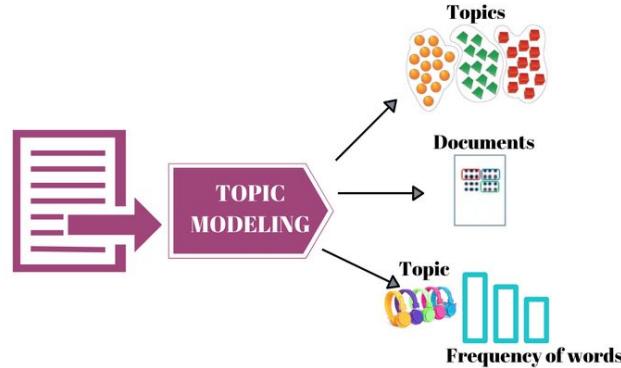


01 August 2023

# Exploring Bangla News Trends: A Comprehensive Topic Modeling Analysis on The Daily Star Bangla News Portal

[https://colab.research.google.com/drive/1DJN8ja2G3BqyL5LLai2E\\_Gw9Jd5y2pDo?usp=sharing](https://colab.research.google.com/drive/1DJN8ja2G3BqyL5LLai2E_Gw9Jd5y2pDo?usp=sharing)



Kazi Shadman Sakib  
4th Year Undergraduate Student

Department of Computer Science & Engineering  
University of Dhaka



# Quick Summary

Conducted an in-depth analysis of the Daily Star Bangla news portal, aiming to gain insights into the most frequent topics discussed over a span of seven days. Leveraging the **Latent Dirichlet Allocation (LDA)** technique, performed topic modeling on three distinct datasets: **daily news**, **peak user hour news (9 PM - 11 PM)**, and **off-peak user hour news (4 AM - 6 AM)**. Additionally, merged the data from all seven days to perform topic modeling on a weekly basis.

## Tools and Technologies :

- Python
- Requests
- BeautifulSoup
- Web Scraping
- Gensim
- NLTK (Natural Language Toolkit)
- BNLP (Bengali Natural Language Processing Toolkit)
- LDA (Latent Dirichlet Allocation)
- PyLDAvis
- Matplotlib
- Googletrans
- Google Colab

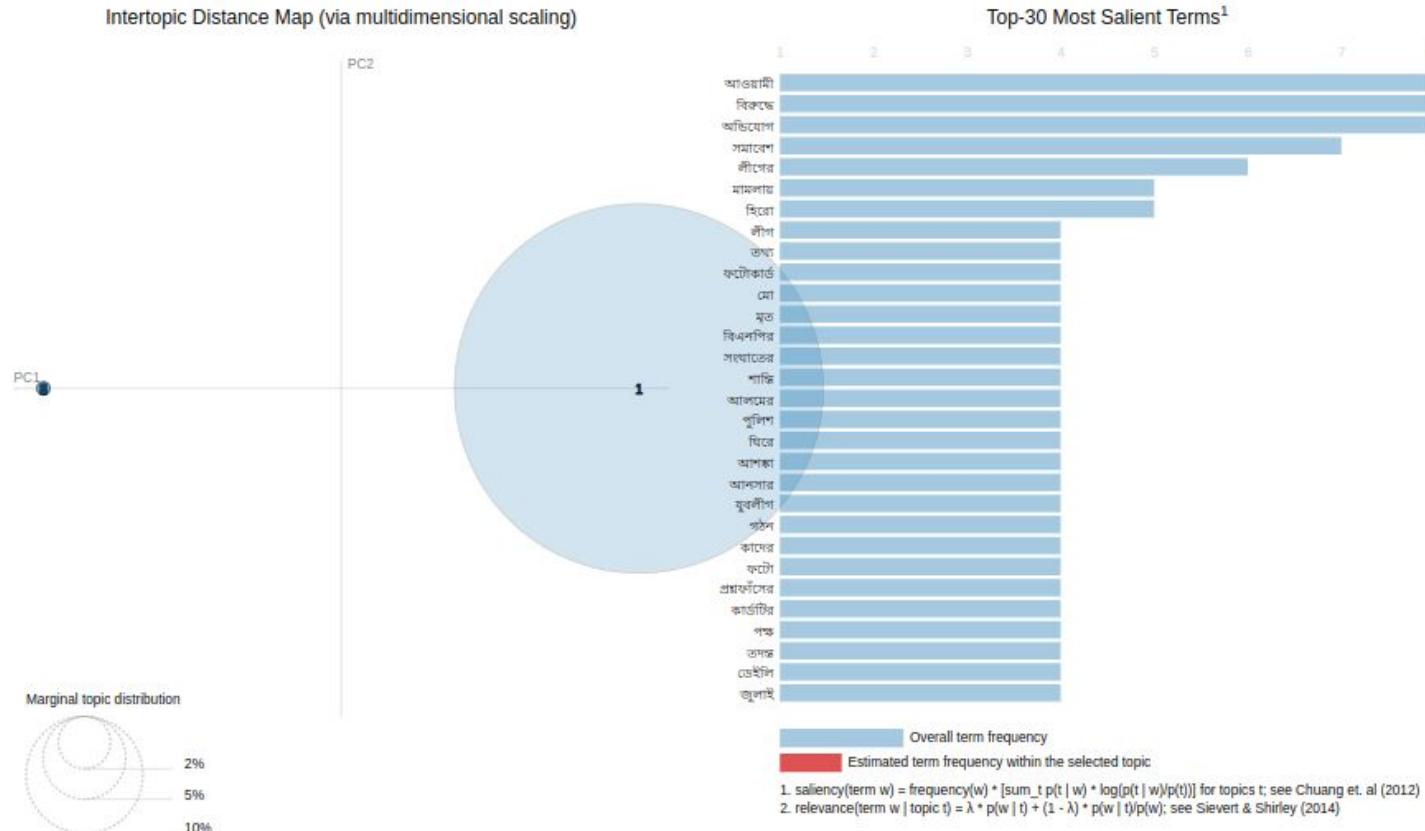
# Daily News

Kazi Shadman Sakib  
4th Year Undergraduate Student

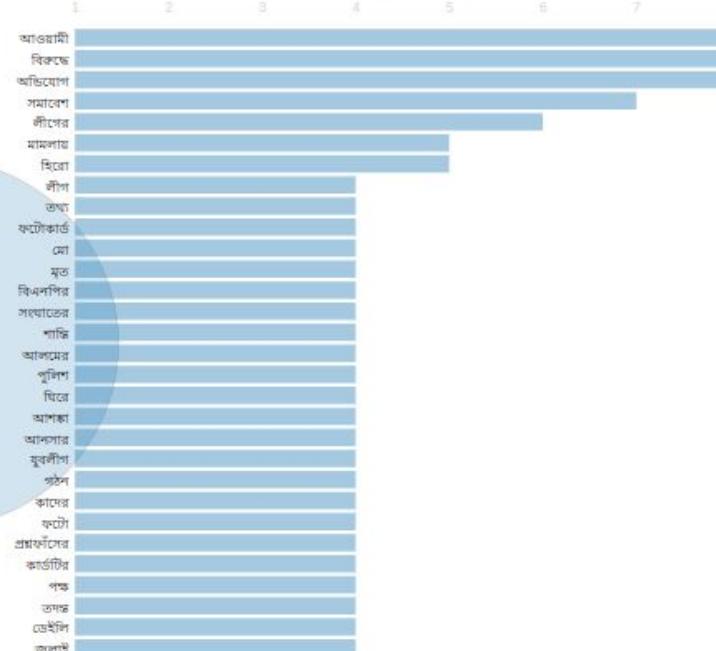
Department of Computer Science & Engineering  
University of Dhaka



# What was the trending topic on July 24, 2023?



Top-30 Most Salient Terms<sup>1</sup>



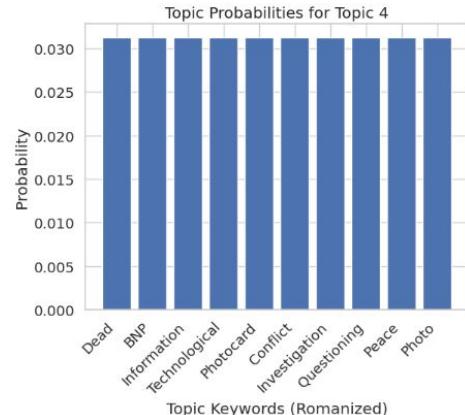
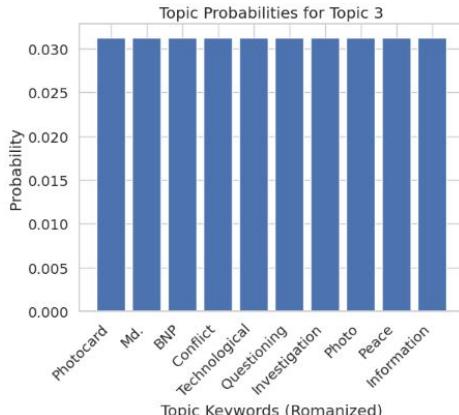
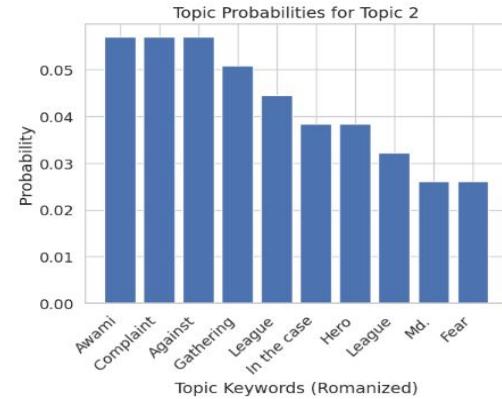
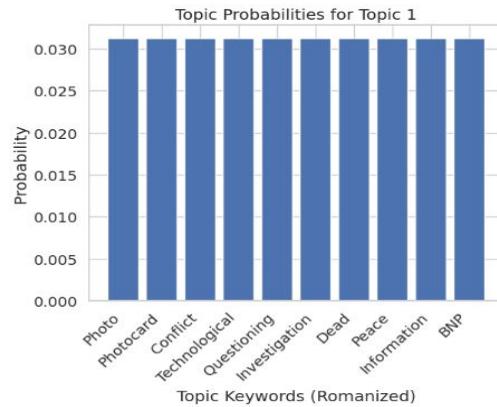
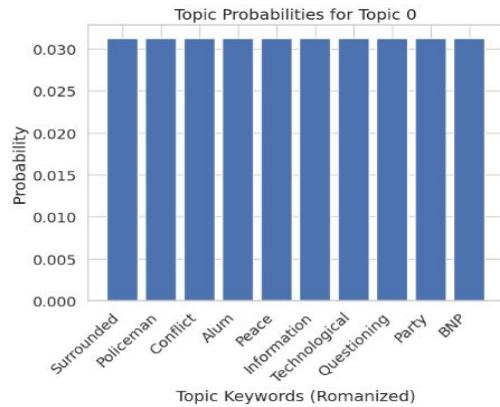
Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)p(t/w)$ ; see Sievert & Shirley (2014)

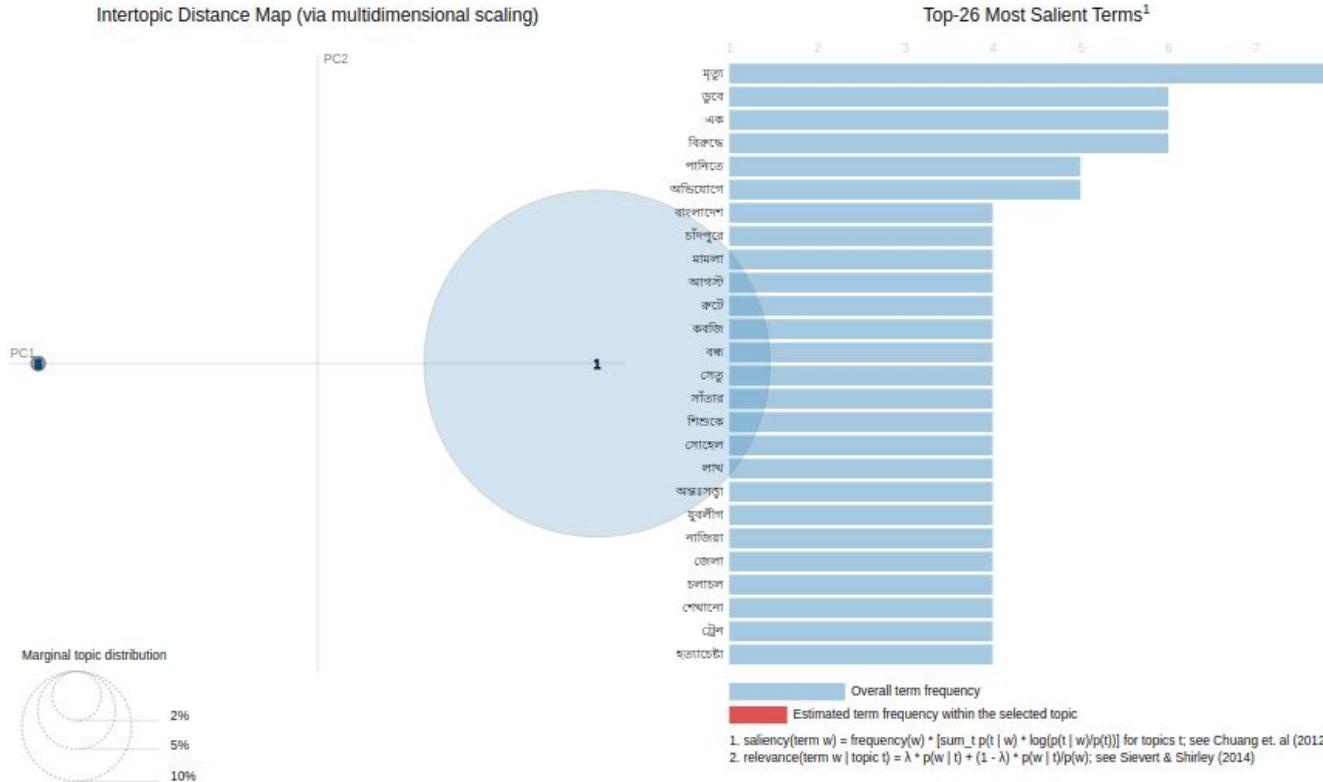
# What was the trending topic on July 24, 2023? (Continued)



