

ReSearchMate

**An Intelligent Recommender System
for Research Papers and Insights**

DS5230

Isha Singh, Nisharg Gosai, Sudarshan Paranjape, Umang Jain



ReSearchMate

An Intelligent Recommender System for Research Papers and Insights

Course: DS5230

Team Members:

- Isha Singh
- Nisharg Gosai
- Sudarshan Paranjape
- Umang Jain

ReSearchMate aims to revolutionize the way researchers navigate and discover academic literature. By leveraging advanced clustering, topic modeling, and recommendation algorithms, it provides personalized insights and impactful paper recommendations tailored to individual research interests.

Contents		
01	02	03
Introduction Problem Statement & Key Objectives	Data Data Source, Attributes, Preprocessing	Clustering Overview Techniques Used & challenges faced
04	05	
Topic Modeling Methods & Visuals	Conclusions Results & Outcomes	

Let's Walk Through the Contents

1. Introduction

We'll start by understanding the **Problem Statement** and the **Key Objectives** of the project. This section outlines the challenges faced by researchers in navigating academic literature and introduces our solution, *ReSearchMate*.

2. Data

Next, we'll dive into the **Data Source**, its **Attributes**, and the preprocessing techniques applied to prepare the dataset for analysis. This includes steps like handling missing values, standardizing columns, and cleaning text for meaningful insights.

3. Clustering Overview

This section covers the **Clustering Techniques Used**, including methods like K-Means, along with the **Challenges Faced** in organizing papers into thematic groups.

4. Topic Modeling

We'll explore **Methods & Visuals** from advanced topic modeling techniques such as Latent Dirichlet Allocation (LDA). This section highlights how we captured nuanced research themes and interdisciplinary trends.

5. Conclusions

Finally, we'll present the **Results & Outcomes**, summarizing the impact of the project and discussing potential future directions for *ReSearchMate*.

Problem Statement



Overwhelming Volume of Research Papers

With millions of papers published annually across disciplines, finding relevant research is increasingly difficult for scholars and researchers.



Limited Tools for Thematic Discovery

Existing applications lack advanced clustering and topic modeling capabilities to uncover thematic trends and interdisciplinary connections effectively.



Need for Personalized Recommendations

Researchers struggle to identify impactful papers tailored to their specific interests and ongoing projects

Problem Statement

The sheer volume of research papers published annually across various disciplines poses a significant challenge for scholars and researchers. With millions of papers added to the growing body of academic literature each year, finding relevant and high-quality research has become an overwhelming task. This abundance of information often leads to frustration and inefficiency, as researchers spend valuable time sifting through irrelevant or less impactful studies.

Existing tools for academic search and discovery fall short in addressing this issue effectively. While they provide basic search functions, they often lack advanced clustering and topic modeling capabilities. These capabilities are crucial for uncovering thematic trends and interdisciplinary connections, which can provide fresh perspectives and foster innovative research. Without these features, researchers are left with limited insights into emerging fields and are unable to fully explore the interconnected nature of modern science.

Furthermore, people face difficulties in identifying impactful papers tailored to their specific interests and ongoing projects. Current recommendation systems often fail to deliver personalized suggestions that align with an individual's unique research focus. This gap leaves many scholars struggling to stay updated on critical developments in their areas of expertise, potentially missing key contributions that could influence their work. Addressing this need for personalized, theme-based, and interdisciplinary discovery tools is essential for advancing academic productivity and innovation.

Objective



Develop a Robust Clustering Framework

Use advanced clustering techniques (e.g., K-Means) to organize research papers into meaningful thematic groups / clusters



Implement Enhanced Topic Modeling

Transition to Latent Dirichlet Allocation (LDA) for capturing diverse and nuanced research topics, improving over LSA



Build a Research Paper Recommender System

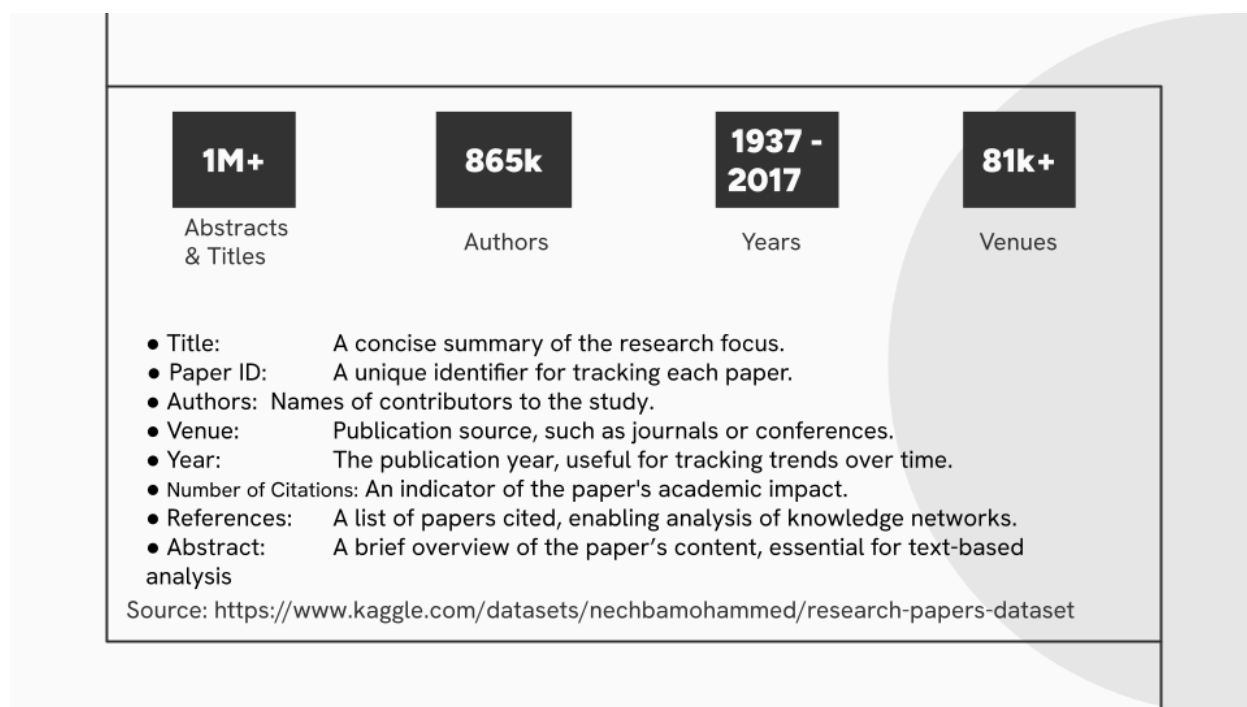
Leverage topic distributions and clustering results to create a smart tool for personalized and impactful paper recommendations.

