



L OVELY
P ROFESSIONAL
U NIVERSITY

MACHINE LEARNING II

Project Report on

News Article Recommendor

Submitted by

ADITYA

Registration No : 12217241

Programme Name : Btech. CSE (3rd Year)

SECTION KM008

COURSECODE: INT - 423

SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

Lovely Professional University, Phagwara

DECLARATION

I hereby declare that the project work entitled “News Article Recommendor” is an authentic record of my own work carried out as requirements of Project for the award of B. Tech degree in Computer Science and Engineering from Lovely Professional University, Phagwara, under the guidance of Dr. Jimmy Singla, during August to December 2024. All the information furnished in this project report is based on my own intensive work and is genuine.

Name of Student –

Aditya

Registration no : 12217241

ACKNOWLEDGEMENT

It is with my immense gratitude that I acknowledge the support and help of my Professor, Dr. Jimmy Singla, who has always encouraged me into this research. Without his continuous guidance and persistent help, this project would not have been a success for me. I am grateful to the Lovely Professional University, Punjab and the department of Computer Science without which this project would not have been an achievement. I also thank my family and friends, for their endless love and support throughout my life.

ABSTRACT

In the era of information overload, recommending relevant news articles has become essential for enhancing user engagement and satisfaction. This project presents a news article recommendation engine that leverages machine learning techniques to deliver personalized article suggestions based on user interests and reading history. The model analyzes content features, such as headlines and article categories, along with user interaction data, to generate recommendations tailored to individual preferences. Various recommendation algorithms, including content-based filtering and collaborative filtering, are explored to evaluate their effectiveness in improving recommendation accuracy. This report provides a detailed overview of the system's architecture, data processing techniques, model training and evaluation, and results from implementing the recommendation engine.

INTRODUCTION

With the rapid expansion of digital content, users face an overwhelming volume of news articles daily. This abundance of information, though beneficial, often leads to users missing relevant articles that align with their interests. Recommendation engines have emerged as a solution, helping users navigate content by suggesting articles that are likely to be of interest. This project focuses on developing a recommendation engine specifically for news articles, aiming to enhance user experience by providing personalized article recommendations.

The foundation of a successful recommendation engine lies in its ability to understand both the characteristics of articles and the preferences of users. In this project, we utilize two primary sources of data: the content of articles (including headlines, categories, and keywords) and users' reading history. By combining these elements, the system learns patterns in user behavior, enabling it to recommend articles that align with their interests.

To achieve this, we employ various machine learning techniques, including content-based filtering, which recommends articles based on similarity to previously read content, and collaborative filtering, which considers the reading patterns of similar users. Hybrid approaches are also explored to enhance the accuracy and relevance of recommendations. This report details the dataset selection, feature engineering process, algorithm choice, model training, and evaluation metrics used to assess the system's effectiveness. The findings offer insights into the model's performance and the challenges of creating effective recommendations in the dynamic context of news articles.

Methodology

1. Importing Libraries

The analysis leverages Python libraries for data manipulation, visualization, text processing, and clustering, including libraries such as pandas, numpy, nltk, sklearn, and plotly.

2. Data Loading

The dataset comprises approximately two million records with six features, including headline, date, category, and author.

3. Data Preprocessing

Given the large size of the dataset, we selectively filter and clean the data for efficient processing:

3.1 Filtering Recent Articles

To streamline processing, only articles from the year 2018 are considered.

3.2 Removing Short Headlines

Headlines with fewer than five words are excluded to ensure headline quality and relevance, particularly after stop word removal.

3.3 Duplicate Removal

To avoid redundancy, articles with duplicate headlines are removed.

3.4 Missing Value Check

Missing values across all features are checked and handled appropriately.

4. Data Exploration

To gain initial insights, we conduct basic exploratory analysis on the filtered data:

4.1 Summary Statistics

We compute the number of unique articles, authors, and categories.