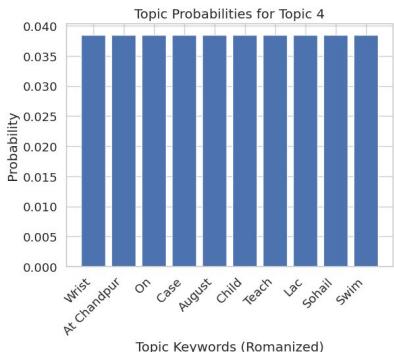
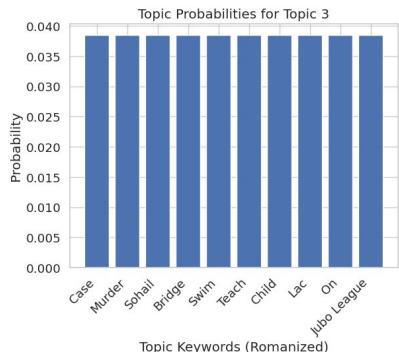
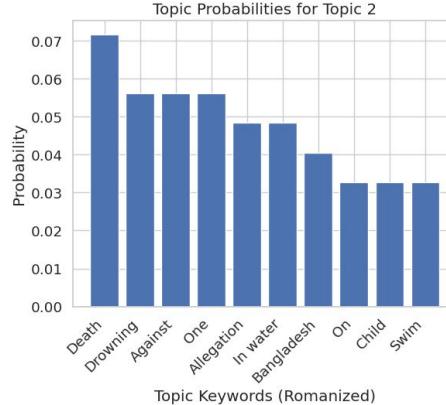
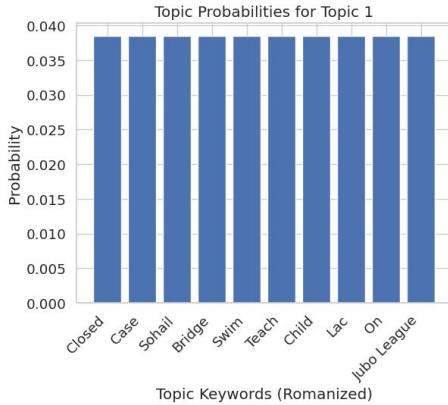
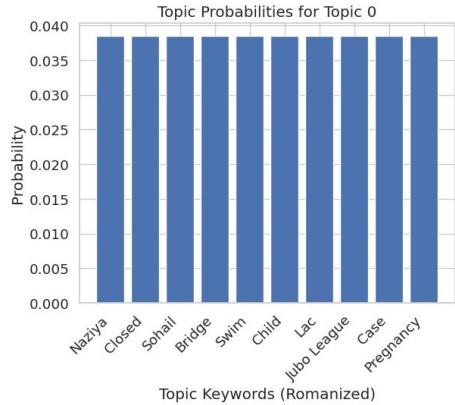
Most trending topic of the day is  
Topic: 2 with probabilities :

$$0.057 * \text{"আওয়ামী"} + 0.057 * \text{"অভিযোগ"} + \\ 0.057 * \text{"বিক্রিদে"} + 0.051 * \text{"সমাবেশ"} + \\ 0.045 * \text{"লীগের"} + 0.038 * \text{"মামলায়"} + \\ 0.038 * \text{"হিরো"} + 0.032 * \text{"লীগ"} + 0.026 * \text{"} \\ মো" + 0.026 * \text{"আশঙ্কা"}$$

# What was the trending topic on July 25, 2023?



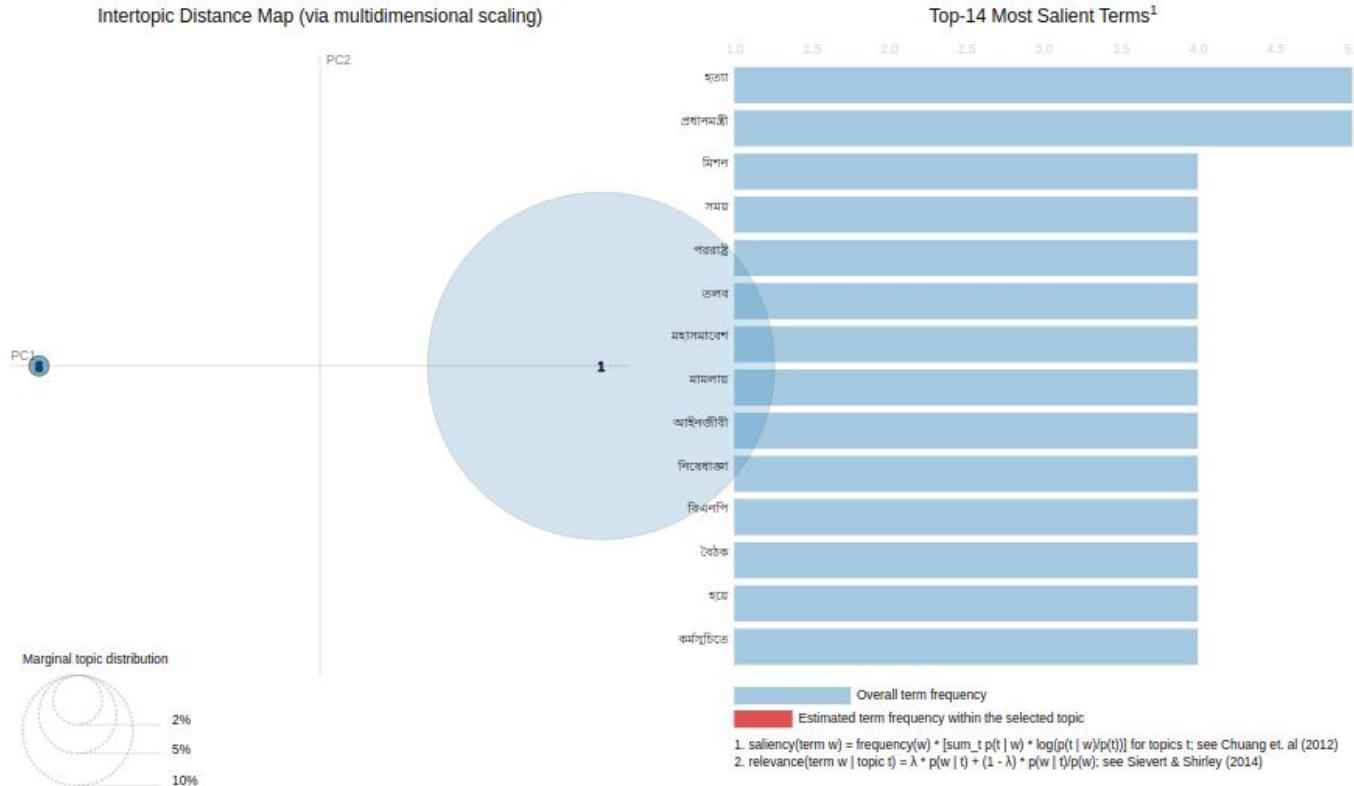
# What was the trending topic on July 25, 2023? (Continued)



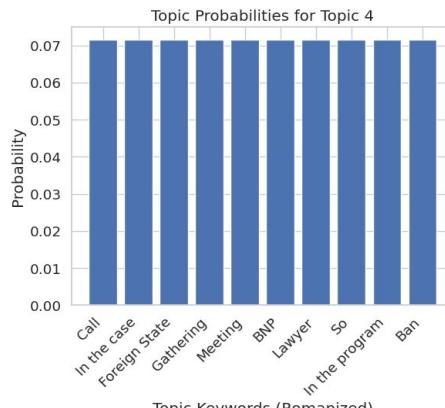
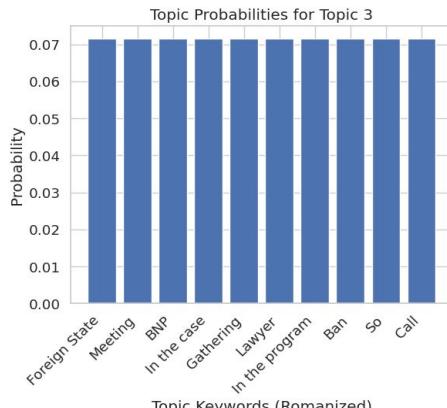
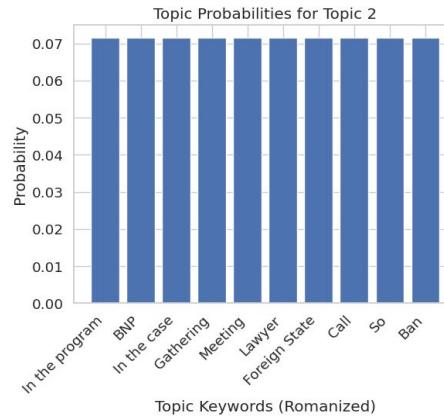
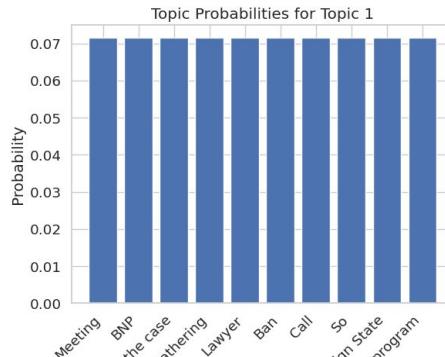
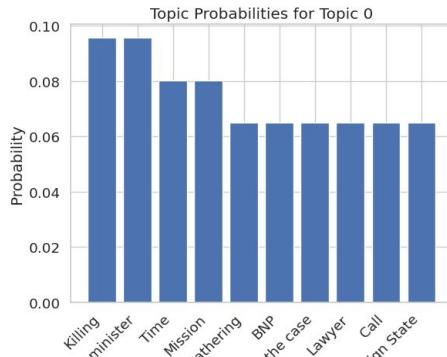
Most trending topic of the day is  
Topic: 2 with probabilities :

$$0.072 * \text{“মৃত্যু”} + 0.056 * \text{“ডুর্বে”} + 0.056 * \text{“বিরক্তে”} + 0.056 * \text{“এক”} + 0.048 * \text{“অভিযোগে”} + 0.048 * \text{“পানিতে”} + 0.041 * \text{“বাংলাদেশ”} + 0.033 * \text{“কটে”} + 0.033 * \text{“শিশুকে”} + 0.033 * \text{“সাঁতার”}$$

# What was the trending topic on July 26, 2023?



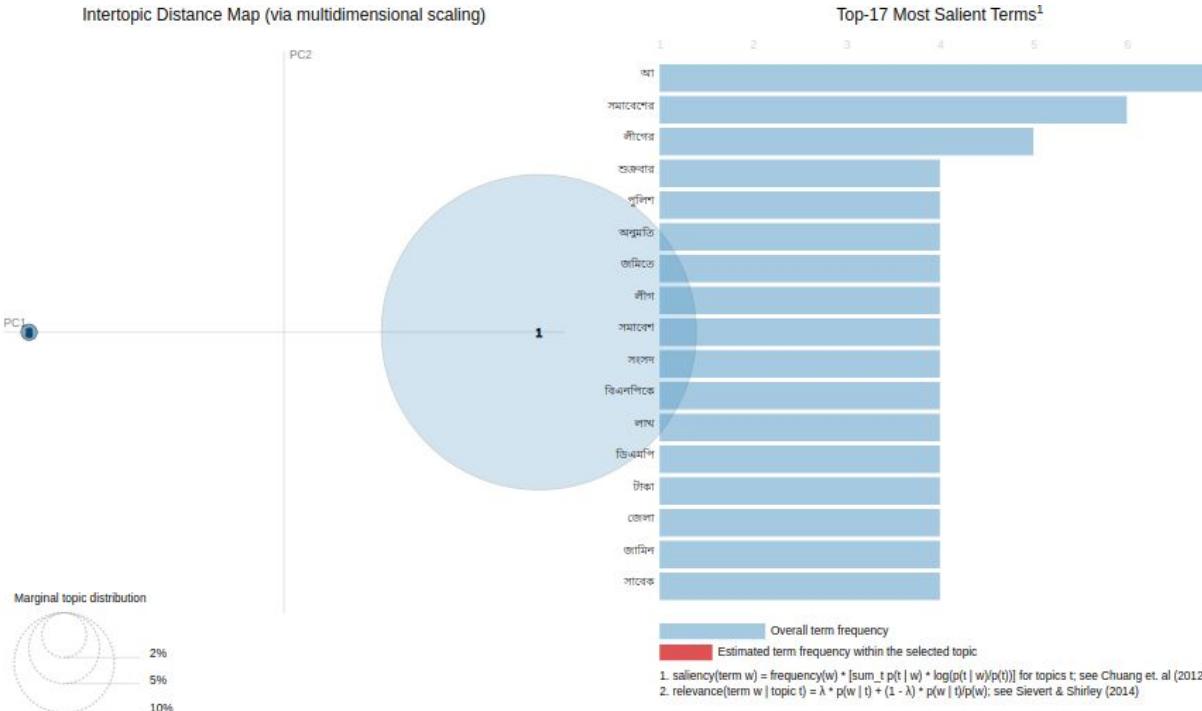
# What was the trending topic on July 26, 2023? (Continued)



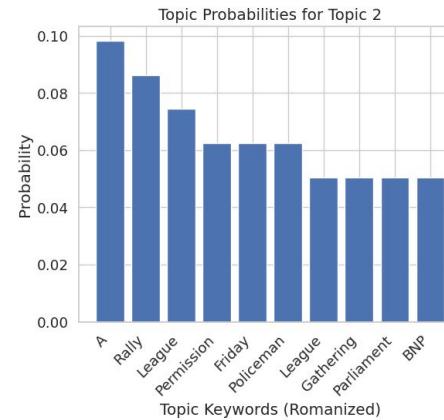
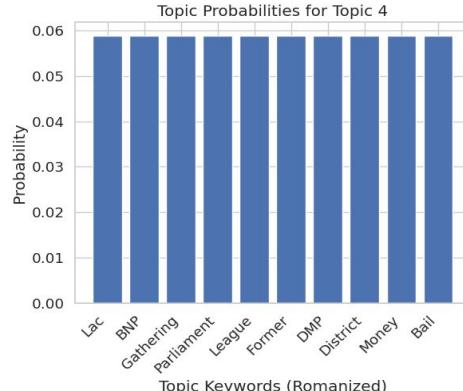
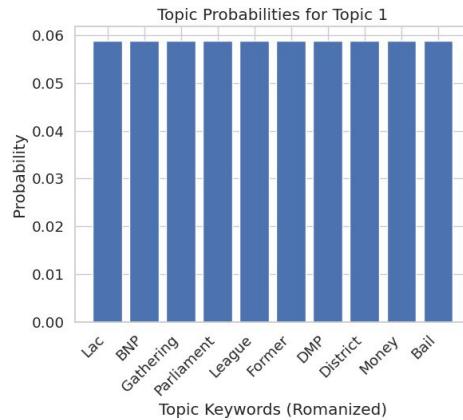
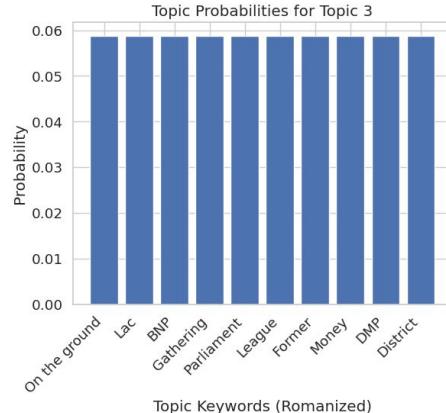
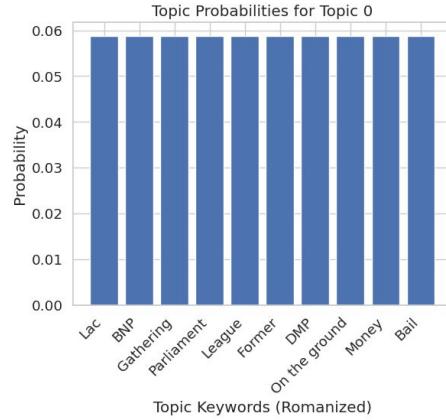
Most trending topic of the day is  
Topic: 0 with probabilities :