Objective

The objective is to develop a robust clustering framework to streamline the organization of research papers into meaningful thematic groups. Advanced clustering techniques such as K-Means will be employed to categorize papers based on similarities, enabling researchers to quickly identify relevant studies within specific domains. This framework will address the challenge of navigating the overwhelming volume of academic literature by providing a structured and intuitive way to explore research.

To enhance the understanding of nuanced topics, the framework will implement Latent Dirichlet Allocation (LDA) for topic modeling, moving beyond traditional methods like Latent Semantic Analysis (LSA). LDA is better suited for capturing the diversity and depth of research topics, making it a valuable tool for uncovering interdisciplinary connections and emerging trends. This advanced topic modeling approach will provide researchers with deeper insights into thematic trends across large datasets.

The final component of the objective is to build a research paper recommender system that leverages the clustering and topic modeling results. By utilizing topic distributions and thematic groupings, the system will deliver personalized recommendations tailored to a researcher's interests and ongoing projects. This smart tool aims to enhance academic productivity by helping researchers discover impactful papers efficiently, fostering innovation and collaboration across disciplines.



Dataset

The dataset contains over 1 million research papers, with 865,000 unique authors and data spanning from 1937 to 2017. It includes over 81,000 abstracts and titles, providing a substantial foundation for analyzing academic trends and connections. The data offers a rich resource for exploring the evolution of research across decades.

Key metadata includes critical details such as the **title**, offering a concise summary of each paper's focus, and a **paper ID**, a unique identifier for easy tracking. Additional fields include **authors** (contributing researchers), the **venue** (publication source such as journals or conferences), and the **year** of publication, which is crucial for analyzing temporal trends. The dataset also tracks the **number of citations**, serving as a measure of each paper's academic impact, and **references**, which allow for network analysis of the knowledge base.

The **abstracts** and **titles** provide an essential textual foundation for advanced analyses, such as clustering and topic modeling. These elements enable researchers to extract thematic insights and develop tools like personalized recommender systems, making the dataset an invaluable resource for scholarly exploration and innovation.

Source :

<https://www.kaggle.com/datasets/nechbamohammed/research-papers-dataset>

Data Cleaning

1. Handling Missing Values:

Dropped rows where key features (abstract) was missing.
Retained rows with missing venue or references since they were less critical.

2. Column Standardization

Standardized text improved model input consistency.

3. Removed Non String Entities

Converted non-string abstract entries to an empty string for uniformity.

Data Cleaning

1. Handling Missing Values

- Rows with missing **abstracts** were dropped, as abstracts are crucial for text-based analysis.
- Rows with missing **venues** or **references** were retained since these features are less critical for the core tasks of clustering and topic modeling.

2. Column Standardization

- Text data was standardized to ensure consistency in model inputs. This included normalizing the case, removing extra spaces, and ensuring uniform formats for all text-based columns.

3. Removed Non-String Entities

- Non-string entries in the **abstract** column were converted to empty strings to maintain uniformity. This ensured that downstream processes could handle the data without errors or inconsistencies.

These cleaning steps aimed to optimize the dataset for clustering, topic modeling, and recommendation system development while preserving the integrity of essential information.

Data Preprocessing

1. **Lowercasing:**
Converted all text in abstract to lowercase.
2. **Removing Punctuation:**
Eliminated special characters to focus on meaningful text.
3. **Removing Numbers:**
Stripped numerical values to avoid bias in text patterns.
4. **Tokenization:**
Split text into individual words (tokens).
5. **Stopword Removal:**
Excluded common English stopwords (e.g., "the", "is", "and").
6. **Lemmatization:**
Reduced words to their base forms (e.g., "running" → "run").
7. **Vectorization:**
Using TF-IDF for title & abstracts and apply SVD

Data Preprocessing

1. **Lowercasing**
 - All text in the **abstracts** was converted to lowercase, ensuring uniformity and reducing redundancy during analysis (e.g., "Research" and "research" treated the same).
2. **Removing Punctuation and numbers**
 - Special characters and punctuation marks were eliminated to focus on meaningful textual content and improve tokenization accuracy.
 - Numerical values were stripped from the text to avoid biasing patterns based on irrelevant numeric sequences.
3. **Tokenization & Lemmatization**
 - Text was split into individual words (tokens), forming the basis for further text processing and analysis. And then were reduced to their base forms
4. **Stopword Removal**
 - Common English stopwords (e.g., "the," "is," "and") were excluded to retain only significant terms, improving topic modeling and clustering effectiveness.
5. **Vectorization**

Text data (titles and abstracts) was transformed into numerical representations using **TF-IDF** (Term Frequency-Inverse Document Frequency).

 - **Singular Value Decomposition (SVD)** was applied to reduce the dimensionality of TF-IDF vectors, retaining the most informative features for clustering and topic modeling.

These steps prepared the text data for advanced analysis by focusing on relevant features, minimizing noise, and optimizing computational efficiency

01

Clustering Overview

Clustering techniques were employed to analyze and categorize the research paper dataset based on **Title and Abstract features**

K-Means Clustering:

Description: K-means is a centroid-based clustering algorithm that partitions data into k clusters by minimizing the intra-cluster distance (variance).

Why K-Means?:

- Effective for large datasets.
- Provides interpretable results through defined cluster centroids.
- Computationally efficient compared to hierarchical methods.

Key Insights:

- K-means was effective in capturing distinct themes in the dataset when applied to both Title and Abstract features.
- The use of evaluation metrics such as WCSS, Silhouette Score, and Davies-Bouldin Index ensured the selection of an optimal k, maximizing clustering quality.

Agglomerative Clustering:

Description: Agglomerative clustering is a hierarchical technique that builds clusters through a bottom-up approach, iteratively merging the closest pairs of clusters based on a linkage criterion.

Why Agglomerative Clustering?:

- Provides a hierarchical view of the data, which is useful for understanding relationships between clusters.
- Capable of handling different cluster shapes and densities.

Key Insights:

- Agglomerative clustering complemented the K-means results by offering a hierarchical structure, allowing further exploration of sub-themes within the data.
- Different Linkage methods were experimented with to optimize results.

