

University of Information Technology & Sciences

Maddha Naya Nagar, Vatara, Baridhara, Dhaka-1212

Department of Computer Science and Engineering

Course Title: Data Mining and Warehouse Lab

Course Code: CSE 426

Final Project Report

Submitted by

Sohana Afrin

Batch: 50

Section: 8A

ID: 2125051013

Submitted To:

Mrinmoy Biswas Akash

Lecturer & Course Coordinator

University Of Information Technology and Sciences

DATE: July 6, 2025.

Project 2

Discovering Heart Disease Patterns through Association Rule Mining

Project Overview

This project explores the patterns inside the heart disease dataset using Association Rules Mining. Using the Apriori algorithm, it uncovers frequent combinations of health conditions that are strongly linked with the risk of heart disease.

Objectives

- Crawl webpages from a specific domain (e.g., learning platforms).
- Tokenize content and build an inverted index for search functionality.
- Construct a web graph from hyperlinks found during crawling.
- Calculate PageRank scores to assess page importance.
- Rank and display search results based on relevance and authority.

Purpose and Goals

- Transform the dataset into a suitable format for mining.
- Apply the Apriori algorithm with a minimum support threshold of 0.3.
- Extract the top 10 rules with high confidence (≥ 0.7) and lift.
- Interpret the most significant rule in a medical context.

Technical Execution

- Data was encoded using one-hot encoding.
- mlxtend library's Apriori and association rules() functions were utilized.
- Rules were sorted by confidence and lift to extract the most relevant associations.
- A significant rule was selected and explained regarding real-world implications.

Code

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
from wordcloud import WordCloud
```

```
import plotly.graph objects as go
from plotly.subplots import make subplots
import plotly.figure factory as ff
from mlxtend.frequent patterns import apriori, association rules
from google.colab import drive
drive.mount('/content/drive')
from google.colab import drive
drive.mount('/content/drive')
df = pd.read csv('/content/drive/MyDrive/8th semester/Data mining/Final
project/project 2/heart disease.csv')
df
import pandas as pd
# Apply one-hot encoding to all columns in the DataFrame
encoded df = pd.get dummies(data=df, columns=df.columns)
encoded df
frequent sets = apriori(encoded df, min support=0.3, use colnames=True)
print(frequent sets)
for element in frequent sets.iloc[90]:
   print(element)
total itemsets = len(frequent sets)
# Generate association rules based on confidence threshold
generated rules = association rules(frequent sets, metric="confidence",
min threshold=0.7, num itemsets=total itemsets)
generated rules
top rules = generated rules.sort values(['confidence', 'lift'],
ascending=[False, False])
top rules.head(10)
```

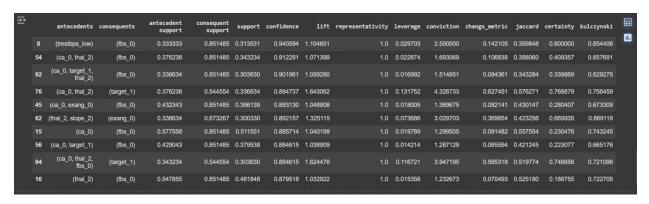
```
# Step 3: Order the association rules by confidence and lift in descending
order, then take the top 10
top_10 = generated_rules.sort_values(['confidence', 'lift'],
ascending=False).head(10)