$0.096 * "হত্যা" + 0.096 * "প্রধানমন্ত্রী" +$   
 $0.080 * "সময়" + 0.080 * "মিশন" + 0.065 * "$   
 $মহাসমাবেশ" + 0.065 * "বিএনপি" + 0.065 * "$   
 $মামলায়" + 0.065 * "আইনজীবী" + 0.065 * "$   
 $ওল্ব" + 0.065 * "পররাষ্ট্র"$

# What was the trending topic on July 27, 2023?



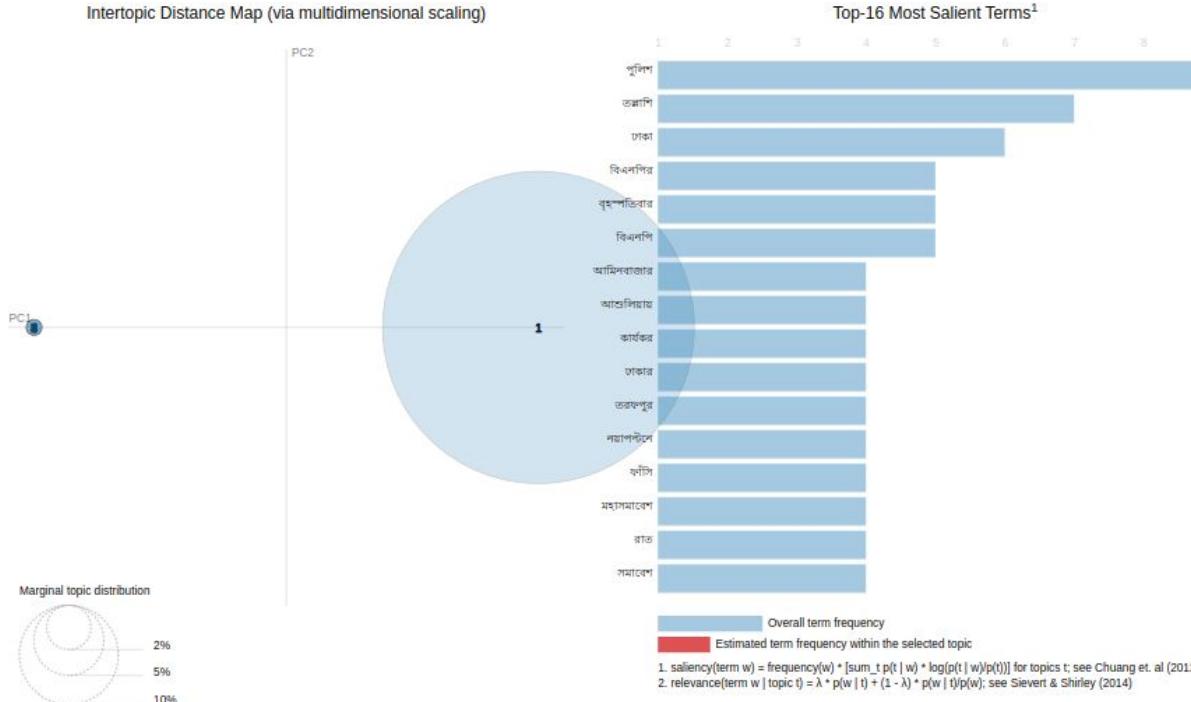
# What was the trending topic on July 27, 2023? (Continued)



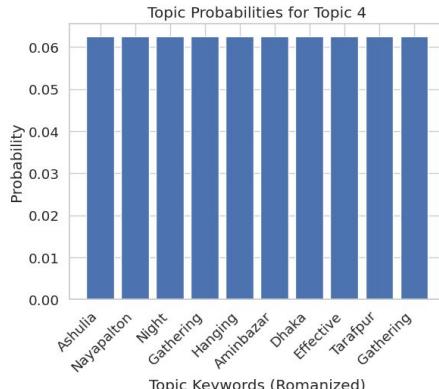
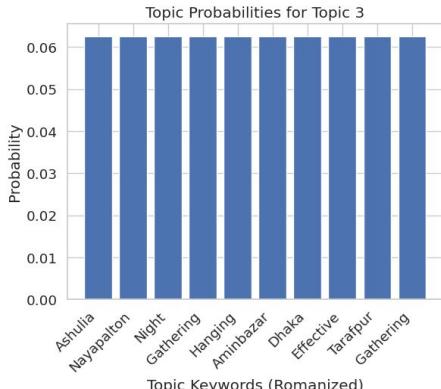
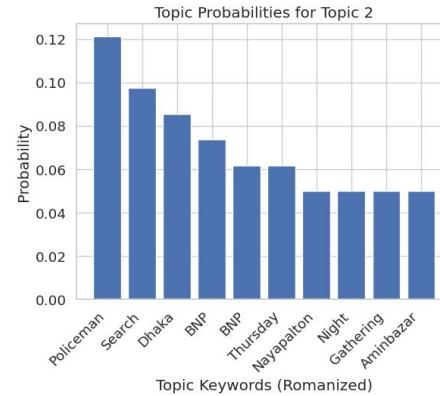
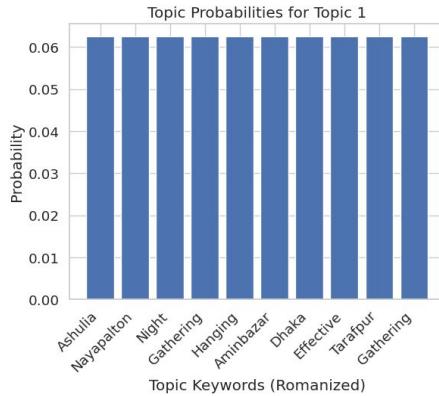
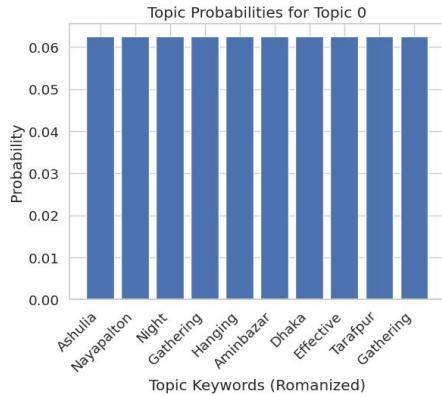
Most trending topic of the day is  
Topic: 2 with probabilities :

$$0.098 * \text{"আ"} + 0.086 * \text{"সমাবেশের"} + \\ 0.074 * \text{"লীগের"} + 0.062 * \text{"অনুমতি"} + \\ 0.062 * \text{"ওক্টোবাৰ"} + 0.062 * \text{"পুলিশ"} + \\ 0.050 * \text{"লীগ"} + 0.050 * \text{"সমাবেশ"} + \\ 0.050 * \text{"সংসদ"} + 0.050 * \text{"বিএনপিৰকে"}$$

# What was the trending topic on July 28, 2023?



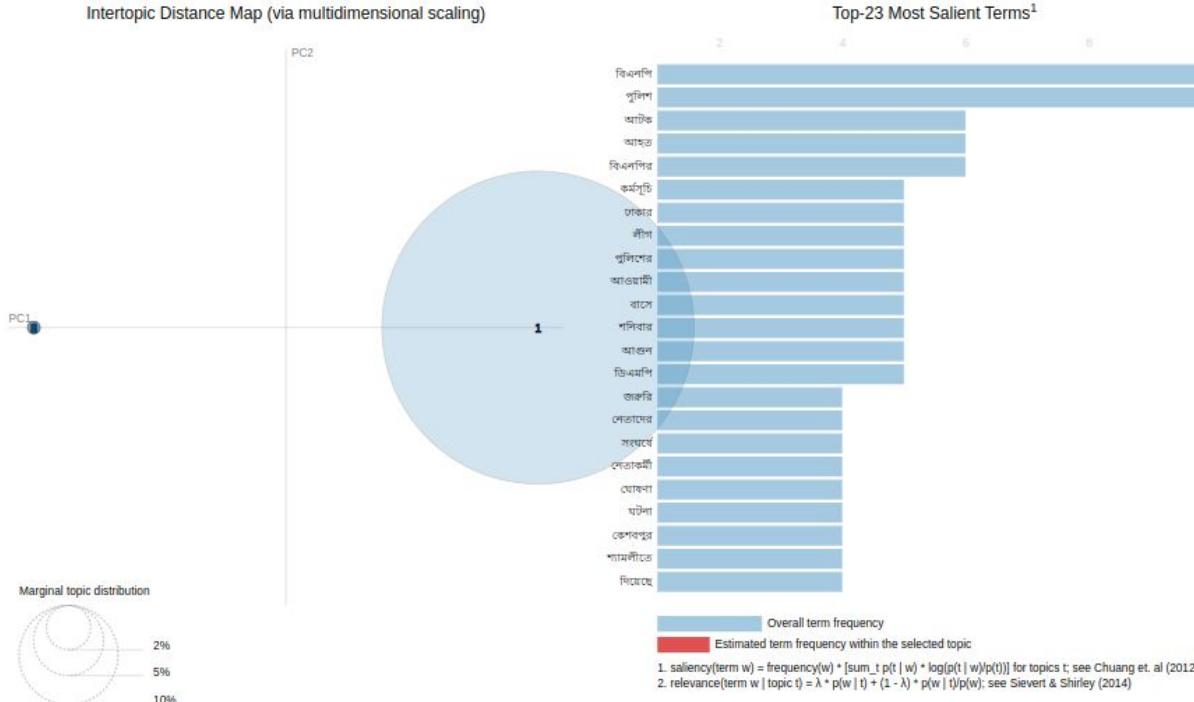
# What was the trending topic on July 28, 2023? (Continued)



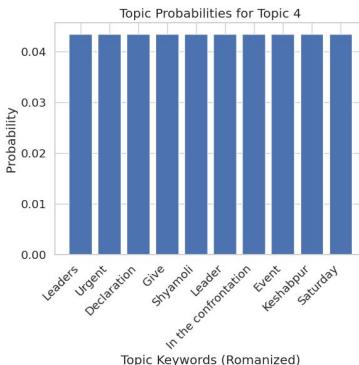
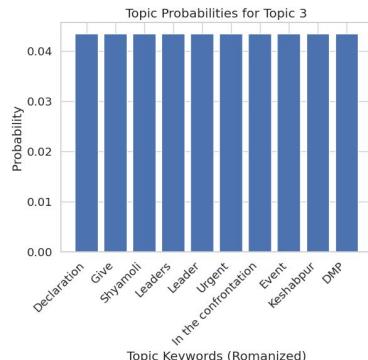
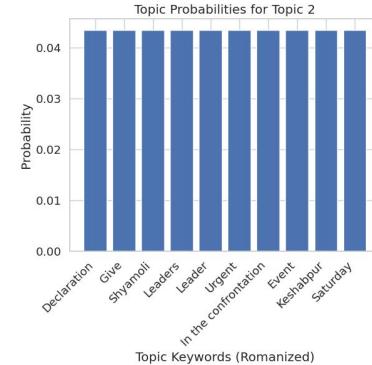
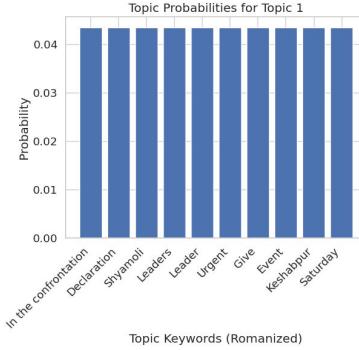
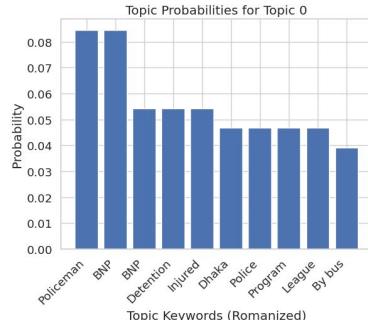
Most trending topic of the day is  
Topic: 2 with probabilities :

$$0.121 * \text{"পুলিশ"} + 0.097 * \text{"তল্লাশি"} + 0.086 * \text{"ঢাকা"} + 0.074 * \text{"বিএনপির"} + 0.062 * \text{"বিএনপি"} + 0.062 * \text{"বৃহস্পতিবার"} + 0.050 * \text{"ন্যাপল্টনে"} + 0.050 * \text{"রাত"} + 0.050 * \text{"মহাসমাবেশ"} + 0.050 * \text{"আমিনবাজার"}$$

# What was the trending topic on July 29, 2023?



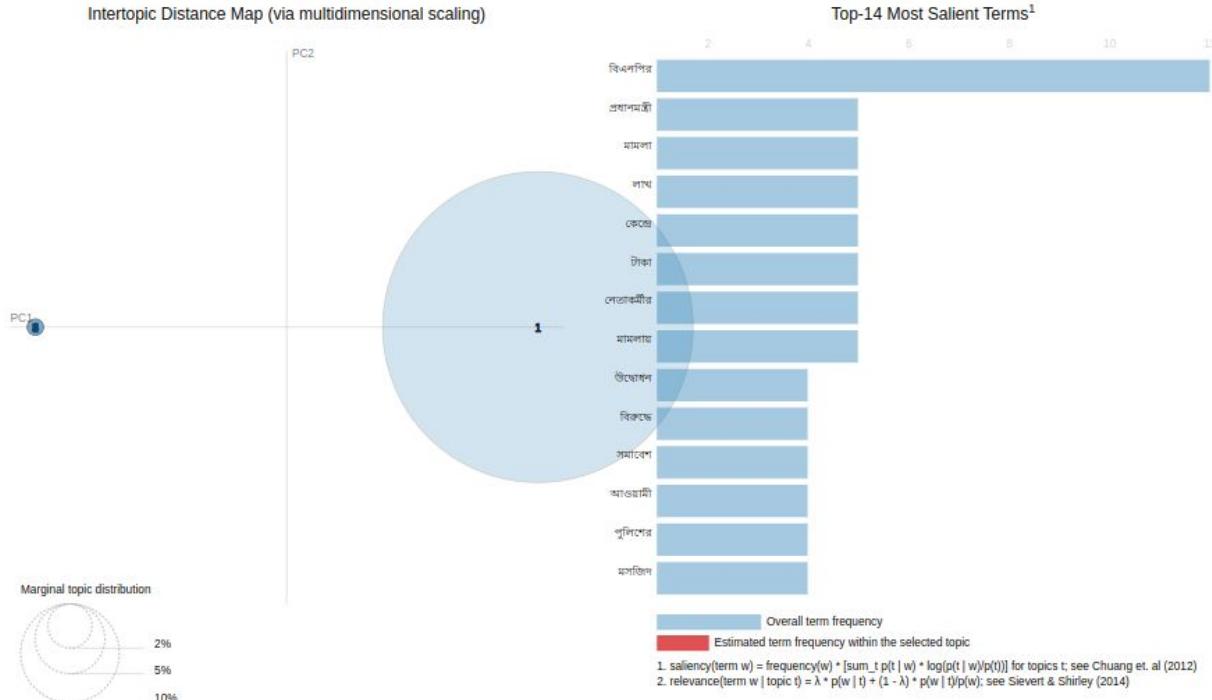
# What was the trending topic on July 29, 2023? (Continued)