K-means clustering

- **Cluster Themes:** The identified themes in title clusters provide a high-level categorization of research papers, potentially useful for organizing large collections of academic literature.
- **Title clustering is more well-defined than abstract clustering, indicated by higher silhouette scores and lower Davies-Bouldin indices.**
 - For **title features**, the elbow and silhouette methods agree on $k=5$, suggesting five clusters effectively capture the structure.
 - For **abstract features**, $k=10$ is optimal per silhouette and Davies-Bouldin scores, indicating higher variability in abstract topics.
 - Top words in clusters highlight key research themes and their focus areas.
- **However, Abstract clustering may benefit from exploring higher numbers of clusters (>10), balancing performance gains against interpretability.**

K-means clustering was applied to the Title and Abstract separately, and experiments were run to find the optimal value of k for both the title and abstract. Experiments were run for **k values from 2 to 10** iteratively and below scores were recorded for each k value

- 1) K value
- 2) WCSS
- 3) Silhouette Score
- 4) Davies-Bouldin Score

A **custom function** was created that provides a structured approach to evaluate the optimal number of clusters using three widely recognized metrics: **Within-Cluster Sum of Squares (WCSS)**, **Silhouette Score**, and **Davies-Bouldin Index**.

Key Steps in the Function:

1. Elbow Method for WCSS:

- Calculate the first differences (deltas) and second differences (double_deltas) of the WCSS values.
- Identify the index where the second difference is maximized, corresponding to the elbow point.
- Set `elbow_k` as the number of clusters at this elbow point.

2. Silhouette Score:

- Find the index of the maximum Silhouette Score and retrieve the corresponding number of clusters.

3. Davies-Bouldin Index:

- Find the index of the minimum Davies-Bouldin Index and retrieve the corresponding number of clusters.

By combining these three metrics, the function provides a comprehensive way to evaluate the optimal number of clusters.

K-means clustering

Top Words in Title Clusters

- **Cluster 0:** Design-focused terms (e.g., "design", "implementation", "analysis").
- **Cluster 1:** Network-related terms (e.g., "network", "wireless", "sensor network").
- **Cluster 2:** Methodology emphasis (e.g., "method", "based", "optimization").
- **Cluster 3:** Application-oriented terms (e.g., "application", "model", "data").
- **Cluster 4:** General technical terms (e.g., "using", "algorithm", "approach").

Top Words in Abstract Clusters

1. **Cluster 0:** Robotics and control systems (e.g., "control", "robot", "motion").
2. **Cluster 1:** Networking and communication (e.g., "network", "protocol", "wireless").
3. **Cluster 2:** General methods and results (e.g., "method", "result", "problem").
4. **Cluster 3:** Data-centric topics (e.g., "data", "query", "gene").
5. **Cluster 4:** User experience and software (e.g., "user", "software", "design").
6. **Cluster 5:** Modeling and processes (e.g., "model", "process", "parameter").
7. **Cluster 6:** Communication channels and coding (e.g., "channel", "scheme", "receiver").
8. **Cluster 7:** Image processing and segmentation (e.g., "image", "object", "feature").
9. **Cluster 8:** Web services and architecture (e.g., "service", "web", "composition").
10. **Cluster 9:** Algorithms and optimization (e.g., "algorithm", "problem", "optimization").

The **output** of the custom function is below

Optimal clusters for Title:

{'elbow_k': 5, 'silhouette_k': 5, 'davies_bouldin_k': 2}.

Optimal clusters for Abstract:

{'elbow_k': 3, 'silhouette_k': 10, 'davies_bouldin_k': 10}

Based on the above output we decided to **select k=5 as the optimal k value for Title** as the elbow point was found to be at k=5 and the silhouette score was also the highest at k=5

For Abstract, k=10 as the optimal k value was selected as the silhouette score was the highest at k=10 and the Davies Bouldin score was the lowest at k=10.

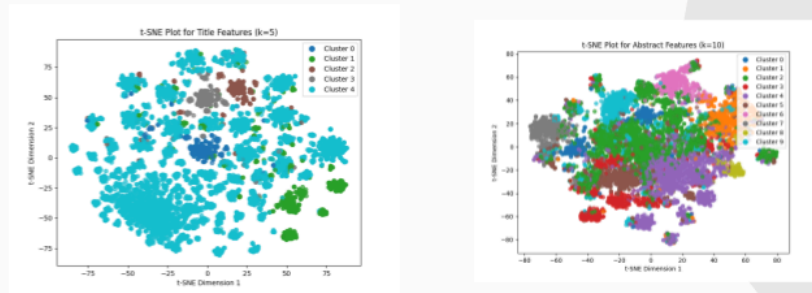
However, it is important to note that since the experiments were run for k values from 2 to 10, the k=10 value is the upper limit and therefore should be considered as an optimal value within the scope.

Future work could include experimenting with higher k values for abstract features.

The above slide shows how k-means captured themes in the data. **The identified themes in the title and abstract clusters provide a high-level categorization of research papers, useful for organizing large collections of academic literature.**

Results:

- Visualizations for both title and abstract clusters were generated using t-SNE, highlighting the distribution and separation of clusters in a reduced dimensional space for $k=5$ and $k=10$



The slide above shows two **t-SNE** (t-distributed Stochastic Neighbor Embedding) visualization plots representing title and abstract feature analyses:

Abstract Features Plot

The plot displays clustering results with $k=10$, showing a complex distribution of data points across 10 distinct clusters. The data points are spread across a 2D space ranging from approximately -80 to +80 on both dimensions. Notable characteristics include:

- Dense cluster formations in multiple regions
- Significant overlap between some clusters, particularly in the central region
- Clear separation of some clusters, especially those at the periphery

Title Features Plot

The plot shows an optimal clustering with $k=5$, revealing:

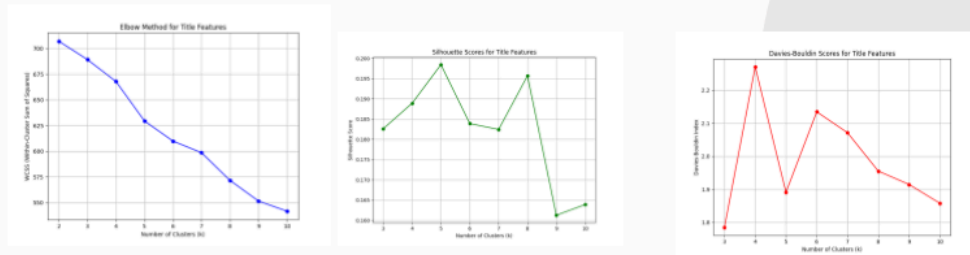
- More concentrated grouping patterns
- Dominant presence of turquoise-colored clusters (Cluster 4)
- Fewer overlapping regions compared to the abstract features
- Data points are distributed in a similar range but with different cluster density patterns
- Clear separation of the green clusters (Cluster 1) in the lower right region

Both visualizations demonstrate dimensionality reduction from high-dimensional space to 2D while preserving local relationships between data points. The different number of clusters ($k=10$ vs $k=5$) suggests that the title features might have more natural grouping tendencies compared to the more complex abstract features.