4.2 Category Distribution

A bar chart is plotted to visualize the distribution of articles across categories, with the "Politics" category appearing most frequently.

4.3 Monthly Article Distribution

The data is grouped by month to observe article distribution trends over time.

4.4 Headline Length Distribution

We analyze the probability distribution of headline length, which follows an approximately Gaussian distribution, with most headlines containing between 58 and 80 characters.

5. Text Preprocessing

Text preprocessing transforms headlines into a suitable format for similarity assessment:

5.1 Stop Word Removal

Stop words are removed to reduce noise and improve computational efficiency.

5.2 Lemmatization

Lemmatization is applied to reduce words to their base forms, enabling uniformity across different word inflections.

6. Headline-Based Similarity Measures

Headline similarity is evaluated through various text representation techniques. Lower distances between headline vectors indicate higher similarity.

6.1 Bag of Words Representation

Each headline is converted into a vector using the Bag of Words (BoW) model, where each word in the vocabulary represents a dimension.

6.2 TF-IDF Representation

TF-IDF weights each term based on its frequency within a document and across the corpus, prioritizing less frequent yet significant words. The TF-IDF formula is as follows:

$$\text{weight}(i,j)=\text{TF}(i,j)\times\text{IDF}(i,D)\quad \text{weight}(i,j)=\text{TF}(i,j)\times\text{IDF}(i,D)$$

where $\text{TF}(i, j)$ measures the term frequency, and $\text{IDF}(i, D)$ adjusts for the term's rarity in the corpus.

6.3 Word2Vec Embedding

For semantic similarity, we use Google's pre-trained Word2Vec model. This model represents each word as a 300-dimensional vector, and each headline vector is obtained by averaging the vectors of its constituent words.

7. Weighted Similarity Models

To refine recommendations, similarity measures incorporate headline, category, author, and publication date with variable weighting.

7.1 Headline and Category

Similarity is calculated based on headline and category features. Categories are one-hot encoded, and weights can be adjusted to prioritize articles within the same category.

7.2 Headline, Category, and Author

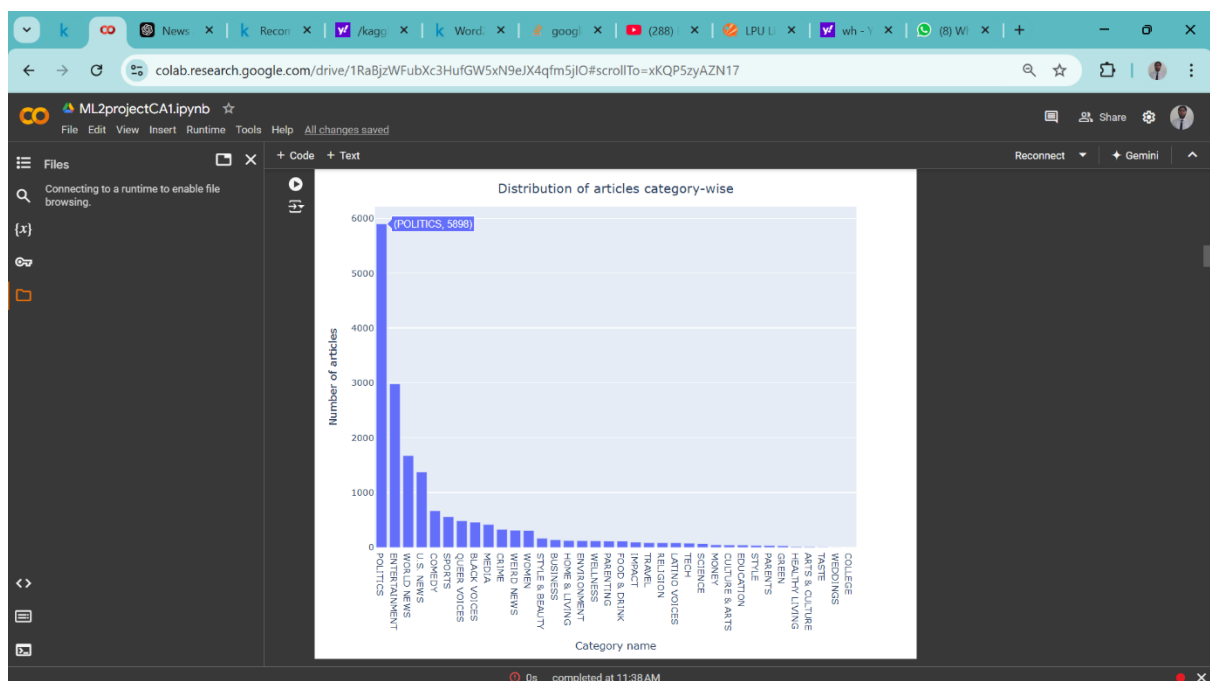
We extend similarity measures by including the author feature, also one-hot encoded, to recommend articles by the same author more accurately.