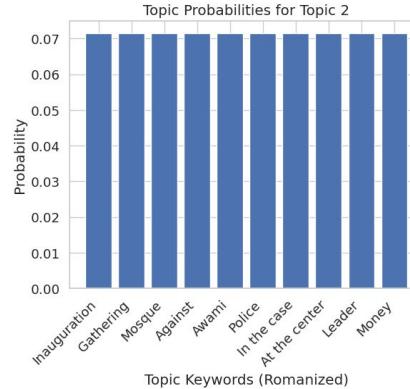
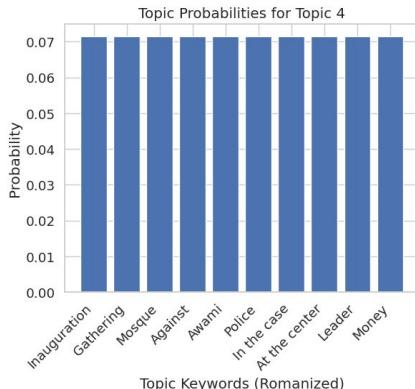
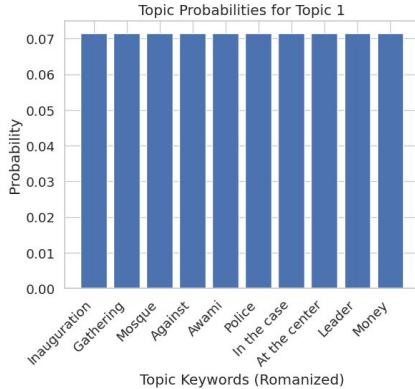
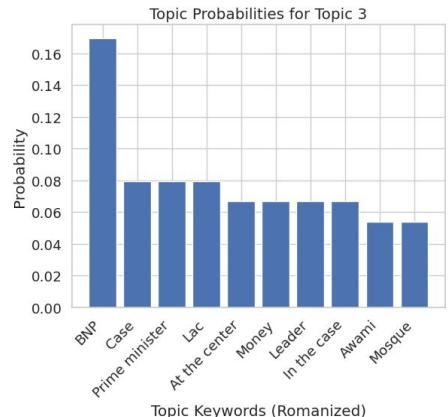
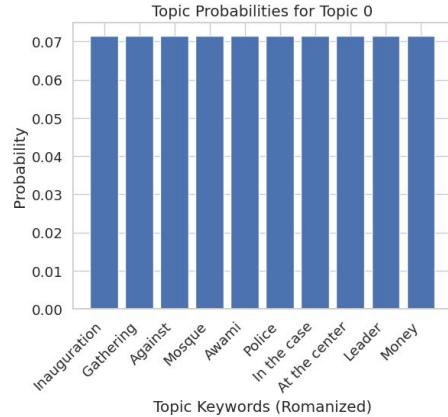
Most trending topic of the day is  
Topic: 0 with probabilities :

$$0.084 * \text{পুলিশ} + 0.084 * \text{বিএনপি} + 0.054 * \text{বিএনপির} + 0.054 * \text{আটক} + 0.054 * \text{আহত} + 0.047 * \text{ঢাকার} + 0.047 * \text{পুলিশের} + \boxed{\phantom{0.047}} + 0.047 * \text{কর্মসূচি} + 0.047 * \text{লীগ} + 0.039 * \text{বামে}$$

# What was the trending topic on July 30, 2023?



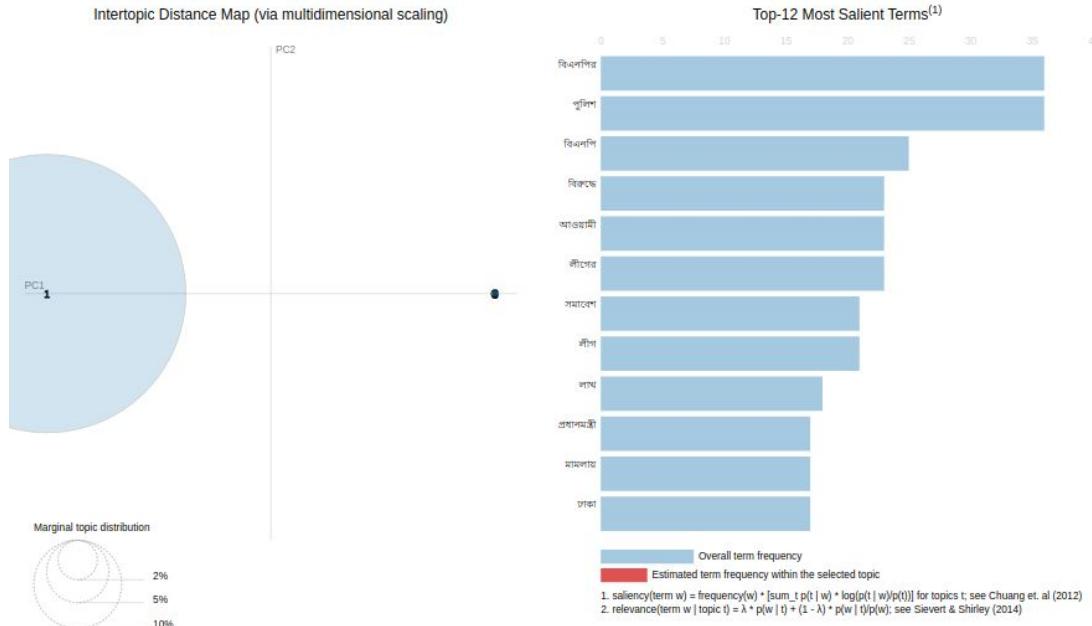
# What was the trending topic on July 30, 2023? (Continued)



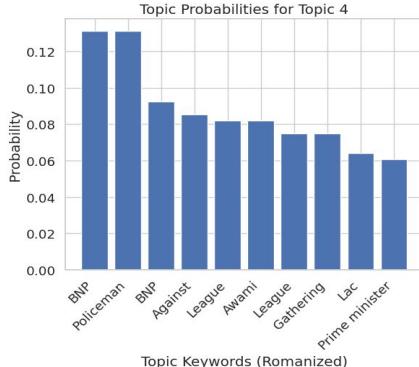
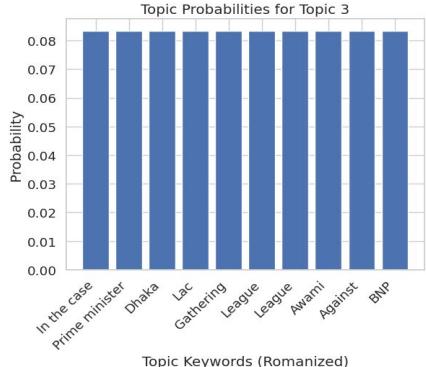
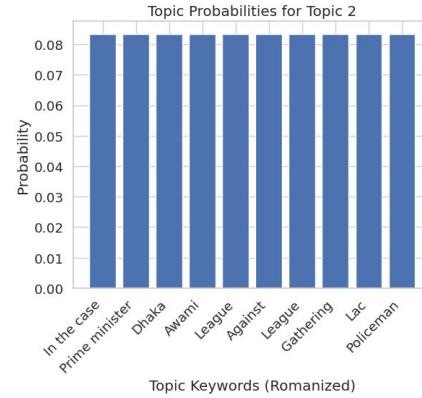
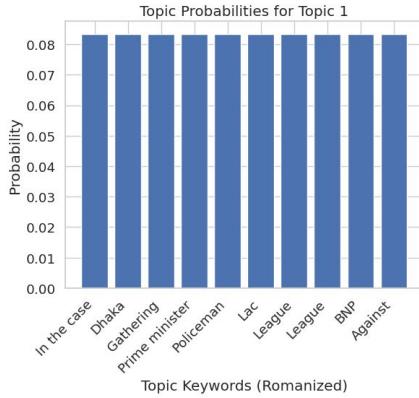
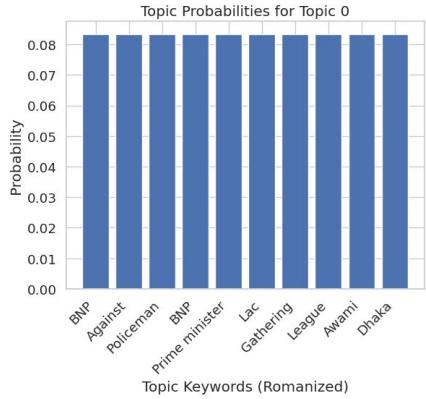
Most trending topic of the day is  
Topic: 3 with probabilities :

$$0.170 * \text{"বিএনপির"} + 0.080 * \text{"মামলা"} + \\ 0.080 * \text{"প্রধানমন্ত্রী"} + 0.080 * \text{"লাখ"} + 0.067 * \text{"কেন্দ্র"} + 0.067 * \text{"টাকা"} + 0.067 * \text{"নেতাকর্মীর"} + 0.067 * \text{"মামলায়"} + 0.054 * \text{"আওয়ামী"} + 0.054 * \text{"মসজিদ"}$$

# What is the trending topic of this week?



# What is the trending topic of this week? (continued)



Most trending topic of the day is  
Topic: 4 with probabilities :

0.131 \* "বিএনপির" + 0.131 \* "পুলিশ" +  
0.092 \* "বিএনপি" + 0.085 \* "বিক্রকে" +  
0.082 \* "লীগের" + 0.082 \* "আওয়ামী" +  
0.075 \* "লীগ" + 0.075 \* "সমাবেশ" + 0.064 \* "  
লাখ" + 0.061 \* "প্রধানমন্ত্রী"

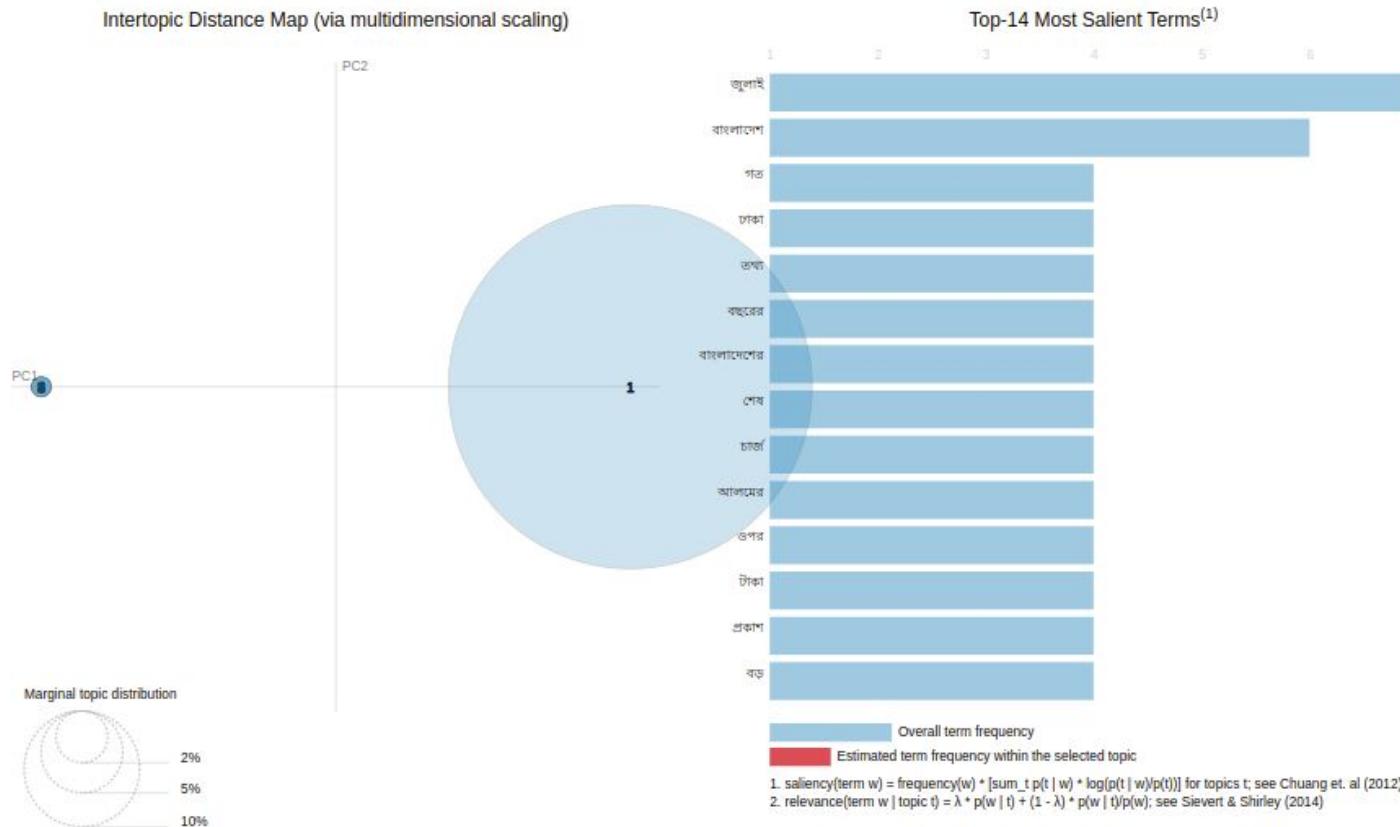
# **Peak User Hour News**

Kazi Shadman Sakib  
4th Year Undergraduate Student

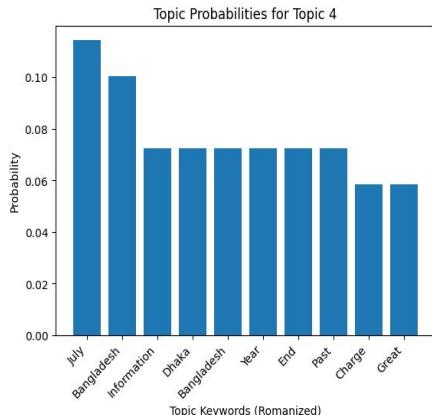
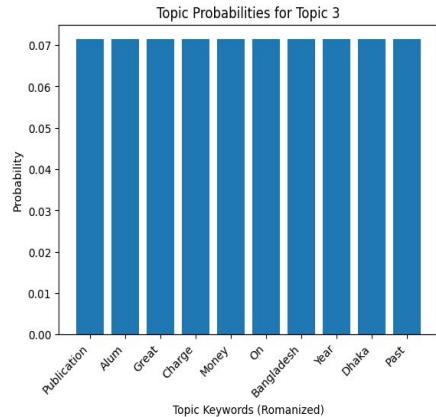
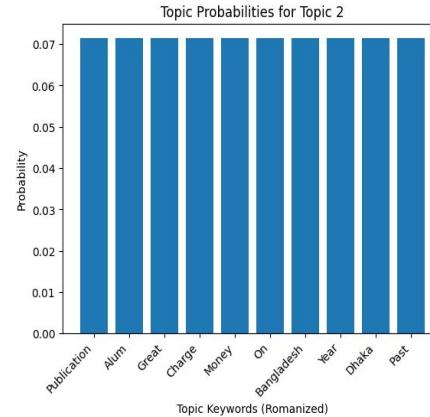
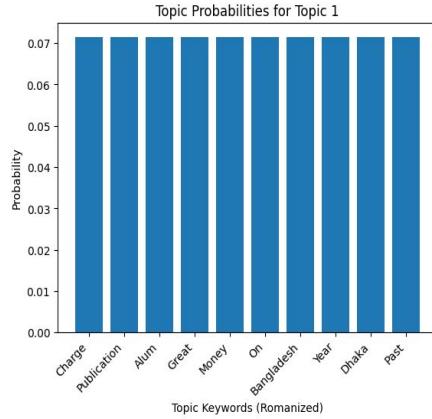
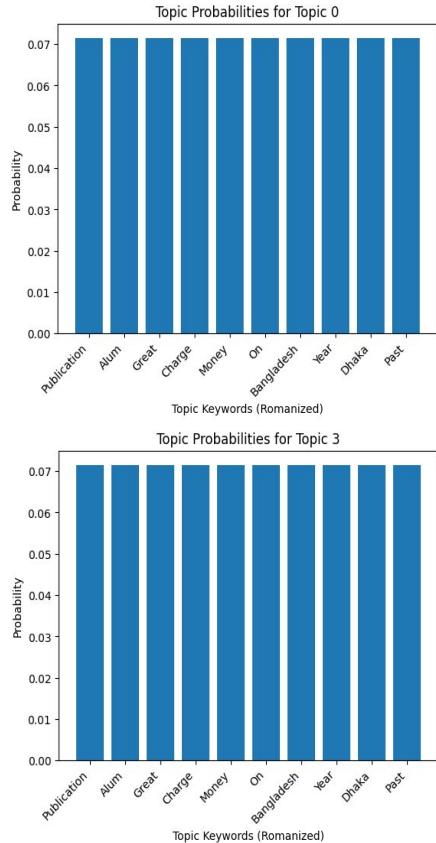
Department of Computer Science & Engineering  
University of Dhaka