The plots represent abstracts showing more complex relationships requiring finer clustering granularity, while titles demonstrate more distinct and separable features.

Results:

- Elbow, Silhouette and DB score visualization for Title clustering experiments for k values from 2 to 10



In the above slide, the three plots show different evaluation metrics for determining the optimal number of clusters (k) for **Title features** in a clustering analysis:

Silhouette Score Analysis

The Silhouette score plot shows two notable peaks:

- First peak at k=5 with a score of ~0.198
- Second peak at k=8 with a score of ~0.195
- Sharp decline after k=8, dropping to ~0.161 at k=9

This suggests that either 5 or 8 clusters could be optimal configurations, with k=5 showing slightly better cohesion and separation.

Elbow Method Analysis

The Within-Cluster Sum of Squares (WCSS) plot shows:

- Steady decrease from k=2 (~705) to k=10 (~540)
- Notable "elbow" point around k=4-5
- Diminishing returns in reduction of WCSS after k=5

This indicates that k=5 might be a good balance between cluster complexity and explained variance.

Davies-Bouldin Score Analysis

The Davies-Bouldin index reveals:

- General downward trend after k=6
- Lowest score at k=10 (~1.85)

Lower Davies-Bouldin scores indicate better clustering, suggesting that higher k values provide better cluster separation, though this should be balanced against model complexity.

Considering all three metrics together, **k=5 appears to be the optimal choice** because of it:

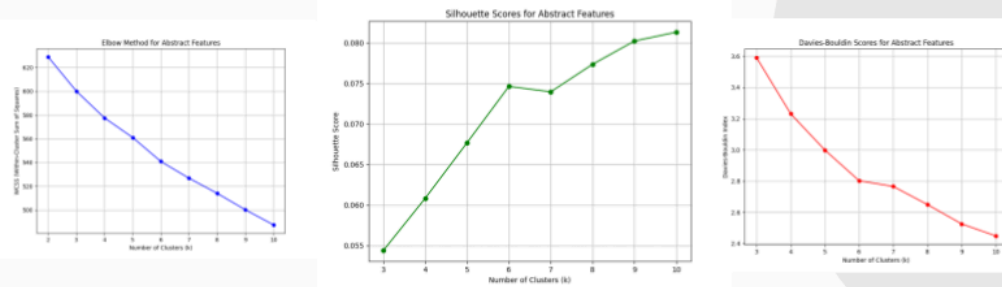
Achieves the highest Silhouette score

Corresponds to the elbow point in the WCSS plot

Offers a reasonable trade-off before the secondary peak in Davies-Bouldin scores

Results:

- Elbow, Silhouette and DB score visualization for Abstract clustering experiments for k values from 2 to 10



The above slide has three plots that show evaluation metrics for determining the optimal number of clusters (k) for **Abstract features** in a clustering analysis:

Elbow Method Analysis

The WCSS plot demonstrates:

- Steady decrease from k=2 (~630) to k=10 (~490)
- Gradual curve without a sharp elbow point
- The most significant drop occurs between k=2 and k=4
- Diminishing returns in reduction after k=6

Silhouette Score Analysis

The Silhouette score shows:

- Continuous upward trend from k=3 (0.054) to k=10 (0.081)
- Notable increase between k=5 and k=6
- Small plateau around k=6-7
- Highest score at k=10, suggesting better cluster separation with more clusters

Davies-Bouldin Score Analysis

The Davies-Bouldin index displays:

- Consistent downward trend from k=3 (3.6) to k=10 (2.4)
- The steepest decline between k=3 and k=5
- More gradual decrease after k=6
- Lower scores indicate better clustering separation

Unlike the title features analysis, the abstract features show a less clear-cut optimal k value. The metrics suggest:

Higher k values (8-10) might be more appropriate for abstract features

More complex clustering structure compared to titles

No single obvious optimal k value, but k=10 shows promising results across metrics

The trade-off between model complexity and cluster quality needs to be considered

Agglomerative Clustering

- Agglomerative Hierarchical Clustering was applied to SVD-reduced features of titles and abstracts with experiments run for k values of 2 to 10 and linkage methods ward, complete, average and single.
- Cluster Quality: Title clustering shows better-defined clusters compared to abstract clustering, as indicated by higher silhouette scores and lower Davies-Bouldin indices.
- Linkage Methods: Average and Single linkage methods produced the highest silhouette scores and lowest Davies-Bouldin indices, suggesting they may be more suitable for this dataset than the Ward method initially used.

Agglomerative clustering was applied to TF-IDF vectorized and SVD-reduced research paper titles and abstracts to identify thematic groupings.

Clustering performance was evaluated for cluster sizes (k=2 to k=10) using four linkage methods: **ward, complete, average, and single**.

The models were assessed using the **Silhouette Score** (measuring cluster cohesion and separation) and the **Davies-Bouldin Index** (evaluating compactness and separation).

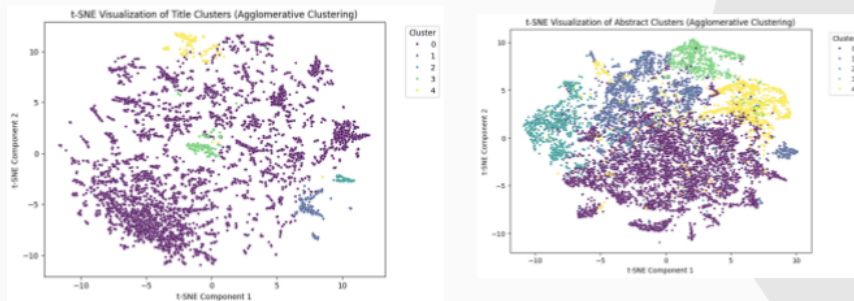
For Title features, the optimal configuration was achieved with **k=2** clusters using the **average linkage method**, yielding a Silhouette Score of 0.634 and a Davies-Bouldin Index of 0.254, indicating well-separated, compact clusters.

For Abstract features, the best Silhouette Score (0.400) was observed with k=2 and the **single linkage method**, while the lowest Davies-Bouldin Index (0.449) occurred at k=3, suggesting a trade-off between cohesion and compactness.

These results demonstrate agglomerative clustering's ability to capture meaningful patterns in the data and identify optimal configurations for grouping research titles and abstracts based on thematic content.