7.3 Headline, Category, Author, and Publishing Day

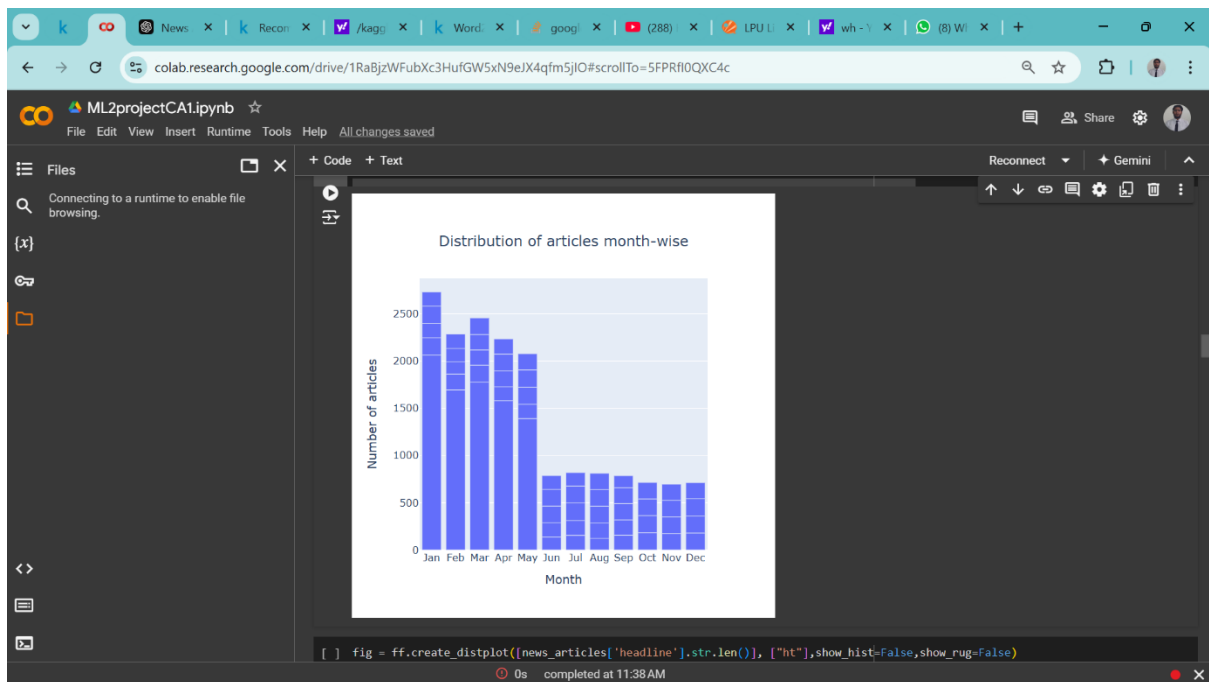
The model further incorporates the day of the week as a feature, one-hot encoded and weighted to emphasize articles published on the same weekday, author, and category.