# What was the most popular topic during the peak user hour on July 24, 2023?



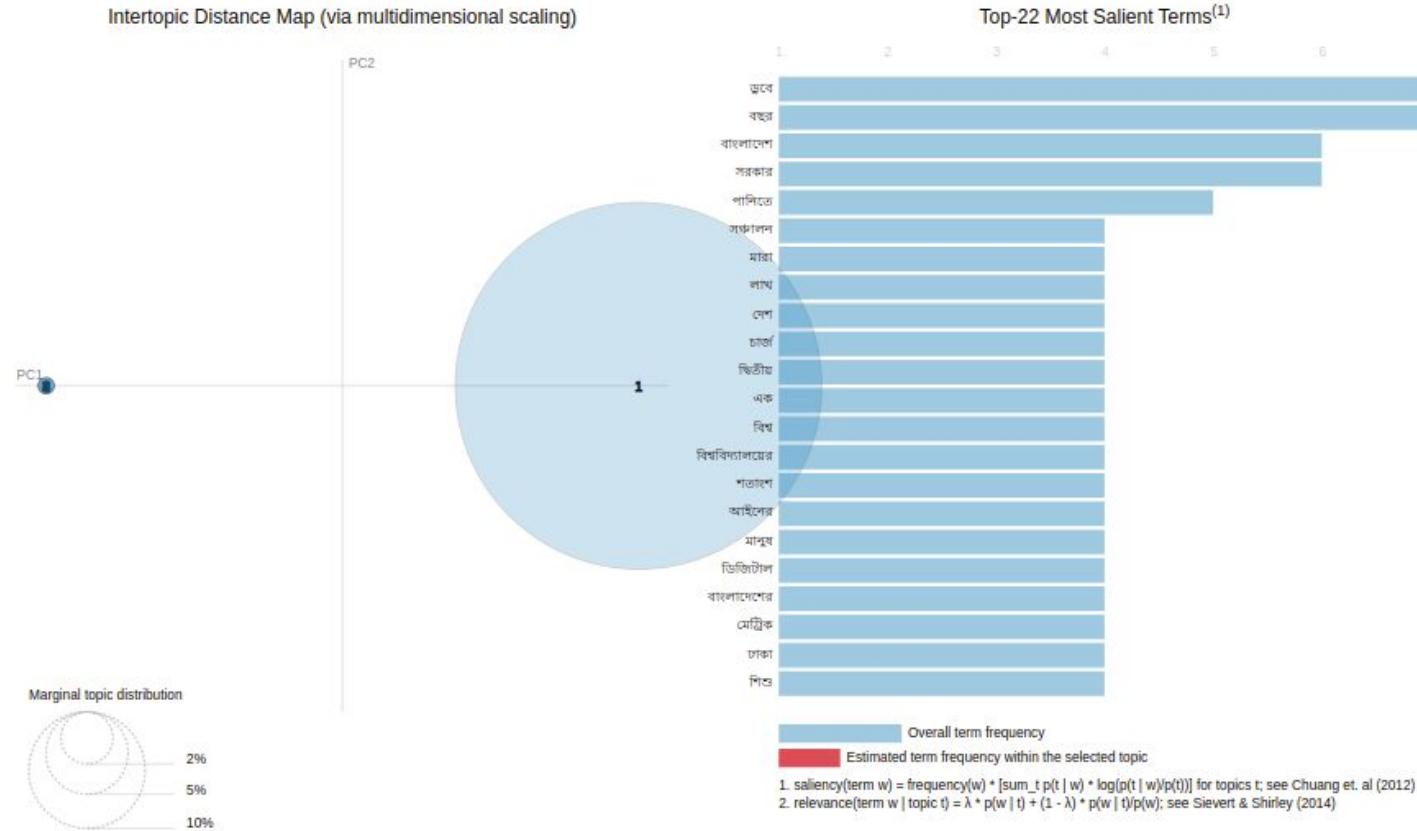
# What was the most popular topic during the peak user hour on July 24, 2023? (Continued)



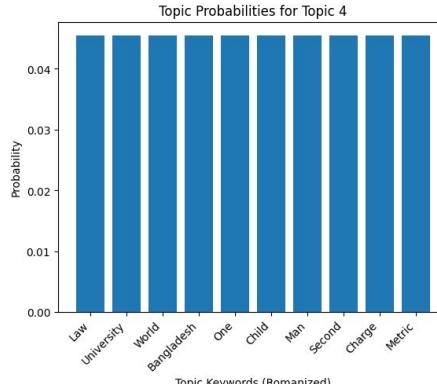
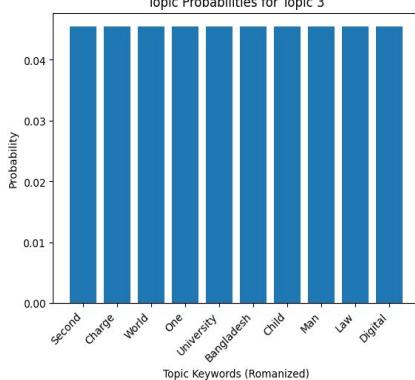
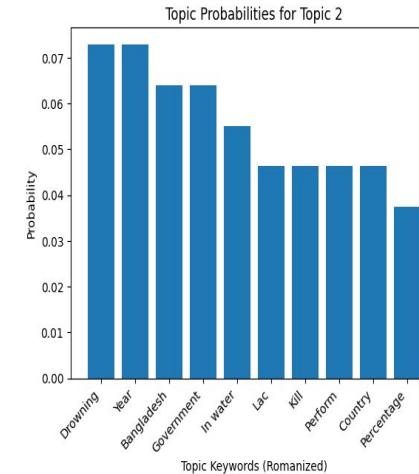
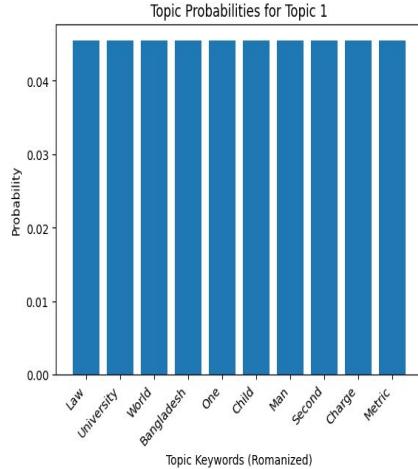
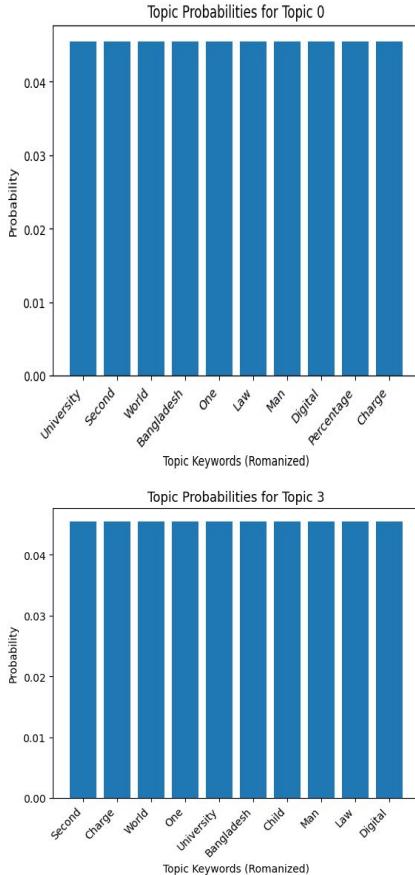
Most trending topic of the day is  
Topic: 4 with probabilities :

$0.114 * \text{"জুনাই"} + 0.100 * \text{"বাংলাদেশ"} +$   
 $0.072 * \text{"তথ্য"} + 0.072 * \text{"ঢাকা"} + 0.072 * \text{"বাংলাদেশের"} + 0.072 * \text{"বছরের"} + 0.072 * \text{"শেষ"} + 0.072 * \text{"গত"} + 0.058 * \text{"চার্জ"} +$   
 $0.058 * \text{"বড়"}$

# What was the most popular topic during the peak user hour on July 25, 2023?



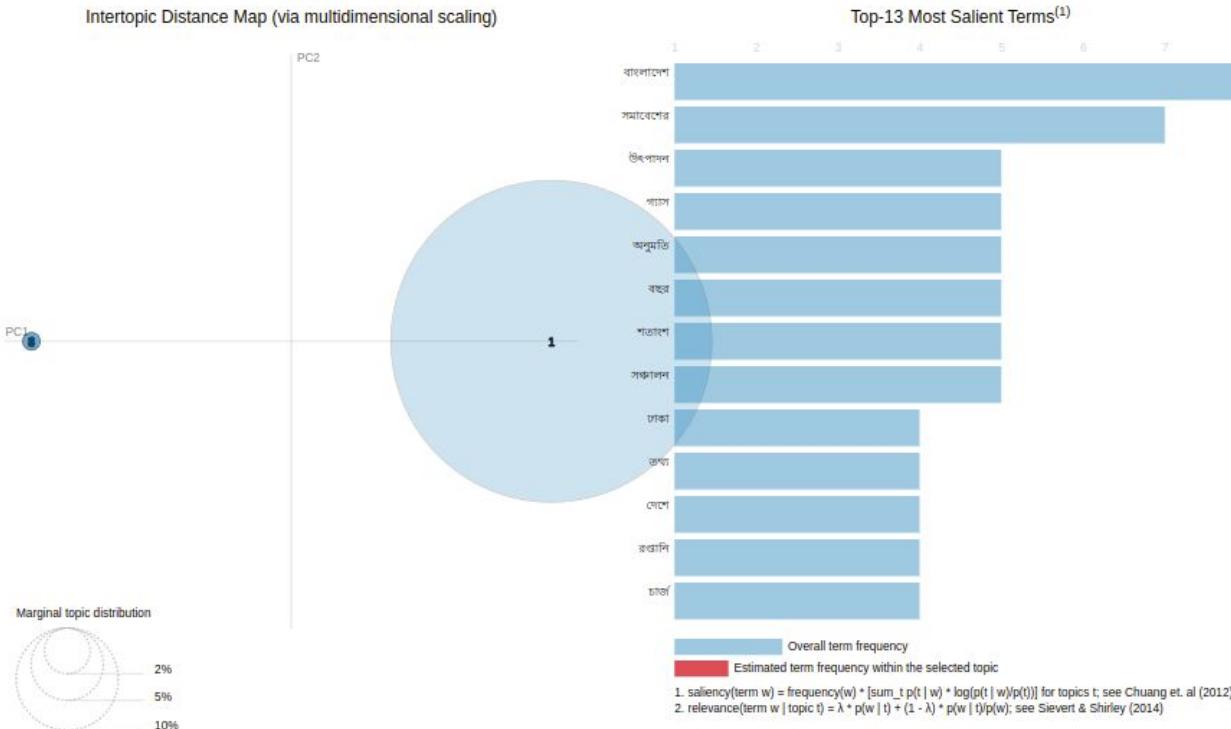
# What was the most popular topic during the peak user hour on July 25, 2023? (Continued)



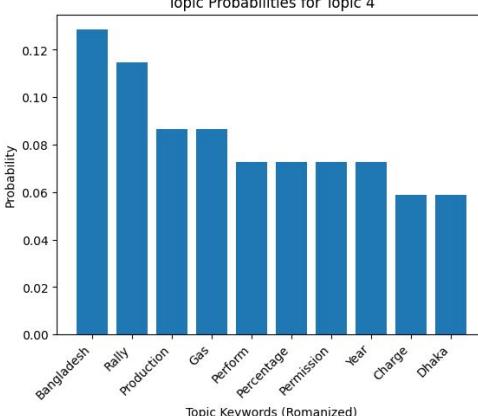
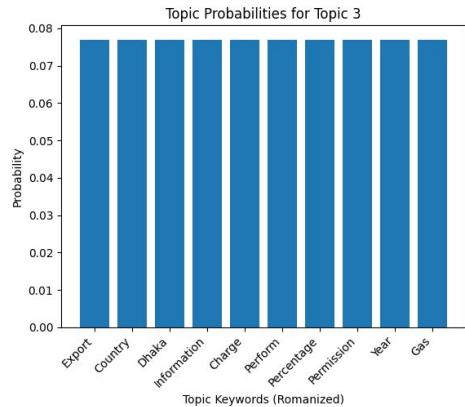
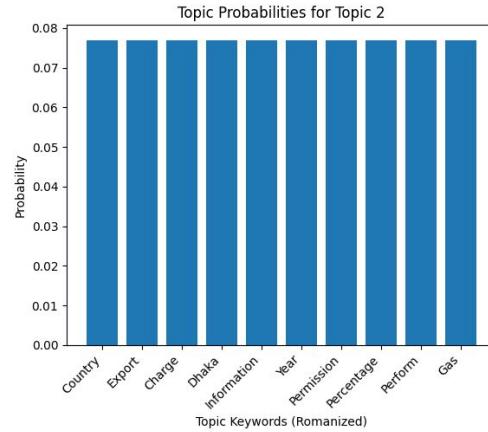
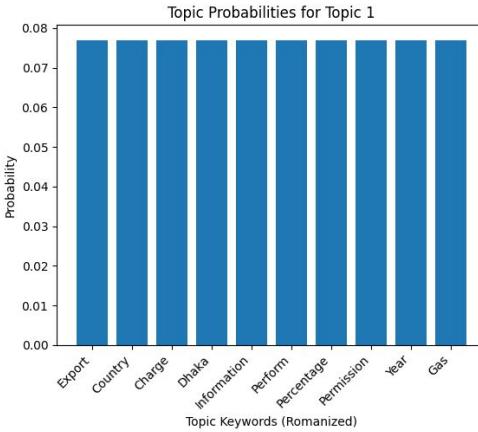
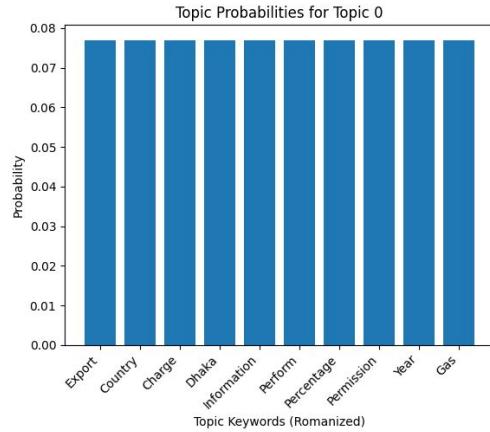
Most trending topic of the day is  
Topic: 2 with probabilities :