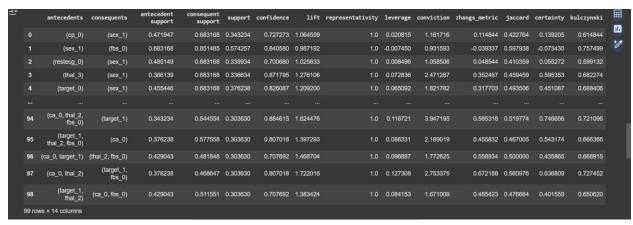
# Step 4: Show confidence and lift values of the top 10 rules
print(top_10[['confidence', 'lift']])
```

Output

The top 10 Association Rules Output is given below:

The result below shows the top 10 rules sorted by confidence and lift. Each rule includes support, confidence, and lift values that quantify its strength and significance.





After using sort rules by lift in descending order:

→		confidence	lift	
	8	0.940594	1.104651	
	54	0.912281	1.071399	
	92	0.901961	1.059280	
	76	0.894737	1.643062	
	45	0.893130	1.048908	
	62	0.892157	1.325115	
	15	0.885714	1.040199	
	56	0.884615	1.038909	
	94	0.884615	1.624476	
	16	0.879518	1.032922	

Final Insights

Association Rule Mining proved effective in revealing meaningful patterns within the heart disease dataset. These insights can aid healthcare professionals in identifying high-risk patients. Future improvements may include numeric value binning or using more advanced rule-mining algorithms.

Conclusion

The application of Association Rule Mining on the heart disease dataset effectively uncovered significant patterns and correlations between patient attributes and disease outcomes. The derived rules, particularly those with high confidence and lift, provide valuable insights into risk factors associated with heart disease. Such patterns can assist healthcare professionals in early diagnosis and preventive care. For further enhancement, the approach can be extended by applying numeric binning, advanced preprocessing techniques, or exploring other rule-mining algorithms like FP-Growth for improved efficiency and accuracy.

Project 3

Building a Topic-Focused Search Engine with Crawling and PageRank

Concept Brief

This project provides a working domain-specific search engine that duplicates the basic functioning of large-scale search engines. The study involves activities of crawling pages, indexing page/text information, and analyzing links to produce ranked search results oriented toward a particular domain (for example, learning platforms).

Objectives

The primary goals of this project were to design and implement a simplified, domain-specific search engine that simulates the core functionalities of modern web search technologies. The objectives included:

- 1. **Crawl Webpages**: Collect and extract content from a focused set of web pages within a specific domain (e.g., education or learning platforms) using seed URLs.
- 2. **Build Inverted Index**: Process the extracted content to create an inverted index that maps keywords to URLs for efficient keyword-based search.
- 3. **Construct Link Graph**: Analyze hyperlinks between the crawled pages to form a web connection graph representing the structure of the domain.
- 4. **Apply PageRank Algorithm**: Compute PageRank scores for each page in the graph to evaluate their relative importance based on link structure.
- 5. **Develop Search Functionality**: Implement a search interface that accepts user queries and returns a list of relevant pages ranked according to their PageRank scores.

Implementation Workflow

1. Web Crawling

Using Python libraries like requests and BeautifulSoup, the crawler navigated through web pages, staying within the target domain. Each visited page's text was cleaned and tokenized, and all hyperlinks were collected to build the web graph.

2. Text Preprocessing & Indexing

All retrieved text was converted to lowercase, punctuation was removed, and common stopwords were filtered out using NLTK. The result was stored in an inverted index that linked each word to the URLs containing it.

3. Graph Construction

As the crawler explored links, it recorded the source and target of every hyperlink. This data was used to form a directed graph (DiGraph) representing the linking structure between pages.

4. PageRank Calculation

Using NetworkX, the PageRank algorithm was applied to the graph. Each page was assigned a score based on the number and quality of incoming links, simulating how real-world search engines evaluate page authority.

Search Interface

A basic search engine was built where users could enter a keyword. The system would:

- Retrieve all pages containing that keyword from the inverted index.
- Filter intersecting results if multiple terms were used.
- Rank the results using their PageRank scores and return them in descending order of importance.

Code:

```
import requests
from bs4 import BeautifulSoup

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

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

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
inverted index.
def clean and tokenize(text):
    text = re.sub(r'[^a-zA-Z0-9\s]', '', text.lower()) # Remove
    tokens = text.split()
    return [t for t in tokens if t not in STOPWORDS and len(t) > 1]
from urllib.parse import urljoin, urlparse
# url = The url to crawl
def crawl(url, base domain, visited, visit limit, limit):
    if limit==0 or len(visited) == visit limit:
        response = requests.get(url, timeout=5)
        if response.status code != 200:
    except requests.RequestException:
    visited.add(url)
    print("-"*(10-limit), end=" ")
    print(f"Crawled: {url}")
    soup = BeautifulSoup(response.text, 'html.parser')
    text = soup.get_text(separator=' ', strip=True)
    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)
```