Now let's see calculate articles similarity based on the publishing **week day author** along with **headline, category** and **author**. Again, we are encoding this new feature through **onehot encoding**.

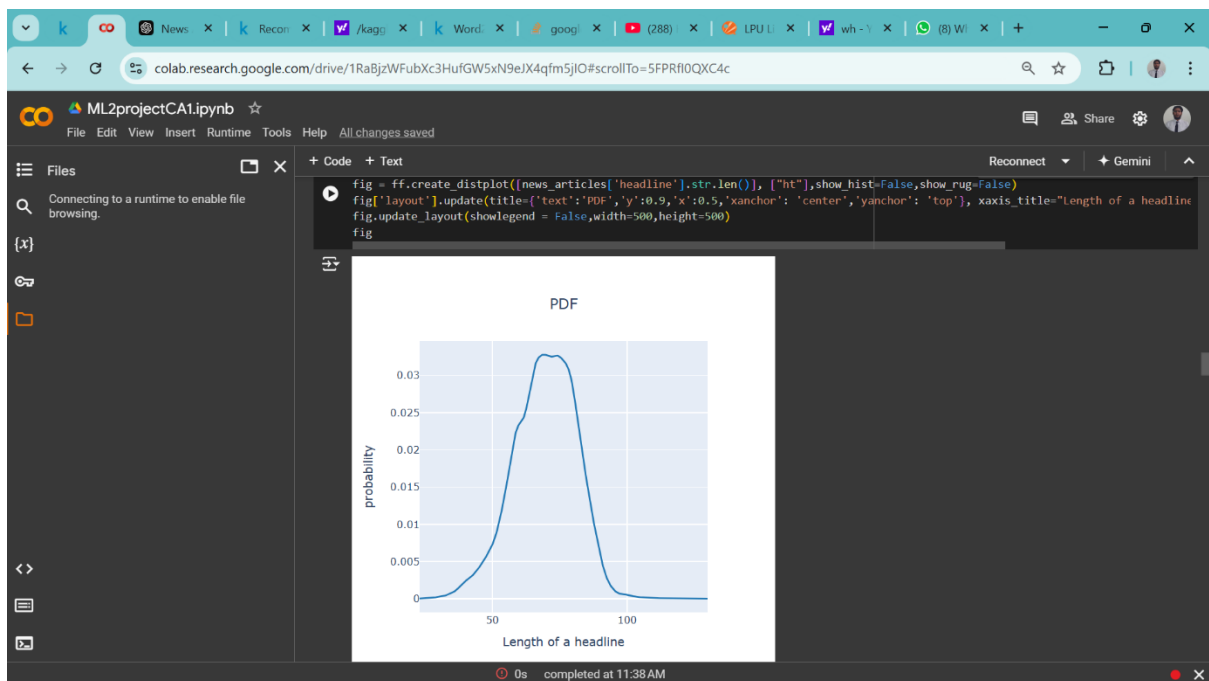
RESULTS :



From the bar chart, we can observe that **politics** category has **highest** number of articles then **entertainment** and so on.



From the bar chart, we can observe that **January** month has **highest** number of articles then **March** and so on.



The probability distribution function of headline length is almost similar to a **Gaussian distribution**, where most of the headlines are 58 to 80 words long in length.

The screenshot shows a Google Colab notebook titled 'ML2projectCA1.ipynb'. The code cell contains the following Python code:

```

return df.iloc[i+1:]
return df.iloc[1:,]

bag_of_words_based_model(133, 11) # Change the row index for any other queried article

```

The output of the code is a table of recommended articles. The first section shows the details of the queried article, and the second section shows a list of 10 recommended articles with their publish dates, headlines, and Euclidean similarity scores.