Results:

- Visualizations for both title and abstract clusters were generated using t-SNE, highlighting the distribution and separation of clusters in a reduced dimensional space for $k=5$ and ward linkage



The above slide shows two t-SNE visualization plots which show the results of agglomerative clustering applied to abstract and title features:

Abstract Features Clustering ($k=5$)

The plot reveals:

- Five distinct clusters with varying sizes and densities
- Large central purple cluster (Cluster 0) showing high density
- Well-separated yellow cluster (Cluster 4) on the right side
- Distinct turquoise cluster (Cluster 2) on the left
- Green cluster (Cluster 3) showing clear separation in the upper right
- Some overlap between clusters, particularly near the central region
- Data points spread across a range of approximately -10 to +10 in both dimensions

Title Features Clustering ($k=5$)

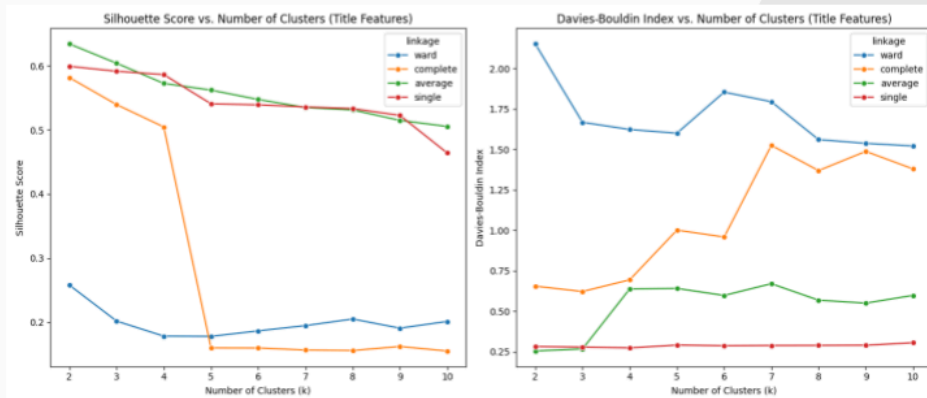
The plot demonstrates:

- More dispersed distribution with less defined boundaries
- Dominant purple cluster (Cluster 0) spread across multiple regions
- Small, isolated clusters for other groups (1-4)
- A higher degree of scatter compared to abstract clustering
- Less cohesive cluster formation
- Similar dimensional range (-10 to +10) as abstract features

The abstract features show more coherent clustering patterns with clearer boundaries between groups, while title features exhibit more fragmented clustering with scattered subgroups. This suggests that abstract content may contain more structured and distinguishable patterns compared to titles, which show higher variability and less distinct grouping tendencies.

Results:

- Silhouette and DB scores for different values of k and linkage methods



The above slide shows a comprehensive analysis of agglomerative clustering performance on title features for different values of k using multiple evaluation metrics:

Silhouette Score Analysis

The left plot shows Silhouette scores for different linkage methods:

- Average and single linkage methods perform best, starting at ~0.65 for k=2
- Complete linkage shows a dramatic drop after k=4
- Ward linkage consistently shows the lowest scores (~0.2)
- All methods show a general declining trend as k increases

Davies-Bouldin Index Analysis

The right plot reveals:

- Ward linkage starts highest (~2.1) and gradually decreases
- Complete linkage shows a significant increase after k=4
- Single linkage remains most stable (~0.3)
- Average linkage maintains moderate scores throughout

The metrics suggest different optimal cluster numbers depending on the method:

Silhouette scores favor lower k values (2-4)

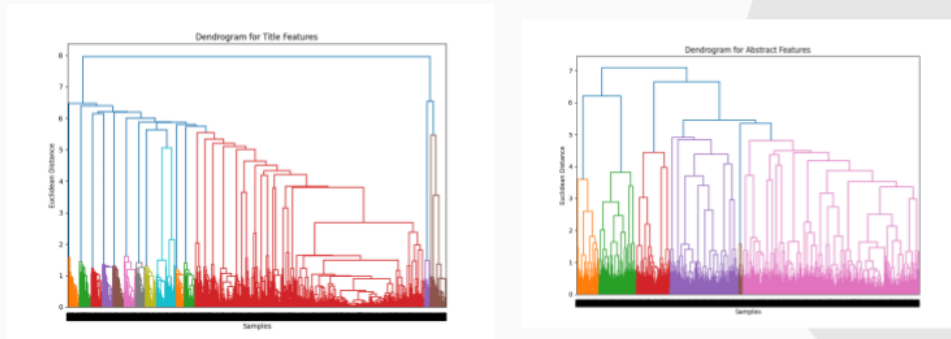
Davies-Bouldin index suggests better separation at higher k values for ward linkage

The elbow plot indicates diminishing returns after k=6

Single and average linkage methods generally outperform ward and complete linkage

Results:

- Dendrogram visualizations for title and abstract features



In the above slide, the two dendrograms reveal distinct clustering patterns between title and abstract features.

Title features exhibit a more complex branching pattern with multiple small clusters merging at various heights, indicating a more nuanced relationship between features. In contrast, the abstract features show more distinct, cohesive groupings with clearer separation between major clusters.

Exploration of Other Techniques

Other clustering methods were also evaluated but were found to have limitations in this specific context:

DBSCAN (Density-Based Spatial Clustering of Applications with Noise):

Strengths: DBSCAN is robust to outliers and works well for clusters of varying densities.

Limitations:

- Struggled with the high dimensionality of the TF-IDF vectorized data.
- Required extensive parameter tuning (e.g., epsilon and min_samples), and results were inconsistent across different subsets of the data.
- Failed to identify meaningful clusters due to the sparse nature of the dataset.

Spectral Clustering:

Strengths: Suitable for data with complex, non-convex cluster structures.

Limitations:

- Computationally intensive for large datasets due to eigenvalue decomposition.
- Did not perform well after dimensionality reduction, as it relies heavily on the affinity matrix, which can degrade in high-dimensional spaces.
- Clusters formed needed to be more well-defined and interpretable for the given features.

Topic Modeling

Topic modeling is an unsupervised machine-learning technique used to extract hidden themes or topics from large collections of textual data. This approach is widely used in applications such as information retrieval, document classification, and recommendation systems.

By grouping similar documents under common topics, topic modeling enables efficient navigation and personalized content recommendations, making it a cornerstone for text-based applications. All topic models are based on the same basic assumption:

- Each document consists of a mixture of topics, and
- Each topic consists of a collection of words.

In other words, topic models are built around the idea that the semantics of our document are governed by some hidden, or “latent,” variables that we are not observing. As a result, the goal of topic modeling is to uncover these latent variables — topics — that shape the meaning of our document and corpus. In this project, we discover 2 types of topic models - LDA and LSA.

Why LSA and LDA?

- LSA was chosen for its efficiency in grouping papers into initial clusters.
- LDA was used as a follow-up to refine and diversify the topic distributions for better recommendations.

By combining LSA and LDA, we leveraged their respective advantages to achieve a balance between computational efficiency and topic modeling accuracy.