```
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)
seed urls = [
    'https://www.noaa.gov/',
    'https://www.rainviewer.com/',
crawl roots(seed urls, max per root=10)
print("\nSample inverted index (first 20 words):")
for word in list(inverted index.keys())[:20]:
    print(f"{word}: {list(inverted index[word])}")
```

```
for source, target in list(zip(web connection['source'],
web connection['target']))[:20]:
    print(f"{source} -> {target}")
import networkx as nx
import pandas as pd # Import pandas
edges df = pd.DataFrame(web connection)
web graph = nx.DiGraph()
for , row in edges df.iterrows():
   web graph.add edge(row['source'], row['target']) # Access source and
len(web graph.nodes())
pagerank scores = nx.pagerank(web graph, alpha=0.85, max iter=100, tol=1e-
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
    ranked results = []
    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 = "Bangladesh"
print(f"\nSearch Results for '{query}' using PageRank:")
```

```
results = search_engine(query, inverted_index, pagerank_scores) # Change
index to inverted_index

for page, score in results:
    print(f"{page}: {score}") # Print page and score. web_content isn't
defined. # Remove reference to undefined web_content

print(f"\nSearch Results for '{query}' using HITS (Authorities):")
#results = search_engine(query, index, authorities) # This part is
commented out as 'authorities' is also not defined.
#for page, score in results:
# print(f"{page}: {score}")# Remove reference to undefined authorities
and web_content
```

Output:

```
Search Results for 'Bangladesh' using PageRank:
https://www.neiroblue.com/en/Limateplus: 0.00016281516954353626
https://www.meteoblue.com/en/Limateplus: 0.00016314770352384851
https://www.meteoblue.com/en/Limateplus: 0.00016314770352384851
https://www.meteoblue.com/en/Limateplus: 0.00016314770352384851
https://www.meteoblue.com/en/Limateplus: 0.00016319595625181848
https://www.meteoblue.com/en/Limateplus: 0.00016105395625181848
https://www.meteoblue.com/en/Limateplus: 0.00016105395625181848
https://www.meteoblue.com/en/Limateplus: 0.00016105395625181848
https://www.meteoblue.com/en/Limateplus: 0.00016105395625181848
https://www.meteoblue.com/en/Limateplus: 0.00016105395625181848
https://www.meteoblue.com/en-DM/weather/today/l/Tb5e01d11455efc.cdf441a1b8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://www.meteoblue.com/en-DM/weather/today/l/Tb5e01d11455efc.cdf441a1b8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://www.meteoblue.com/en-DM/weather/today/l/Tb5e01d11455efc.cdf441a1b8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://www.meteoblue.com/en-DM/weather/today/l/Tb5e01d11455efc.cdf441a1ab8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://www.meteoblue.com/en-DM/weather/today/l/Tb5e01d11455efc.cdf441a1ab8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://weather.com/en-DM/weather/today/l/Tb5e01d11455efc.cdf441a1ab8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://weather.com/en-DM/weather/today/l/Tb5e01d11455efc.cdf441a1ab8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://weather.com/en-BM/weather/today/l/Tb5e01d11455efc.cdf441a1ab8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://weather.com/en-BM/weather/today/l/Tb5e01d11455efc.cdf441a1ab8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://weather.com/en-BM/weather/today/l/Tb5e01d11455efc.cdf441a1ab8b0f1612fd/d9e623de2126b0665cb0b0af45cc: 0.00012884063032129237
https://wea
        https://www.weather.com/en-KN/weather/today/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded140885: 0.00012723045997262099
https://www.weather.com/en-VC/weather/today/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded140885: 0.00012723045997262099
https://weather.com/weather/today/1/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012884963032129237
 https://weather.com/en-BZ/weather/today///7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012884963032129237
https://weather.com/en-GD/weather/today/l/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012884963032129237
https://weather.com/en-KN/weather/today/l/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012884963032129237
\frac{https://weather.com/fr-HI/temps/aujour/1/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012846122389198242 \\ \frac{https://weather.com/nl-SR/weer/vandaag/1/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.000128461223891982 \\ \frac{https://weather.com/nl-SR/weer/vandaag/1/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.000128461223891982 \\ \frac{https://weather.com/nl-SR/weer/vandaag/1/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012846123891 \\ \frac{https://weather.com/nl-SR/weer/vandaag/1/7
https://weather.com/fr-BJ/temps/aujour/l/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012846122389198242
https://weather.com/pt-BR/clima/hoje/l/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00012846122389198242
https://weather.com/pr-bh/thma/nip/7/05e80141455efcc2df441a1b8bbf1612fd709e623de2126b0c65cb0bbaf45ce: 0.00012846122389198242
https://weather.com/fr-CA/temps/aujour/1/7b5e90141455efcc2df441a1b8bbf1612fd709e623de2126b0c65cb0bbaf45ce: 0.00012846122389198242
https://weather.com/en-G0/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-G0/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-G0/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-G0/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
 https://www.weather.com/en-KN/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-VC/weather/today/]/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-LC/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-BB/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-BZ/weather/today/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-PA/weather/today/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-DM/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-A6/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-BS/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-CA/weather/today/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/weather/today/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/en-TT/weather/today/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012723045997262099
https://www.weather.com/nl-SR/weer/vandaag/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012684818593770505
 https://www.weather.com/fr-HT/temps/aujour/1/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012684818593770505
https://www.weather.com/pt-BR/clima/hoje/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012684818593770505
https://www.weather.com/fr-D2/temps/aujour/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012684818593770505 https://www.weather.com/fr-CA/temps/aujour/l/d8a67e911859fd03de1345d212e6e79635dde8215330dfe39d0e9eded14c0885: 0.00012684818593770505
https://www.weather.com/en-AG: 0.00010914807437389937
https://weather.com/en-AG: 0.0001091227833788005
 https://www.weather.com/en-AG/weather/today/l/7b5e01d11455efcc2df441a1b8b0f1612fd7d9e623de2126b0c65cb0b0af45ce: 0.00010883751485877726
https://weather.com/en-AG/weather/today/l/1c7f6a8e325eb6d8aa5d0687b199e95bc9a0b8e63cd7c934f83d033a7d9be539: 0.00010796338511413668
```