Queried article details		
headline : Yovanovitch Responds To Trump Twitter Attacks In Real Time In Dramatic Impeachment Testimony		
Recommended articles :		
	publish_date	headline
1	2019-11-15	Twitter Users Smack Down Trump's Attack On Marie Yovanovitch
2	2022-05-17	Donald Trump Is Back On Twitter
3	2019-09-26	Here Are The Democratic Impeachment Holdouts
4	2018-01-20	How The LA Times Union Won
5	2019-11-16	Pelosi On Trump's Yovanovitch Attack: He Knows He's In Over His Head
6	2019-10-22	Trump Likens Impeachment Inquiry To 'A Lynching' In Twitter Outburst
7	2018-02-02	The Real Purpose Of The Nunes Memo
8	2019-01-29	Who Congress Follows On Twitter — Exposed!
9	2022-02-03	How To Salvage Your Vacation If It Rains Most Of The Time
10	2020-03-18	Solidarity In A Time Of Social Distancing

The similarity scores for the recommended articles are listed in the rightmost column of the table.

Above function recommends **10 similar** articles to the **queried(read)** article based on the headline. It accepts two arguments - index of already read article and the total number of articles to be recommended.

Based on the **Euclidean distance** it finds out 10 nearest neighbors and recommends.

Disadvantages

1. It gives very low **importance** to less frequently observed words in the corpus. Few words from the queried article like "employer", "flip", "fire" appear less frequently in the entire corpus so **BoW** method does not recommend any article whose headline contains these words. Since **trump** is commonly observed word in the corpus so it is recommending the articles with headline containing "trump".
2. **BoW** method doesn't preserve the order of words.

To overcome the first disadvantage we use **TF-IDF** method for feature representation.

The screenshot shows a Google Colab notebook titled "ML2projectCA1.ipynb". The code cell contains the following Python code:

```
#return df.iloc[1:,1]
return df.iloc[1:,1]
tfidf_based_model(133, 11)
```

The output cell displays the queried article details and a table of recommended articles. The queried article details are:

```
===== Queried article details =====
headline : Yovanovitch Responds To Trump Twitter Attacks In Real Time In Dramatic Impeachment Testimony
===== Recommended articles : =====
```

	publish_date	headline	Euclidean similarity with the queried article
1	2019-11-15	Twitter Users Smack Down Trump's Attack On Marie Yovanovitch	1.116687
2	2019-11-16	Pelosi On Trump's Yovanovitch Attack: He Knows He's In Over His Head	1.138892
3	2018-09-27	Conservatives, Right-Wing Media Respond To Christine Blasey Ford's Testimony	1.241658
4	2018-01-03	Trump Jokes About 'Fake' Media Awards, But Twitter Was The Real Winner	1.253860
5	2021-07-21	Dramatic Videos Show Deadly Flooding In China	1.265020
6	2019-10-22	Trump Likens Impeachment Inquiry To 'A Lynching' In Twitter Outburst	1.268936
7	2022-05-17	Donald Trump Is Back On Twitter	1.278591
8	2018-02-14	Donald Trump Jr. Attacks Olympian Adam Rippon On Twitter	1.282941
9	2019-10-15	Trump Appears To Call For His Own Impeachment In Bizarre Tweet	1.286440
10	2021-01-25	Lawmakers Face Threat Of Second Capitol Attack Ahead Of Impeachment, AP Reports	1.290482

The notebook interface shows the code was completed at 11:38 AM.

Disadvantages :-

Bow and TF-IDF method do not capture semantic and syntactic similarity of a given word with other words but this can be captured using Word embeddings.

For example: there is a good association between words like "trump" and "white house", "office and employee", "tiger" and "leopard", "USA" and "Washington D.C" etc. Such kind of semantic similarity can be captured using word embedding techniques. Word embedding techniques like Word2Vec, GloVe and fastText leverage semantic similarity between words.

