinio
editorials: ['https://www.cricbuzz.com/cricket-team/denmark/185', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.cricbuzz.com/cricket-series/8881/ic
latest: ['https://www.news18.com/cricket/videos/', 'https://www.cricbuzz.com/cricket-team/denmark/185', 'https://www.news18.com/cricketnext/live-score/team-squads/delhi-capitals-vs-royal-challengers-bengaluru-ddbc04272025
topics: ['https://www.cricbuzz.com/cricket-team/denmark/185', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.news18.com/cricketnext/live-score/team-squads/delhi-capitals-vs-royal-challengers-bengaluru-ddbc0427202
spotlight: ['https://www.icc-cricket.com/media-releases/', 'https://www.cricbuzz.com/cricket-team/denmark/185', 'https://indianexpress.com/article/world/surreal-bizarre-first-100-days-of-donald-trump-2-0-9973175/', 'https://www.cricb
opinions: ['https://www.cricbuzz.com/cricket-team/denmark/185', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://www.cricbuzz.com/cricket-series/8881/icc-
specials: ['https://www.cricbuzz.com/cricket-team/denmark/185', 'https://www.cricbuzz.com/cricket-team/zimbabwe/12', 'https://www.cricbuzz.com/cricket-news/editorial/editorial-list', 'https://indianexpress.com', 'https://www.cricbuzz
```

```
# Print first 20 connections

for source, target in list(zip(web_connection['source'], web_connection['target']))[:20]:
    print(f"{source} -> {target}")
```

```
https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim -> https://plus.google.com/104502282508811467249
https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim -> Javascript:void(0)
https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim -> Javascript:void(0)
https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim -> https://www.cricbuzz.com/
https://www.cricbuzz.com/ -> https://plus.google.com/104502282508811467249
https://www.cricbuzz.com/ -> Javascript:void(0)
https://www.cricbuzz.com/ -> Javascript:void(0)
https://www.cricbuzz.com/ -> https://www.cricbuzz.com/
https://www.cricbuzz.com/ -> https://www.cricbuzz.com/cricket-match/live-scores
https://www.cricbuzz.com/cricket-match/live-scores -> https://plus.google.com/104502282508811467249
https://www.cricbuzz.com/cricket-match/live-scores -> Javascript:void(0)
https://www.cricbuzz.com/cricket-match/live-scores -> Javascript:void(0)
https://www.cricbuzz.com/cricket-match/live-scores -> https://www.cricbuzz.com/
https://www.cricbuzz.com/cricket-match/live-scores -> https://www.cricbuzz.com/cricket-match/live-scores
https://www.cricbuzz.com/cricket-match/live-scores -> https://www.cricbuzz.com/cricket-schedule/upcoming-series/international
https://www.cricbuzz.com/cricket-schedule/upcoming-series/international -> https://plus.google.com/104502282508811467249
https://www.cricbuzz.com/cricket-schedule/upcoming-series/international -> Javascript:void(0)
https://www.cricbuzz.com/cricket-schedule/upcoming-series/international -> Javascript:void(0)
https://www.cricbuzz.com/cricket-schedule/upcoming-series/international -> https://www.cricbuzz.com/
https://www.cricbuzz.com/cricket-schedule/upcoming-series/international -> https://www.cricbuzz.com/cricket-match/live-scores
```

```
[ ] import networkx as nx

    web_graph = nx.DiGraph()
    for i in range(len(web_connection['source'])):
        web_graph.add_edge(web_connection["source"][i], web_connection["target"][i])

[ ] len(web_graph.nodes)

    7711

[ ] pagerank_scores = nx.pagerank(web_graph, alpha=0.85, max_iter=100, tol=1e-6)
    print("\nPageRank Scores:", pagerank_scores)

    PageRank Scores: {'https://www.cricbuzz.com/cricket-news/134203/confident-that-we-can-put-bangladesh-under-pressure-dion-ebrahim': 0.00012998097259886197, 'https://plus.google.com/104502282508811467249': 0.0001785810369616234, 'Javas...
```