$0.073 * \text{"ডুবে"} + 0.073 * \text{"বছর"} + 0.064 * \text{"বাংলাদেশ"} + 0.064 * \text{"সরকার"} + 0.055 * \text{"পানিতে"} + 0.046 * \text{"লাখ"} + 0.046 * \text{"মারা"} + 0.046 * \text{"সঞ্চালন"} + 0.046 * \text{"দেশ"} + 0.037 * \text{"শতাংশ"}$

# What was the most popular topic during the peak user hour on July 26, 2023?



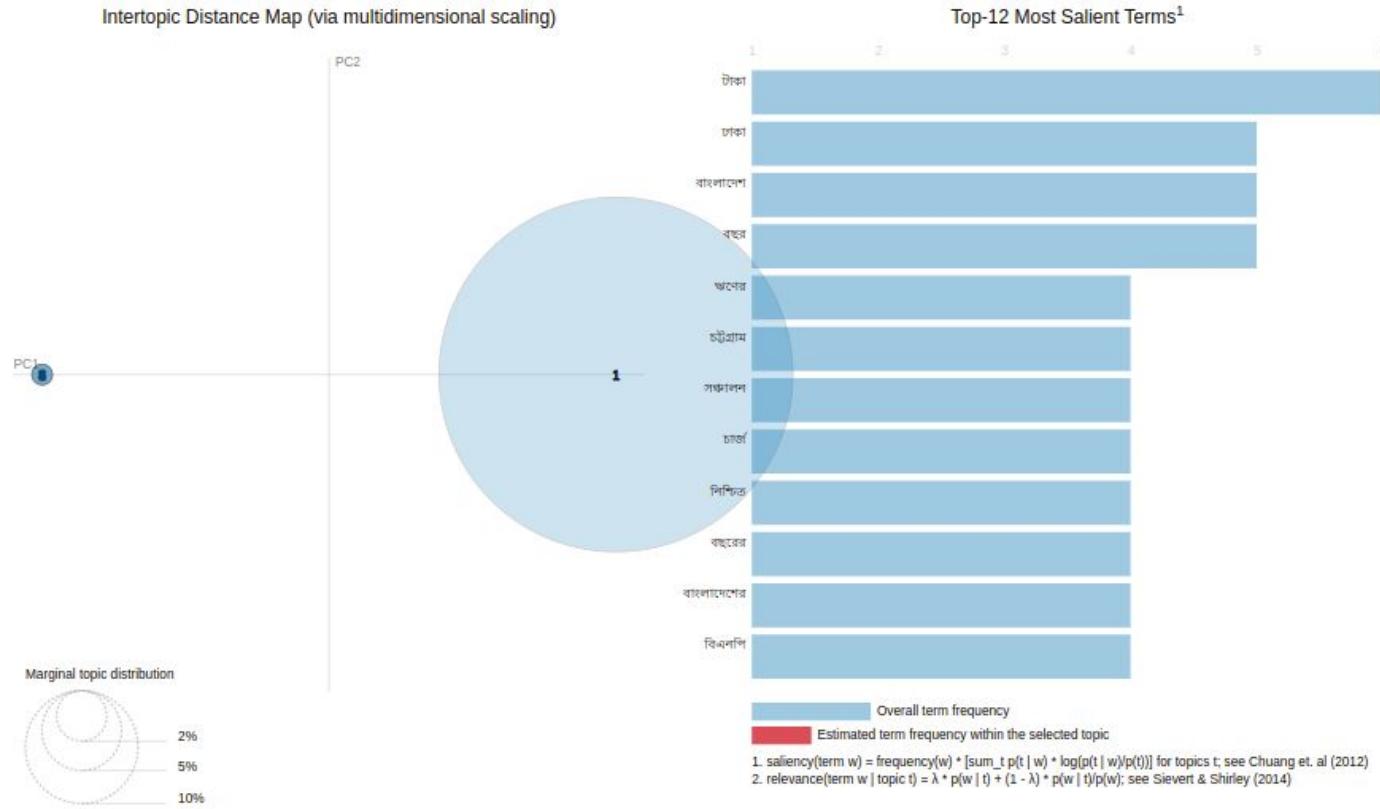
# What was the most popular topic during the peak user hour on July 26, 2023? (Continued)



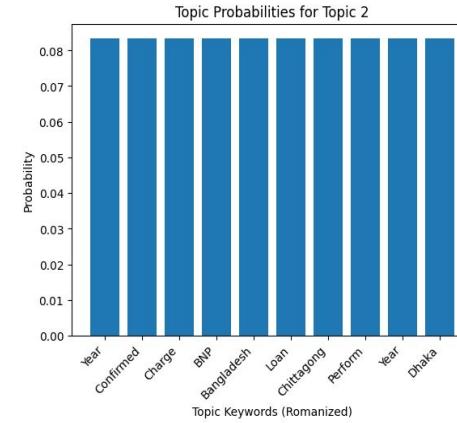
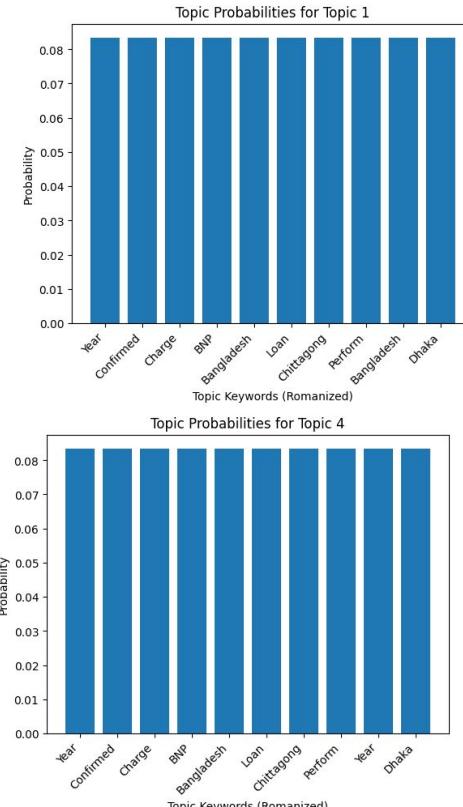
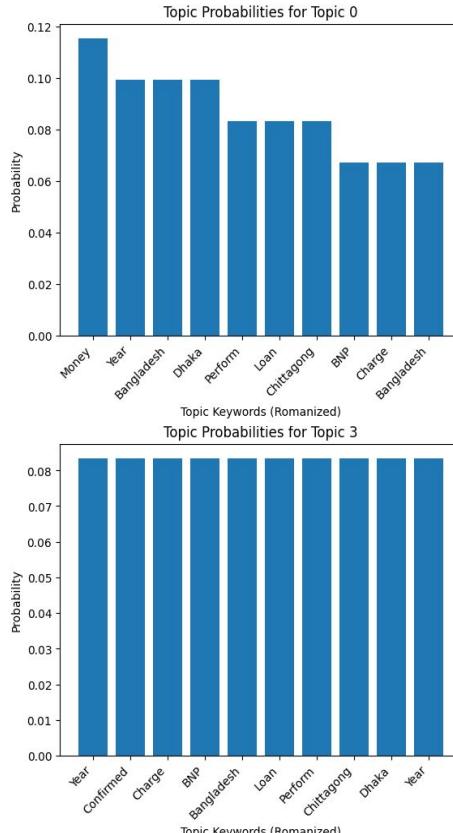
Most trending topic of the day is  
Topic: 4 with probabilities :

0.128 \* "বাংলাদেশ" + 0.115 \* "সমাবেশের" +  
0.087 \* "উৎপাদন" + 0.087 \* "গ্যাস" +  
0.073 \* "সঞ্চালন" + 0.073 \* "শতাংশ" +  
0.073 \* "অনুমতি" + 0.073 \* "বছর" +  
0.059 \* "চার্ট" + 0.059 \* "ঢাকা"

# What was the most popular topic during the peak user hour on July 27, 2023?



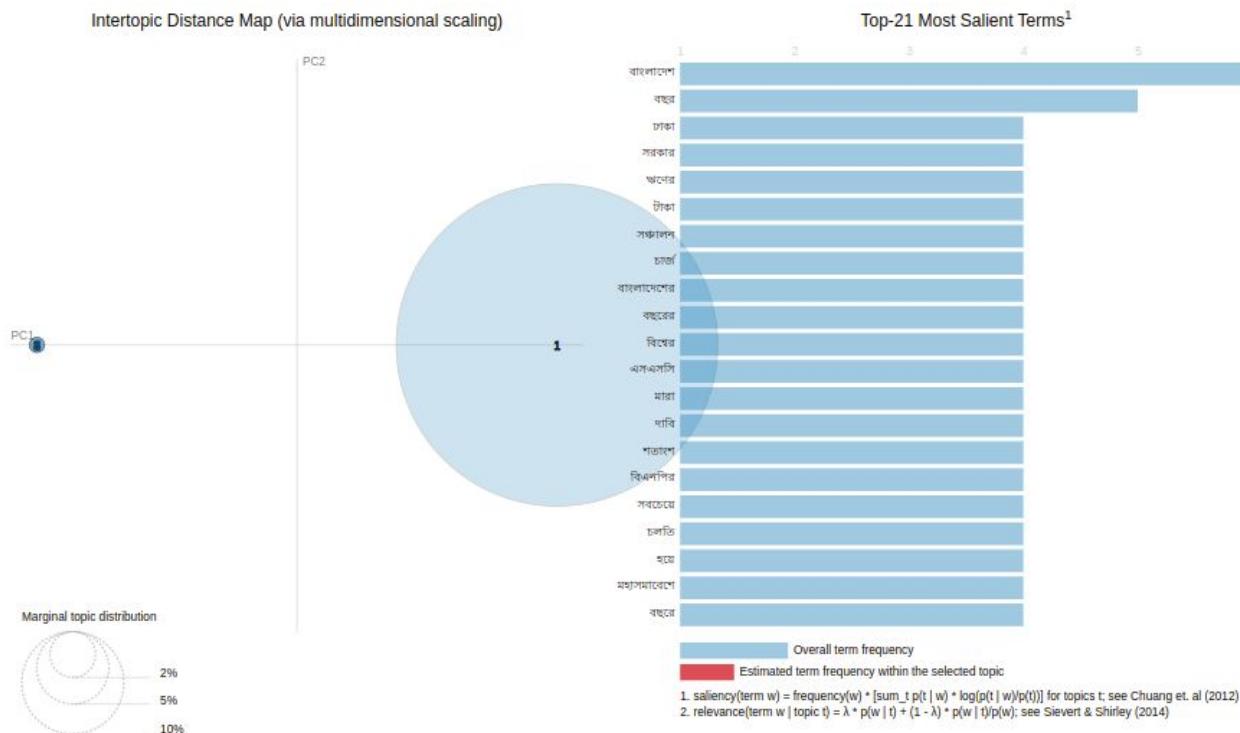
# What was the most popular topic during the peak user hour on July 27, 2023? (Continued)



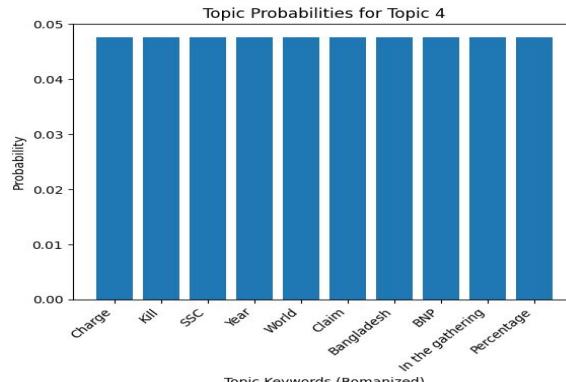
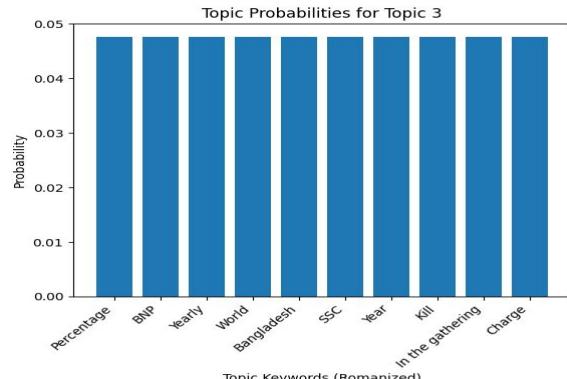
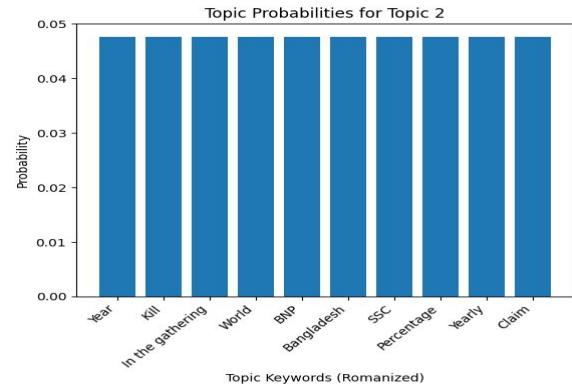
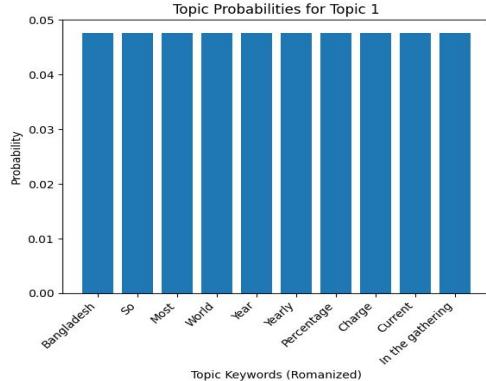
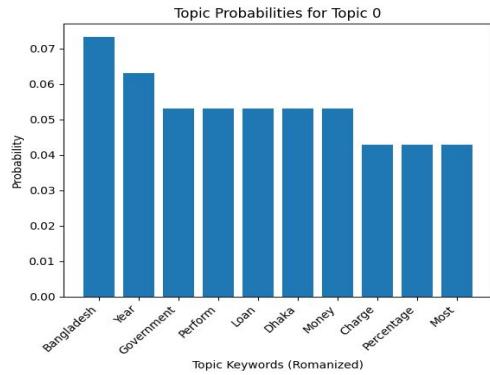
Most trending topic of the day is  
Topic: 0 with probabilities :