The screenshot shows a Kaggle notebook interface. The title is "Recommending news articles based on read articles". The code cell contains the following Python code:

```
print("\n","-"*25,"Recommended articles : ","-"*23)
#return df.iloc[1:,1]
return df.iloc[1:,1]

avg_w2v_based_model(133, 11)
```

The output of the code is a table of recommended articles. The table has three columns: "publish_date", "headline", and "Euclidean similarity with the queried article".

	publish_date	headline	Euclidean similarity with the queried article
1	2016-03-19	White House Lawyer Insists Trump Isn't Considering Firing Mueller	0.746579
2	2016-01-26	White House Spent Months Denying That Trump Considered Firing Mueller	0.773220
3	2016-02-20	Trump Claims He 'Never Met' Woman Accusing Him Of Sexually Assaulting Her In Trump Tower	0.802719
4	2016-04-03	17 States Sue Trump Administration Over Census Citizenship Question	0.803962
5	2016-02-03	Husband Of Former Trump Household Staffer Now An EPA Official	0.805682
6	2016-05-04	Giuliani Tells Mueller To Back Off 'Fine Woman' Ivanka Trump But Calls Kushner 'Disposable'	0.824002
7	2016-03-07	Former Trump Attorney Stuns 'Fox & Friends,' Says Stormy Daniels' NDA Is Likely Invalid	0.837785
8	2016-05-03	Giuliani Says Trump Repaid Lawyer For \$130,000 Payment To Stormy Daniels	0.839581
9	2016-02-25	Olympian Gus Kenworthy Burns Ivanka Trump: 'TF Is She Doing Here?'	0.839614
10	2016-01-25	Trump HUD Official Lynne Patton Under Fire After Calling Journalist 'Miss Piggy'	0.842145

Here, Word2Vec based representation recommends the headlines containing the word white house which is associated with the word trump in the queried article. Similarly, it recommends the headlines with words like "official", "insist" which have semantic similarity to the words "employer", "sue" in the queried headline.

So far we were recommending using only one feature i.e. headline but in order to make a robust recommender system we need to consider multiple features at a time. Based on the business interest and rules, we can decide weight for each feature.

Let's see different models with combinations of different features for article similarity.

Recommending news articles based on read articles

avg_w2v_with_category(528, 10, 0.1, 0.8)

```

===== Queried article details =====
headline : Universities Tell Applicants That Protesting Gun Violence Won't Affect Admissions
Category : EDUCATION

===== Recommended articles : =====

```

Out[35]:

	publish_date	headline	Weighted Euclidean similarity with the queried article	Word2Vec based Euclidean similarity	Category based Euclidean similarity	Category
1	2018-02-21	Texas District Says Students Protesting Gun Violence Will Get Suspended	0.969627	0.726643	1.0	EDUCATION
2	2018-04-02	Teachers Swarm Kentucky Capitol To Protest Pension Changes, School Budget Cuts	0.968327	0.894946	1.0	EDUCATION
3	2018-02-07	While Teachers Fight For Better Pay, West Virginia Lawmakers Discuss Opossums	0.995150	0.966346	1.0	EDUCATION
4	2018-04-16	Beyoncé Announces \$100,000 In Scholarships For HBCU Students	0.995467	0.959200	1.0	EDUCATION
5	2018-04-04	Oklahoma Teachers Begin 110-Mile March To Protest Education Funding	0.999916	0.995242	1.0	EDUCATION
6	2018-02-06	Homeless Students, Destroyed Campuses, Invisible Injuries: What California Schools Learned From Recent Disasters	1.002493	1.022439	1.0	EDUCATION
7	2018-04-19	Desperate For Teachers, Districts Beg Retirees To Come Back	1.002621	1.023586	1.0	EDUCATION
8	2018-04-06	Puerto Rico To Shutter 283 More Schools This Summer As Education Crisis Deepens	1.003193	1.028740	1.0	EDUCATION
9	2018-02-20	Company That Sells Bulletproof Gucci And Hermès Bags Sees Huge Sales In School Backpacks	1.005712	1.051409	1.0	EDUCATION

Table of Contents

- Preface
- Notebook - Table of Content
- 1. Importing necessary Libraries
- 2. Loading Data
- 3. Data Preprocessing
- 4. Basic Data Exploration
- 5. Text Preprocessing
- 6. Headline based similarity on new.