Latent Semantic Analysis (LSA)

- LSA is a technique that captures latent relationships between terms and documents.
- The core idea is to take a matrix of what we have — documents and terms — and decompose it into a separate **document-topic matrix** and a topic-term matrix.
- LSA uses **Singular Value Decomposition (SVD)** for dimensionality reduction, identifying hidden relationships between documents and terms.

The diagram shows the formula $w_{i,j} = t f_{i,j} \times \log \frac{N}{d f_j}$ with arrows pointing to its components: $t f_{i,j}$ is labeled 'tf-idf score' (red), N is labeled '# total documents' (blue), and $d f_j$ is labeled '# documents containing word' (purple). A green label '# occurrences of term in document' points to $f_{i,j}$.

LSA Overview

Latent Semantic Analysis (LSA) is a linear algebra-based method for topic modeling that leverages Singular Value Decomposition (SVD). It reduces high-dimensional text data into a lower-dimensional latent space. It identifies patterns in the matrix to discover hidden structures and relationships between terms and documents.

Characteristics:

- Efficient for dimensionality reduction and initial exploration of themes.
- Provides interpretable results by ranking words and documents based on their contributions to topics.

The core idea is to take a matrix of what we have — documents and terms — and decompose it into a separate document-topic matrix and a topic-term matrix.

LSA is particularly useful for dimensionality reduction, clustering, and improving the interpretability of text data. However, its reliance on linear relationships and inability to capture contextual nuances limit its applicability to complex semantic data.

Steps involved in LSA

Step 1: Document-Term Matrix

Constructed a matrix where rows represent documents and columns represent terms, weighted using TF-IDF.

Step 2: Dimensionality Reduction

Used truncated Singular Value Decomposition (SVD) to reduce the matrix and uncover latent topics.

Step 3: Topic Extraction & Visualization

Selected the top 5 topics by analyzing singular values and corresponding term distributions.

LSA was used to extract topics and organize papers thematically. Here's how it was applied:

Step 1) Document-Term matrix:

The first step is generating our document-term matrix. Given m documents and n words in our vocabulary, we can construct an $m \times n$ matrix A in which each row represents a document and each column represents a word. In the simplest version of LSA, each entry can simply be a raw count of the number of times the j -th word appeared in the i -th document. In practice, however, raw counts do not work particularly well because they do not account for the significance of each word in the document. Hence, our LSA model replaces raw counts in the document-term matrix with a tf-idf score. **Tf-idf, or term frequency-inverse document frequency**, assigns a weight for term j in document I as shown:

$$w_{i,j} = \underset{\substack{\uparrow \\ \text{tf-idf score}}}{tf_{i,j}} \times \log \frac{\underset{\substack{\uparrow \\ \text{\# documents containing word}}}{N}}{\underset{\substack{\uparrow \\ \text{\# total documents}}}{df_j}}$$

The diagram illustrates the components of the tf-idf formula. The term $tf_{i,j}$ is labeled as the 'tf-idf score' with an upward arrow. The term N is labeled as the '# total documents' with an upward arrow. The term df_j is labeled as the '# documents containing word' with an upward arrow. The formula shows that the weight $w_{i,j}$ is the product of the term frequency in the document and the inverse of the document frequency.

Step 2) Dimensionality Reduction:

Since our document matrix A is very sparse, dimensionality reduction can be performed using truncated SVD. **SVD, or singular value decomposition**, is a technique in linear algebra that factorizes any matrix M into the product of 3 separate matrices:

- U : Document-topic matrix.
- S : Singular values representing topic importance.
- V : Term-topic matrix.

Step 3) Topic Extraction and Visualization:

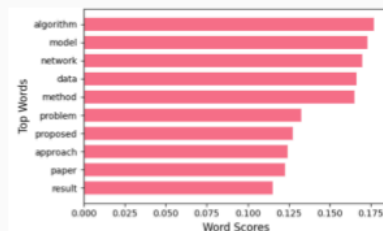
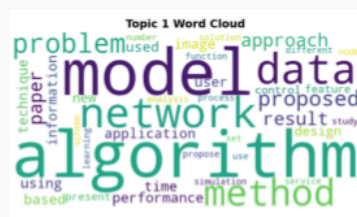
The final step is to select the top 5 topics by analyzing singular values and corresponding term distributions and visualize each topic in detail using word clouds and bar charts.

LSA Results & Visuals

- Extracted 5 dominant topics with corresponding top keywords.

	Topic	Count	Name	Representation	Representative Docs
0	0	778432	0_algorithm_model_network_data_method	[algorithm, model, network, data, method]	[study connection orthogonal polynomial reproducibility]
1	1	18784	1_network_service_node_protocol_wireless	[network, service, node, protocol, wireless]	[mobile communication network need public key cryptography]
2	2	7782	2_algorithm_channel_network_scheme_proposed	[algorithm, channel, network, scheme, proposed]	[paper investigate signal-to-noise ratio based on channel capacity]
3	3	13774	3_image_network_node_feature_sensor	[image, network, node, feature, sensor]	[algorithm represents hdr image using dual exposure]
4	4	8761	4_algorithm_network_problem_graph_node	[algorithm, network, problem, graph, node]	[let h undirected graph list homomorphism problem]

- Created word clouds and topic-word bar charts to visualize relationships.



Here is a summary of the results obtained by performing LSA:

- 1. Summary Dataframe:** Provides an overview of the topics identified through LSA.
 - **Topics:** Sequentially numbered unique identifiers (e.g., Topic 0, Topic 1).
 - **Count:** Number of documents associated with each topic.
 - **Name:** Top 5 dominant words defining the core theme of each topic.
 - **Representation:** A broader set of words providing context about the topic.
 - **Representative Documents:** Sample documents illustrating the content linked to each topic.
- 2. Word Clouds:** The word clouds help visualize the most dominant terms in each topic, showcasing the keywords that define the central themes of the respective topics. Larger words in the cloud correspond to terms that are more strongly associated with the respective topic.
- 3. Bar Chart:** The bar charts represent the top words associated with each of the five identified topics in the dataset. Each chart highlights the most significant terms that define a topic, based on their scores derived from the term-topic matrix. The height of each bar reflects the 'weight' or importance of a word in describing the respective topic. A higher score indicates that the word contributes more to the semantic meaning of the topic. For example,
 - **Topic 1:** Dominated by terms like 'algorithm,' 'model,' and 'network.' Suggests a focus on computational methods or systems.
 - **Topic 2:** Includes words like 'network,' 'service,' and 'protocol.' Indicates this topic revolves around communication technologies or services.