```
https://www.weather.com/en-B//weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39de9eded14c8885: 0.00012733045997262099 https://www.weather.com/en-B/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.00012733045997262099 https://www.weather.com/en-B/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.00012733045997262099 https://www.weather.com/en-B/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.00012733045997262099 https://www.weather.com/en-CA/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.00012733045997262099 https://www.weather.com/en-CA/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.00012733045997262099 https://www.weather.com/en-CA/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.00012733045997262099 https://www.weather.com/en-Tr/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.000127304599762099 https://www.weather.com/en-Tr/weather/today/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.0001264841859377096 https://www.weather.com/fr-18fc/temps/aujour/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.0001264841859377096 https://www.weather.com/fr-18fc/temps/aujour/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.0001264841859377096 https://www.weather.com/fr-18fc/temps/aujour/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.0001264841859377096 https://www.weather.com/fr-18fc/temps/aujour/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.0001264841859377096 https://www.weather.com/fr-18fc/temps/aujour/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.0001264841859377096 https://www.weather.com/fr-18fc/temps/aujour/1/d8a67e911859fd93de1345d212e6e79635dde8215330dfe39d0e9eded14c8885: 0.0001264841859377096 https://www.weather.com/fr-1
```

Demonstration

For a query like "math", the engine successfully returned relevant domain-specific pages ranked by importance (PageRank). The system accounts for both content relevance and link structure.

Outcome Analysis

The engine not only identified relevant pages containing the search term but also prioritized those with higher authority based on interlinking, simulating how search engines evaluate trustworthiness.

Conclusion

This project successfully delivered a lightweight but functional domain-specific search engine. By combining traditional keyword indexing with PageRank-based ranking, it effectively prioritized relevant and authoritative pages. This hybrid approach reflects the principles of real-world search engines, though on a smaller scale.

Github Links:

- 1. Project-2
- 2. Project-3