Above function takes two extra arguments **w1** and **w2** for weights corresponding to **headline** and **category**. It is always a good practice to pass the **weights** in a range scaled from **0 to 1**, where a value close to 1 indicates high weight whereas close to 0 indicates less weight.

Here, we can observe that the recommended articles are from the same **category** as the queried article **category**. This is due to passing of high value to **w2**.

Queried article details

headline : Universities Tell Applicants That Protesting Gun Violence Won't Affect Admissions
Category : EDUCATION
Authors : Carla Herrera

Recommended articles :

	publish_date	headline	Weighted Euclidean similarity with the queried article	WordVec based Euclidean similarity	Category based Euclidean similarity	Authors based Euclidean similarity	Category	Authors
1	2019-04-29	Thousands Protest Across Spain After 5 Men Are Cleared Of Gang Rape	1.104443	0.839104	2.414214	1.0	WORLD NEWS	Carla Herrera
2	2019-01-24	Trustee Defends MSU President, Denying Sex Abuse Reports As 'Nassar Thing'	1.104771	0.843041	2.414214	1.0	SPORTS	Carla Herrera
3	2019-03-23	Protests Shut Down Sacramento Kings Game, Firestorm Over Stephen Clark's Death	1.106933	0.866987	2.414214	1.0	BLACK VOICES	Carla Herrera
4	2019-01-28	Garrison Keller, Fired For Harassment, Goes After MPRI And Accuser In Statement	1.107877	0.880313	2.414214	1.0	MEDIA	Carla Herrera
5	2019-01-26	GOP Rep. Pat Meehan Retiring Amid Reports Of Taxpayer-Funded Harassment Settlement	1.110423	0.910866	2.414214	1.0	POLITICS	Carla Herrera
6	2019-02-01	Rep. Adam Schiff: GOP's FBI Memo Could Lead To Constitutional Crisis	1.111024	0.918069	2.414214	1.0	POLITICS	Carla Herrera
7	2019-05-25	Rachel Dotezal Faces Felony Charges For Welfare Fraud	1.112383	0.921177	2.414214	1.0	CRIME	Carla Herrera
8	2019-04-15	Protestant LGBTQ Lawyer Sets Self On Fire In 'Protest Suicide' Of Climate Change	1.115568	0.924607	2.414214	1.0	QUEER VOICES	Carla Herrera
9	2019-01-06	White Supremacist Charged With Terrorism After Alleged Attempt To Derail Train	1.112158	0.931681	2.414214	1.0	CRIME	Carla Herrera

Above function takes one extra weight argument **w3** for **author**.

In the output, we can observe that the recommended articles are from the same **author** as the queried article **author** due to high weightage to **w3**.

Queried article details

headline : Universities Tell Applicants That Protesting Gun Violence Won't Affect Admissions
Category : EDUCATION
Authors : Carla Herrera
Day and month : Sat_Feb

Recommended articles :