Issue with LSA

Topic 0 has a significantly higher count (778,432) compared to the other topics, making it dominant. This happens due to the following:

- 1) Sensitivity to Frequent Terms:** LSA overemphasizes terms that appear frequently across the corpus, leading to dominant, broad topics.
- 2) Lack of Contextual Understanding:** LSA doesn't distinguish between different meanings of words, causing generalized topics.
- 3) Dimensionality Reduction Smoothing:** SVD reduces complexity, but it may blur finer topic distinctions, favoring broad topics over specific ones.

To address these issues, we used more advanced methods like LDA (Latent Dirichlet Allocation) which handles topics in a probabilistic way.

Topic Dominance: Topic 0, with a count of 778432, is significantly larger than other topics, highlighting its prominence in the dataset. The dominance of Topic 0 in LSA can be attributed to some inherent drawbacks of the LSA method:

1. Sensitivity to Frequent Terms (High Term Frequency):

LSA is heavily influenced by the frequency of terms across the entire corpus. In cases where certain terms like "algorithm," "network," and "data" appear frequently, they dominate the singular value decomposition (SVD) process, leading to a topic that is overly broad and dominant. This is because LSA relies on singular values that capture global patterns, and frequently occurring terms often dominate these patterns, overshadowing more specific topics.

2. Lack of Semantic Context (No Word Sense Disambiguation):

LSA does not account for multiple meanings of a word (polysemy). For example, the word "network" can refer to various concepts such as computer networks, social networks, or biological networks. LSA treats these words as a single, global term without distinguishing between their different senses, which can result in overly general topics that do not capture domain-specific nuances.

3. Dimensionality Reduction Limitations:

In LSA, dimensionality reduction via SVD is used to find latent structures in the data. While this is effective in identifying broad patterns, it can also smooth out finer, domain-specific distinctions. As a result, general topics that have broad term distributions (like Topic 0) can dominate, while niche or specialized topics with more concentrated term distributions may fail to emerge clearly.

Topic Modeling Using LDA

- **Latent Dirichlet Allocation (LDA)** is a generative probabilistic model used for topic modeling.
- It identifies hidden topics within a corpus of text documents and associates words with these topics.
- LDA assigns the most likely topic to each document, creating a **document-topic mapping**.

The use of topic modeling and Latent Dirichlet Allocation (LDA) was driven by the objective of uncovering latent topics in a collection of document abstracts and segmenting them based on shared topics. This approach provides valuable insights into the hidden motifs and trends within the data.

LDA was chosen for several reasons:

1. **Probabilistic Modeling:**
LDA analyzes the words in each document and assigns probabilities to various topics, enabling the automatic discovery of themes without requiring predefined labels.
2. **Soft Clustering:**
This method allows documents to belong to multiple topics, offering a more flexible and precise understanding of the data.
3. **Unsupervised Technique:**
LDA does not rely on prior knowledge of the topics. The model independently deduces the topics from the dataset.

By employing the LDA technique, documents were organized into logical groups, revealing underlying topics. This automated approach yielded insights that would be speculative or challenging to achieve through manual analysis.

Working of LDA

Input: A document-term matrix, where rows represent documents, and columns represent word counts or frequencies.

Process:

- Assign each word in the document to a random topic.
- Iteratively refine the assignment of words to topics using a process like Gibbs Sampling, based on the likelihood of word-topic and topic-document distributions.

Output:

- Topic distributions for each document.
- Word distributions for each topic.

The functioning of Latent Dirichlet Allocation (LDA) can be described in three major steps:

1. **Input:**

LDA begins with a document-term matrix, where each row represents a document, and each column represents the frequency of words within those documents.

2. **Process:**

Initially, each word in the document is assigned to a topic at random. The algorithm then refines these initial assignments using Gibbs Sampling, a Bayesian technique. This iterative process updates the word-topic and topic-document assignments along with their probabilities, improving the results over time.

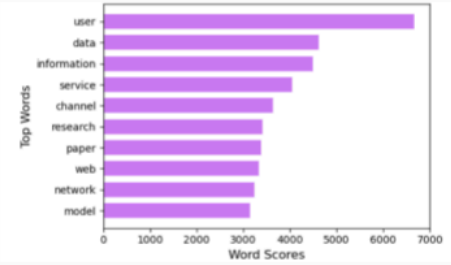
3. **Output:**

After several iterations, the output consists of:

- **Topic Distributions:** Indicating the proportion of each topic within every document.
- **Word Distributions:** Representing the probabilities of each word occurring within a specific topic.

In summary, LDA unveils the underlying themes in the text and establishes connections between words and topics.

Topic	Count	Name	Representation	Representative_Docs
0	0	196666_0_network_design_control_performance_algorithm	[network, design, control, performance, algo...	[paper proposes novel methodology fpga trust z...
1	1	67378_1_network_node_graph_protocol_problem	[network, node, graph, protocol, problem]	[car equipped gps wi fi transmitter becomes ea...
2	2	230980_2_method_algorithm_image_model_problem	[method, algorithm, image, model, problem]	[investigate approximation ability multilayer ...
3	3	134563_3_robot_method_image_model_object	[robot, method, image, model, object]	[backgroundinthe immune system multifaceted st...
4	4	197892_4_user_data_information_service_channel	[user, data, information, service, channel]	[aim research explore use crowdsourcing espec...



- **Table:** A table was created for each topic, including the count of documents, keywords describing the theme, and representative documents which are a sample of documents that best exemplify each topic.
- **Word Cloud:** The **word cloud** is a visual representation of the most frequent keywords within **Topic 5**. Larger words indicate higher frequencies, emphasizing their importance in defining the topic. Key terms such as **"user," "data,"** and **"information"** are dominant, highlighting their central role in this topic. Other significant words like **"technology," "service,"** and **"research"** suggest sub-themes or related contexts, giving a broad view of the topic's focus.
- **Bar Graph:** The **bar plot** displays the **top word scores** for Topic 5, reflecting the relevance of individual words to the topic. The **x-axis** represents the word scores, which indicate how strongly each word is associated with the topic. The **y-axis** lists the top words, such as **"user," "data," "information,"** and others. This quantitative representation complements the word cloud by providing a measurable perspective on keyword importance.

LDA Over LSA?

Topic Distribution:

- LDA: Balanced topics (67,378-230,980).
- LSA: Uneven, one dominant topic (778,432).

Word Coherence:

- LDA: Coherent and theme-aligned.
- LSA: Generalized and overlapping.

Interpretability:

- LDA: Clear and interpretable.
- LSA: Lacks clear structure.