0.115 \* "টাকা" + 0.099 \* "বছর" + 0.099 \* "বাংলাদেশ" + 0.099 \* "ঢাকা" + 0.083 \* "সঞ্চালন" + 0.083 \* "ঝরেন" + 0.083 \* "চট্টগ্রাম" + 0.067 \* "বিএনপি" + 0.067 \* "চার্জ" + 0.067 \* "বাংলাদেশের"

# What was the most popular topic during the peak user hour on July 28, 2023?



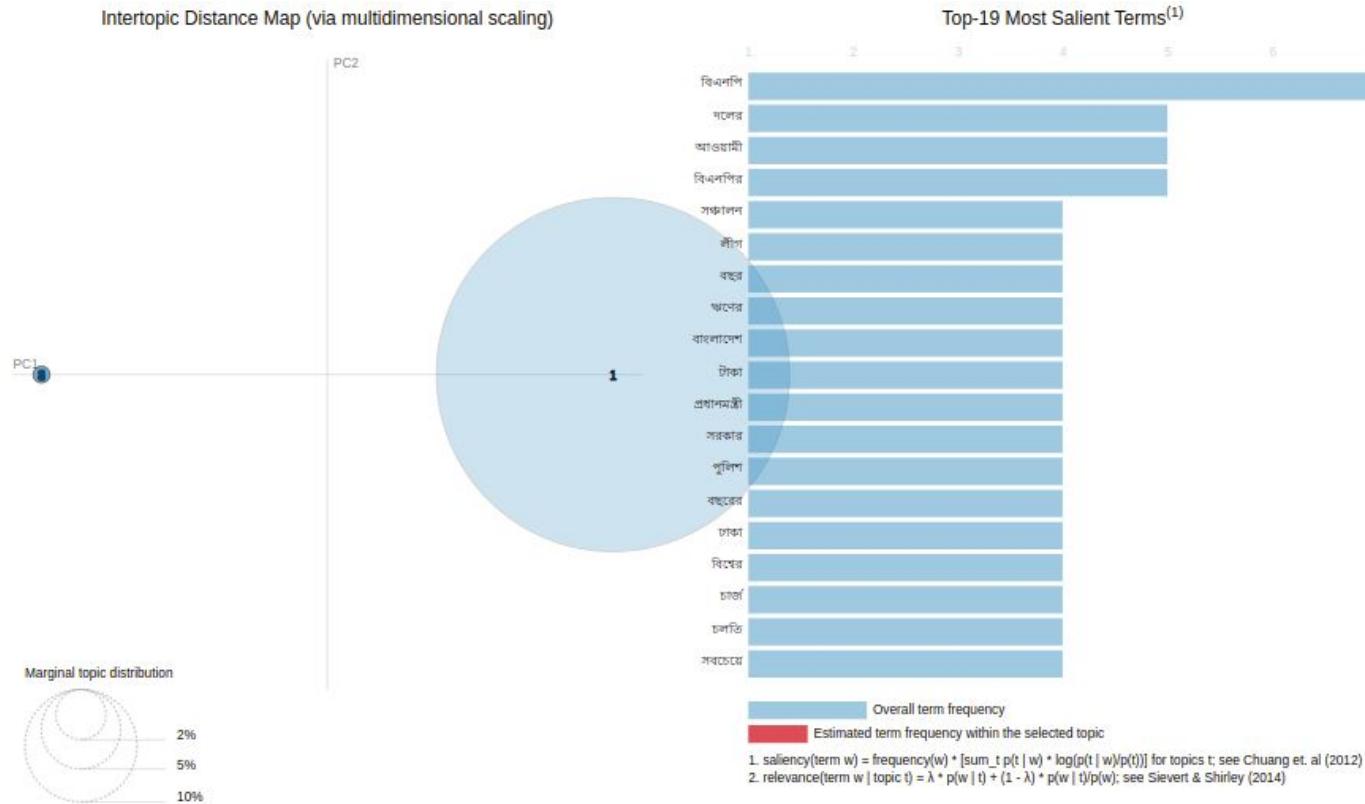
# What was the most popular topic during the peak user hour on July 28, 2023? (Continued)



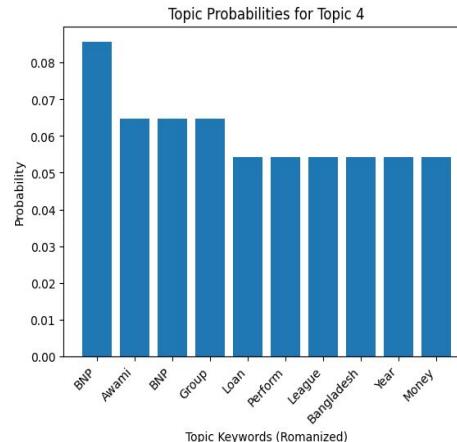
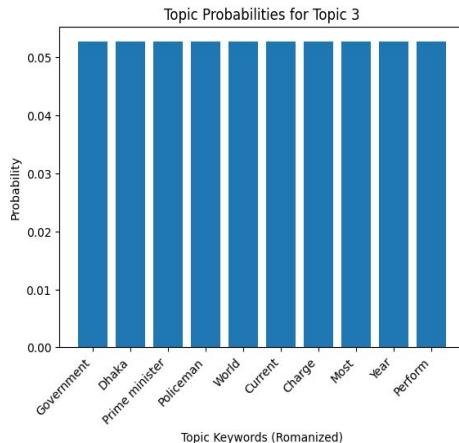
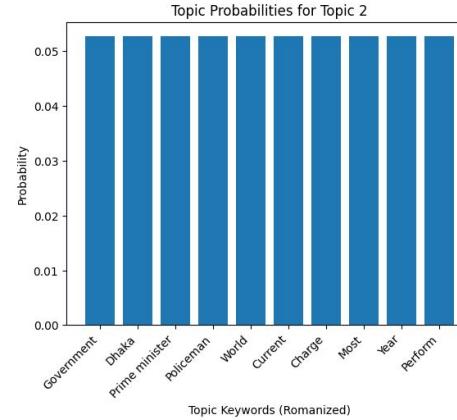
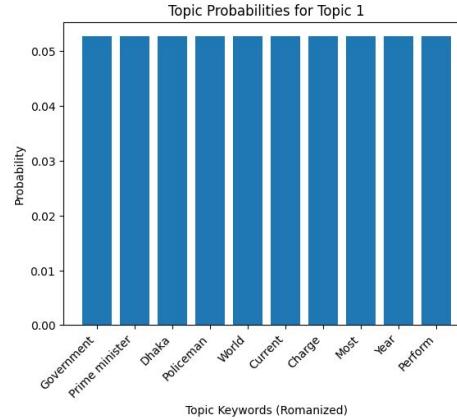
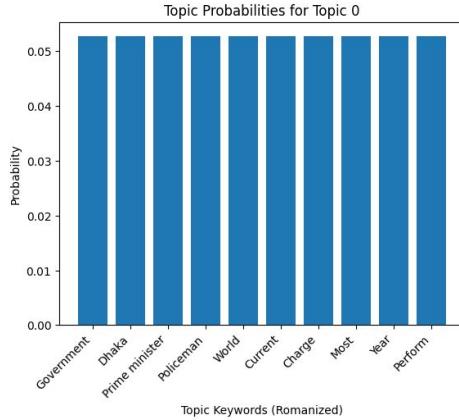
Most trending topic of the day is  
Topic: 0 with probabilities :

$$0.073 \times \text{"বাংলাদেশ"} + 0.063 \times \text{"বছর"} + \\ 0.053 \times \text{"সরকার"} + 0.053 \times \text{"সঞ্চালন"} + \\ 0.053 \times \text{"খণ্ডের"} + 0.053 \times \text{"ঢাকা"} + \\ 0.053 \times \text{"টাকা"} + 0.043 \times \text{"চাজ"} + \\ 0.043 \times \text{"শতাংশ"} + 0.043 \times \text{"সরচেয়ের"}$$

# What was the most popular topic during the peak user hour on July 29, 2023?



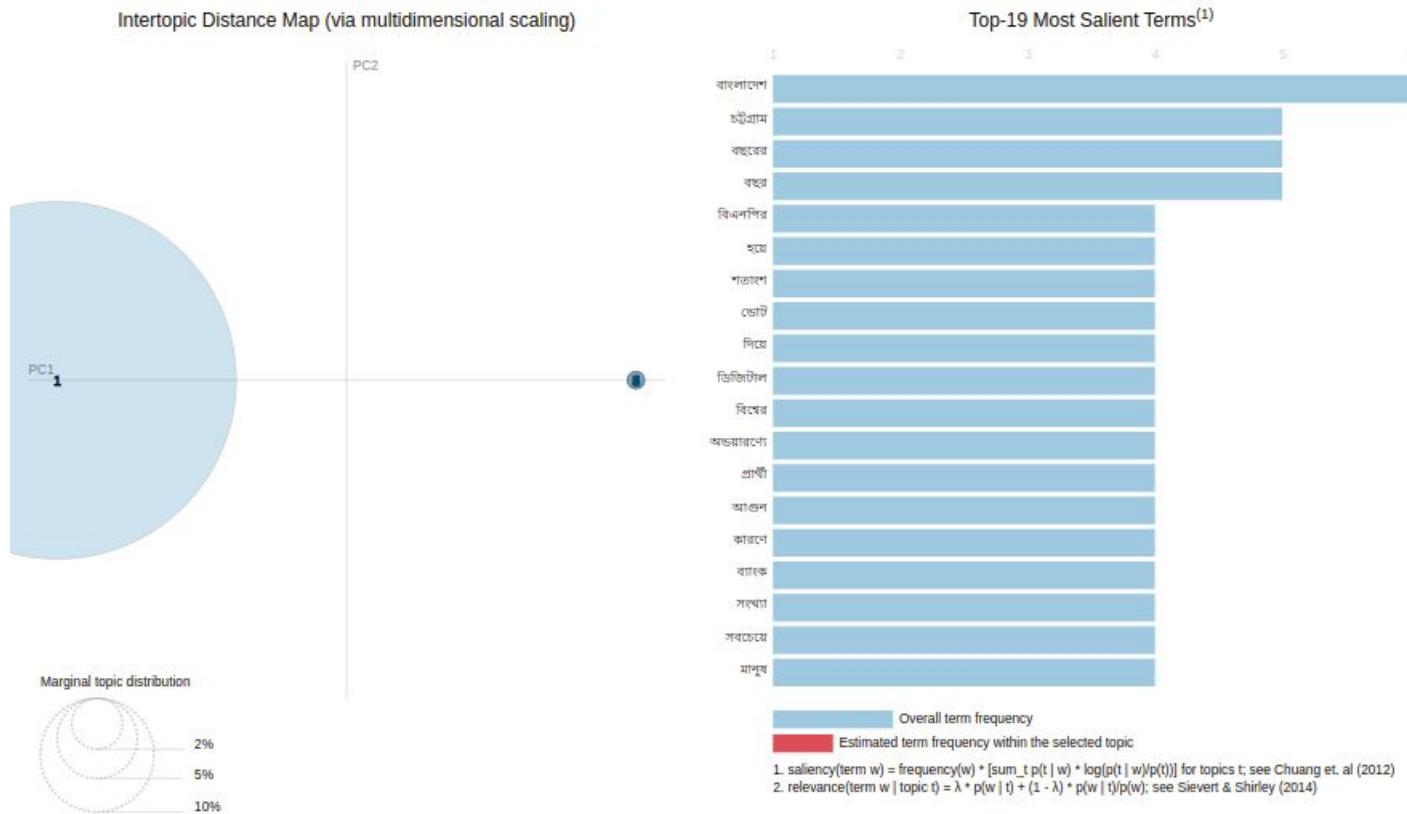
# What was the most popular topic during the peak user hour on July 29, 2023? (Continued)



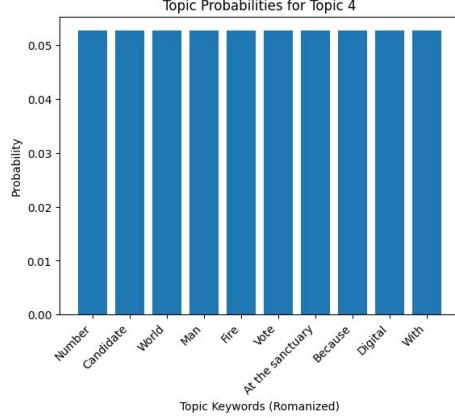
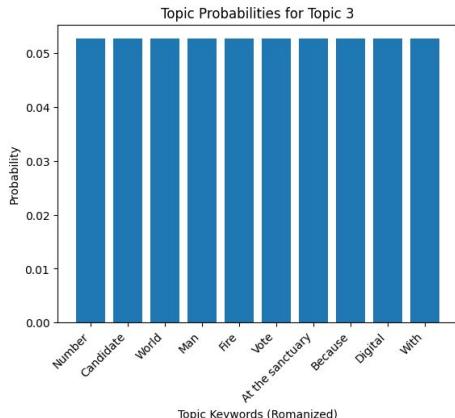
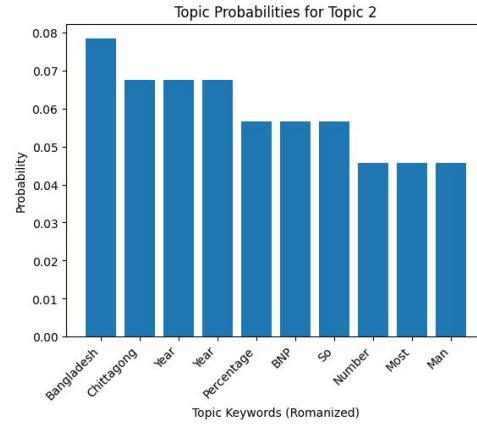
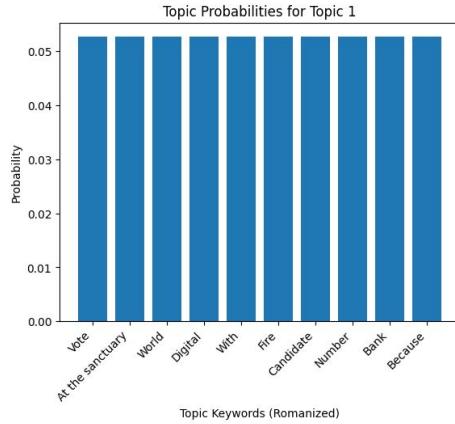
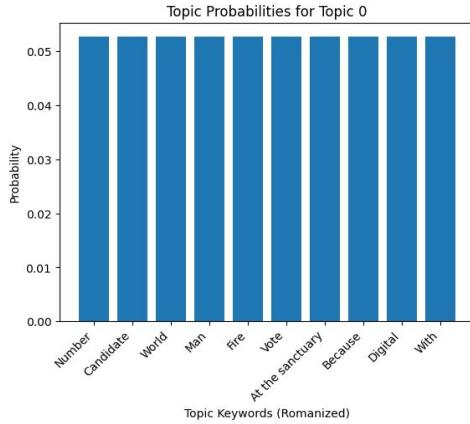
Most trending topic of the day is  
Topic: 4 with probabilities :

$$0.086 * \text{"বিএনপি"} + 0.065 * \text{"আওয়ামী"} + \\ 0.065 * \text{"বিএনপির"} + 0.065 * \text{"দলের"} + \\ 0.054 * \text{"ঞ্চণের"} + 0.054 * \text{"সঞ্চালন"} + \\ 0.054 * \text{"গীগ"} + 0.054 * \text{"বাংলাদেশ"} + \\ 0.054 * \text{"বছর"} + 0.054 * \text{"টাকা"}$$