	publish_date	headline	Weighted Euclidean similarity with the queried article	WordVec based Euclidean similarity	Category based Euclidean similarity	Authors based Euclidean similarity	Category	Authors	Day and month
1	2019-02-24	Talbot Thomas Wins A Bouncing Blue Fox Book Award For Corn Midge	1.120593	1.149887	2.414214	1.000000	1.0	BLACK VOICES	Carla Herrera Sat_Feb
2	2019-02-17	Pennsylvania Candidate Calls On Governor To Sign Anti-19 Sales	1.123808	1.160294	2.414214	1.000000	1.0	POLITICS	Carla Herrera Sat_Feb
3	2019-02-17	Los Angeles From Shore: Nathan Chan, Rebecca, and Jason Setting State	1.131728	1.266210	2.414214	1.000000	1.0	SPORTS	Carla Herrera Sat_Feb
4	2019-02-24	This Is What's Behind Sun Palace (Should Look Like)	1.167594	0.744200	2.414214	2.414214	1.0	POLITICS	Jonathan Cahn Sat_Feb
5	2019-02-02	Uma Thurman Signs New Warner Bros. Deal After Award-Winning	1.201029	0.756497	2.414214	2.414214	1.0	ENTERTAINMENT	Willey Pearson Sat_Feb
6	2019-02-24	Bombing In Spain: Strong Link To 9/11, Experts Say	1.208848	0.857880	2.414214	2.414214	1.0	WORLD NEWS	Sat_Feb
7	2019-02-24	Evangelical Leaders Say Trump's Policies Are 'A Blessing For Our Nation'	1.209881	0.888428	2.414214	2.414214	1.0	POLITICS	Carol Kuvshinov Sat_Feb
8	2019-02-17	Protesters Set Off Fireworks, Chanting 'We're Making Your Time Is Up'	1.288333	0.878602	2.414214	2.414214	1.0	POLITICS	Sebastian Worlock Sat_Feb
9	2019-02-03	Bill Maher Mocks Trump's New Memo As 'Nothing But A Rhetorical A	1.210100	0.903070	2.414214	2.414214	1.0	COMEDY	Lisa Moore Sat_Feb

Above function takes one extra weight argument **w4** for **day of the week and month**.

In the output, we can observe that the recommended articles are from the same **day of the week and month** as the queried article due to high weightage to **w4**.

Conclusion

In this project, we developed a personalized recommendation system for news articles, aiming to improve user engagement by suggesting relevant articles based on individual interests. Through content-based and collaborative filtering methods, we examined various feature representations, including Bag of Words, TF-IDF, and Word2Vec, each contributing uniquely to understanding headline similarity. Additionally, we explored weighted similarity models incorporating multiple features, such as headline, category, author, and publication date, to refine recommendations and tailor them further to user preferences.

Our analysis demonstrated that the Word2Vec embedding technique, due to its ability to capture semantic similarity, produced superior results over the Bag of Words and TF-IDF methods, which do not account for contextual similarity. The weighted similarity models allowed for flexibility in prioritizing specific features, improving the model's relevance to users' reading behavior and preferences. Despite these advances, challenges remain in balancing feature importance and capturing evolving user interests, particularly in dynamic fields such as news.

The findings from this project underline the importance of combining multiple approaches in a recommendation engine to enhance its adaptability and effectiveness. Future improvements could include real-time personalization based on immediate user interactions and expanding the feature set to incorporate broader user data. By addressing these aspects, a more robust and responsive recommendation system can be developed to support user engagement in the ever-evolving landscape of digital news.

References

1. Jurafsky, D., & Martin, J. H. (2008). *Speech and Language Processing*. Pearson Prentice Hall.
 2. Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
 3. Aggarwal, C. C. (2016). *Recommender Systems: The Textbook*. Springer.
 4. Salton, G., Wong, A., & Yang, C. S. (1975). A vector space model for automatic indexing. *Communications of the ACM*, 18(11), 613–620.
 5. Ramos, J. (2003). Using TF-IDF to determine word relevance in document queries. In *Proceedings of the First Instructional Conference on Machine Learning*.
 6. Google News pre-trained Word2Vec model, available from Google Code Archive.
 7. Lops, P., De Gemmis, M., & Semeraro, G. (2011). Content-based recommender systems: State of the art and trends. In *Recommender Systems Handbook* (pp. 73–105). Springer.
 8. Koren, Y., Bell, R., & Volinsky, C. (2009). Matrix factorization techniques for recommender systems. *Computer*, 42(8), 30–37.
 9. Lovely Professional University (2024). Course materials for INT-423 Machine Learning II, School of Computer Science and Engineering.
-