```
[ ] # Query and display results
    query = "Virat Kohli"
    print(f"\nSearch Results for '{query}' using PageRank:")
    results = search_engine(query, inverted_index, pagerank_scores)

    for page, score in results:
        print(f"{page}: ({score})")

    Search Results for 'Virat Kohli' using PageRank:
    https://www.icc-cricket.com/tournaments/world-test-championship/: (0.00018118724310883016)
    https://www.icc-cricket.com/tournaments/champions-trophy-2025: (0.0001798901412613605)
    https://www.icc-cricket.com/tournaments/t20cricketworldcup/: (0.00017982517305346763)
    https://www.cricbuzz.com/cricket-news/editorial/interviews: (0.0001785810369616234)
    https://www.cricbuzz.com/cricket-news/latest-news: (0.0001785810369616234)
    https://www.cricbuzz.com/cricket-series/9237/indian-premier-league-2025: (0.0001785810369616234)
    https://www.icc-cricket.com/videos/: (0.00017845837620627267)
    https://www.news18.com/cricket/videos/: (0.0001709785633539953)
    https://www.news18.com/notifications/: (0.0001709785633539953)
    https://www.news18.com/movies/: (0.0001709785633539953)
    https://www.news18.com/: (0.0001709785633539953)
    https://www.news18.com/cricket/live-score/: (0.0001709785633539953)
    https://www.news18.com/cricket/cricket-schedule/: (0.0001709785633539953)
    https://www.news18.com/cricket/ipl/orange-cap-holder/: (0.0001709785633539953)
    https://www.news18.com/cricket/results/: (0.0001709785633539953)
    https://www.news18.com/cricket/ipl/purple-cap-holder/: (0.0001709785633539953)
    https://www.news18.com/cricket/ipl/points-table/: (0.0001709785633539953)
    https://www.news18.com/cricket/: (0.0001709785633539953)
    https://www.news18.com/cricket/photogallery/: (0.0001709785633539953)
    https://www.news18.com/news/: (0.0001709785633539953)
    https://www.news18.com/cricket/ipl/: (0.0001709785633539953)
    https://www.news18.com/cricket/latestnews/: (0.0001709785633539953)
    https://www.news18.com/livetv/: (0.0001709785633539953)
    https://indianexpress.com/section/sports/: (0.0001626703152631307)
    https://indianexpress.com/section/sports/cricket/: (0.0001626703152631307)
    https://www.news18.com/cricket/suryakumar-yadav-creates-history-becomes-fastest-batter-in-the-world-to-9315214.html: (0.00015007787834640844)
    https://www.icc-cricket.com/news/team/4: (0.0001497359169463637)
    https://www.icc-cricket.com/news/category/world-test-championship: (0.0001497359169463637)
    https://www.news18.com/cricketnext/: (0.00012780617962510066)
```

## Results:

1. Successfully crawled and collected relevant cricket content from trusted websites.

2. Built a link graph and applied PageRank to rank important pages.

3. Improved search relevance by focusing only on cricket-related pages.

4. Achieved better accuracy and efficiency compared to general search engines.

5. Test queries returned high-quality, topic-specific results.

## Future Improvements:

1. Use NLP to improve query understanding.

2. Add real-time content updates through live crawling.

3. Include images and videos in search results.

4. Implement user feedback to refine results.

5. Create a mobile-friendly interface or app.

## Conclusion:

This project successfully demonstrated the development of a domain-specific search engine focused on cricket. By using focused crawling and link analysis techniques like PageRank, the system was able to gather, rank, and deliver relevant cricket-related content from multiple trusted sources. The project highlights the effectiveness of narrowing the search domain to improve accuracy, relevance, and efficiency in information retrieval. This approach can be extended to other domains for better targeted search solutions.