The comparison between Latent Dirichlet Allocation (LDA) and Latent Semantic Analysis (LSA) highlights their effectiveness as follows:

Balance of Topics:

- **LDA:** Demonstrates a uniform distribution of topics, with documents ranging in size from approximately 67,000 to 231,000, ensuring balanced topic representation.
- **LSA:** Displays a skewed distribution, with one dominant topic (~778,000 documents) and other significantly smaller topics (~7,800 documents), resulting in an imbalance.

Word Representations:

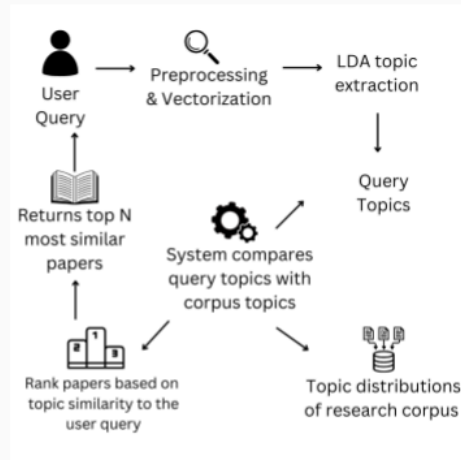
- **LDA:** Provides clear and expressive relationships between words and topics. Words are associated with distinct themes and exhibit minimal overlap, improving interpretability and topic coherence.
- **LSA:** Produces non-specific associations of words, leading to overlapping themes and reduced clarity, making it harder to differentiate between topics.

Conclusion:

LDA offers superior topic balance and clearer topic representations, making it more effective for clustering and discriminating topics in text data. These attributes ensure better visualization and interpretability compared to LSA.

Recommendation System

How it works?



Results

User Query = "image segmentation in neural networks"

```
Top Recommendations:
Title: Subband image segmentation using VQ for content-based image retrieval
Similarity Score: 0.8453
Abstract: Retrieving images from a large image dataset using image content as a ke
Authors: ['Junchul Chun', 'George C. Stockman']
Venue: acm multimedia | Year: 2001

Title: Image Clustering Using Color and Texture
Similarity Score: 0.8338
Abstract: With the advancement in image capturing device, the image data been gene
Authors: ['Manish Maheshwari', 'Sanjay Silakari', 'Mahesh Motwani']
Venue: computational intelligence | Year: 2009

Title: Block-oriented image decomposition and retrieval in image database systems
Similarity Score: 0.8128
Abstract: We investigate approaches to support effective and efficient retrieval o
Authors: ['Edward Remias', 'Gholamhosein Sheikholeslami', 'Aidong Zhang']
Venue: nan | Year: 1996
```

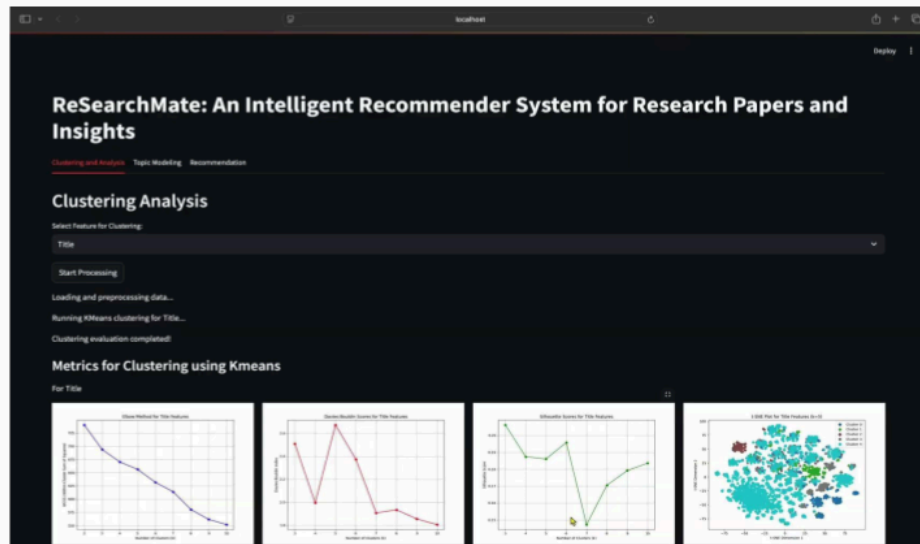
A recommendation system was built on top of the extracted LDA topics to suggest relevant research papers based on user input.

Process:

1. **User Query:** Users input a query, such as keywords or a paper title.
2. **Vectorization and Topic Extraction:** The query is preprocessed, vectorized using TF-IDF, and mapped into the LDA topic space.
3. **Similarity Search:** The system compares the user's topic vector to the topic vectors of the research papers using a FAISS index which uses cosine similarity for fast similarity computation.
4. **Recommendation Generation:** Papers with the highest similarity scores are ranked and returned to the user.
5. **Outputs:** Recommendations include paper titles, abstracts, authors, and similarity scores.

This approach leverages LDA's ability to capture diverse and nuanced research topics, potentially improving upon the LSA implementation by providing more accurate and contextually relevant recommendations.

Demo



The demo is the project in action. It shows the results from Clustering Analysis, Topic Modeling and personalized recommendation system. Here is the link to the demo:

[Link to demo](#) (at Minute 8:47)

Conclusion

- **Clustering Results:** Grouped research papers into meaningful clusters, highlighting dominant themes and emerging research areas.
- **Topic Modeling Insights:** Extracted balanced and coherent topics using LDA, overcoming limitations of LSA and providing actionable insights into academic trends.
- **Personalized Recommendations:** Enabled researchers to identify impactful papers aligned with their interests, enhancing productivity and interdisciplinary exploration

To summarize the project report:

Clustering Results:

The clustering process successfully grouped research papers into distinct, meaningful clusters based on thematic similarities. These clusters highlighted dominant themes in the dataset, such as advancements in algorithms, sensor networks, or image processing. The use of techniques like t-SNE for dimensionality reduction and K-Means for clustering ensured that the papers were visually and semantically well-organized.

Topic Modeling Results:

By transitioning from Latent Semantic Analysis (LSA) to Latent Dirichlet Allocation (LDA), the topic modeling process overcame key challenges, improving the representation of overlapping and nuanced themes while avoiding the dominance of generic topics often seen with LSA. LDA's probabilistic approach resulted in higher coherence scores and more balanced topics, offering a clearer understanding of the corpus.

Personalized Recommendations:

The recommendation system, built on top of the topic modeling outputs, offered personalized paper suggestions to researchers. By analyzing the topics most relevant to a user's query or interests, the system identified papers that were closely aligned with their research focus. This personalized approach not only saves time but also empowers researchers to explore impactful papers across disciplines.

Overall, the system enhanced productivity by tailoring recommendations to the user's academic needs and providing seamless access to relevant research.