# What was the most popular topic during the peak user hour on July 30, 2023?



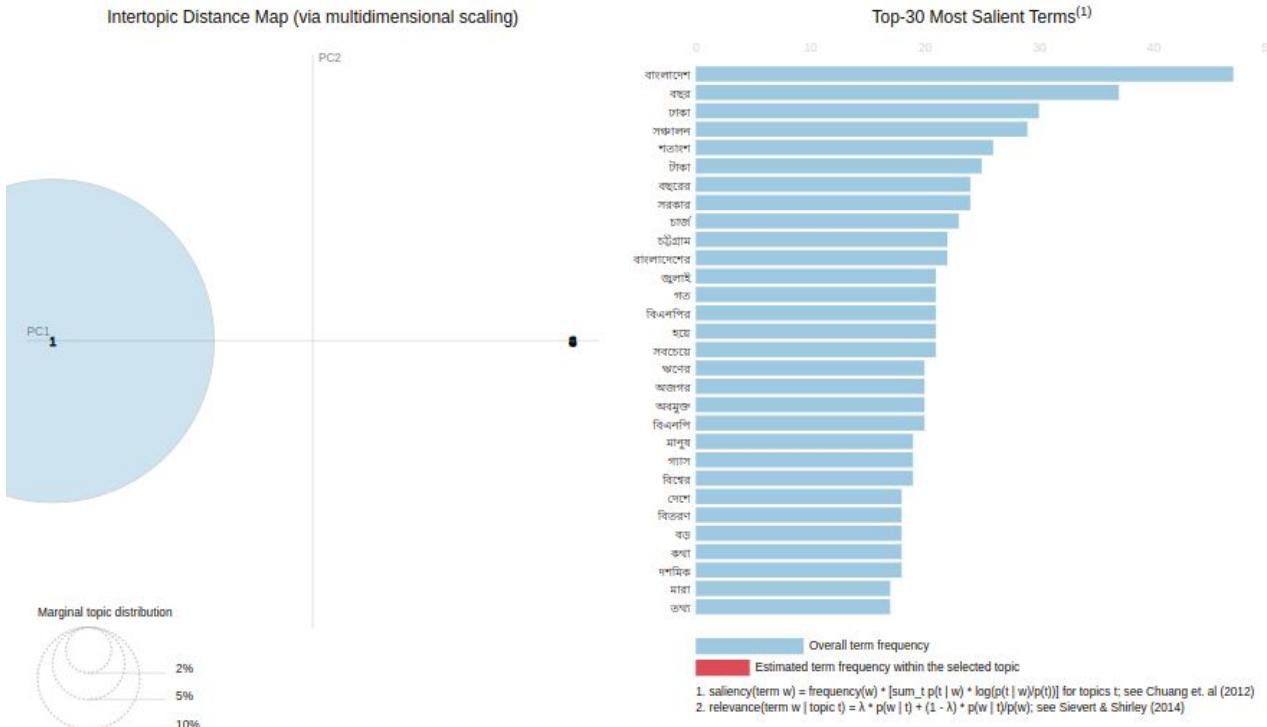
# What was the most popular topic during the peak user hour on July 30, 2023? (Continued)



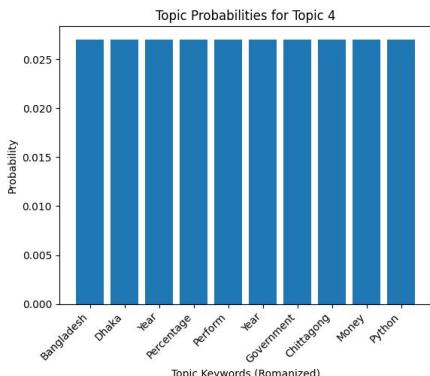
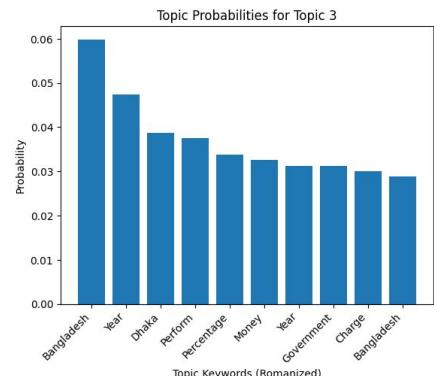
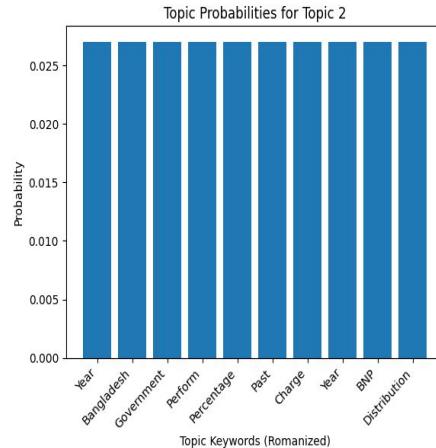
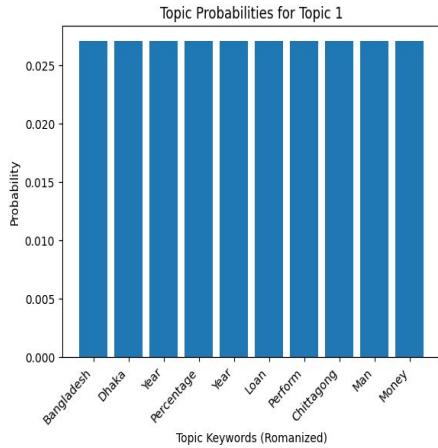
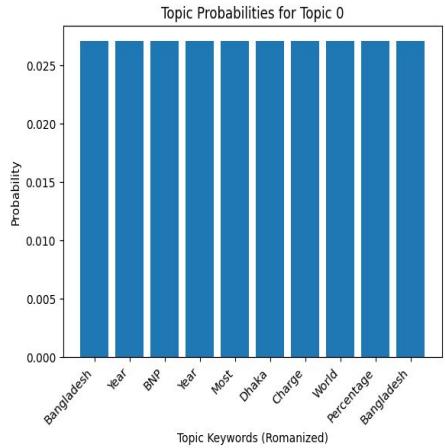
Most trending topic of the day is  
Topic: 2 with probabilities :

$$0.078 * \text{"বাংলাদেশ"} + 0.068 * \text{"চট্টগ্রাম"} + \\ 0.068 * \text{"বছর"} + 0.068 * \text{"বছরের"} + \boxed{0.057} \\ 0.057 * \text{"শতাংশ"} + 0.057 * \text{"বিএনপির"} + \boxed{0.057} \\ 0.057 * \text{"হয়ে"} + 0.046 * \text{"সংখ্যা"} + 0.046 * \text{"সবচেয়ে"} + 0.046 * \text{"মানুষ"} \boxed{0.046}$$

# What is the most popular topic during the peak user hour of this week?



# What is the most popular topic during the peak user hour of this week? (continued)



Most trending topic of the day is  
Topic: 3 with probabilities :

```
0.060 * "বাংলাদেশ" + 0.047 * "বছর" +  
0.039 * "টাকা" + 0.037 * "সঞ্চালন" +  
0.034 * "শতাংশ" + 0.033 * "টাকা" +  
0.031 * "বছরের" + 0.031 * "সরকার" +  
0.030 * "চাজ" + 0.029 * "বাংলাদেশের"
```

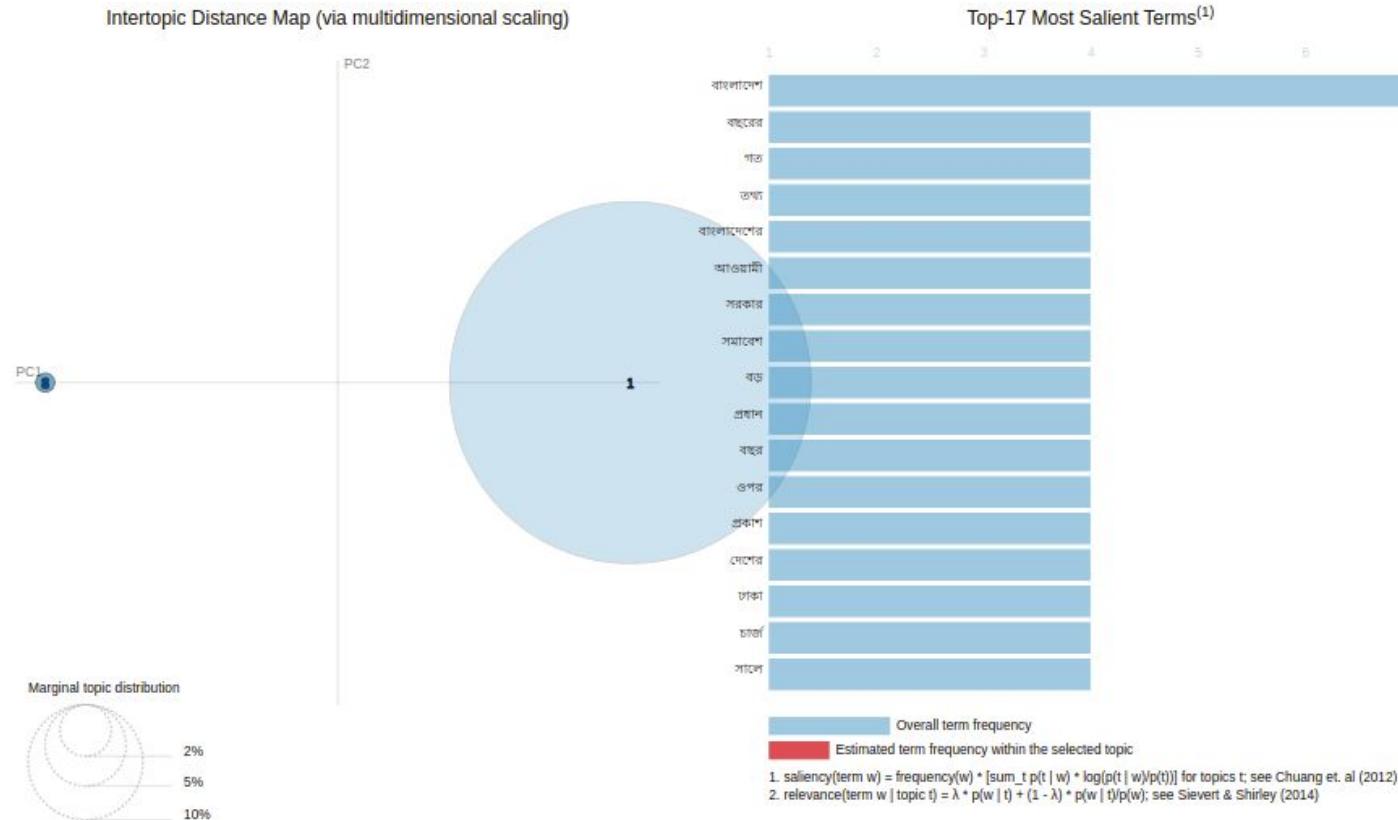
# Off Peak User Hour News

Kazi Shadman Sakib  
4th Year Undergraduate Student

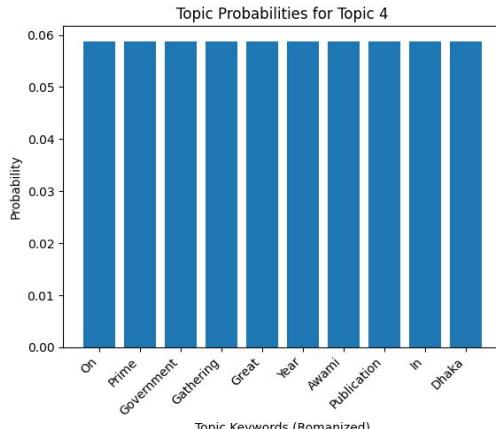
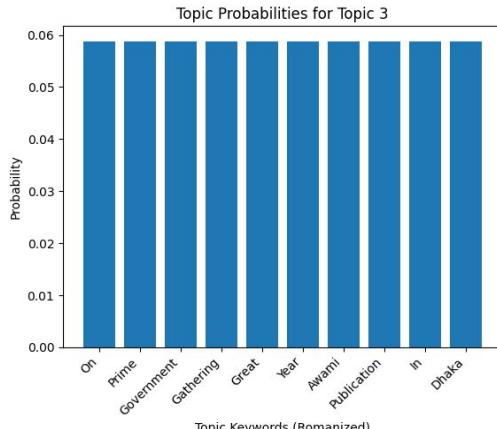
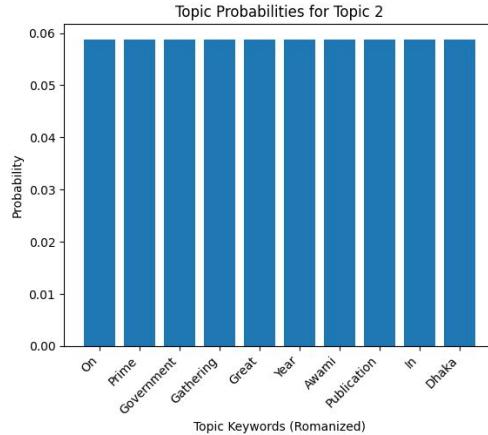
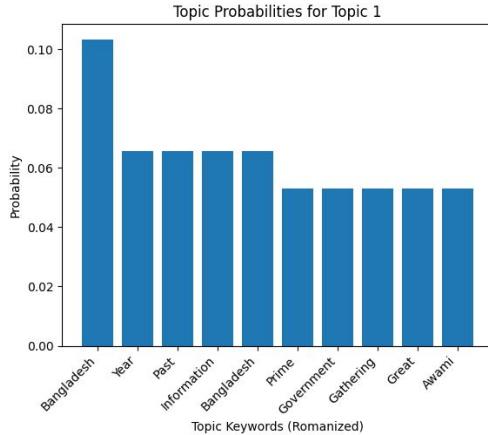
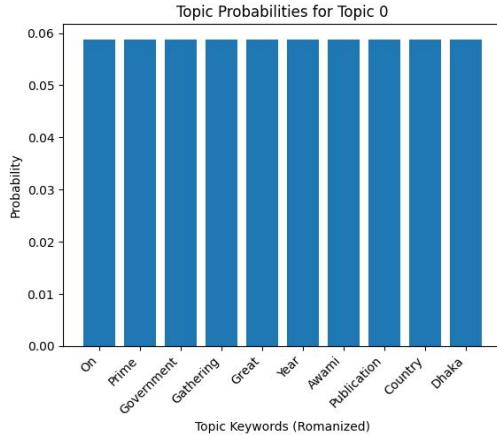
Department of Computer Science & Engineering  
University of Dhaka



# What was the most popular topic during the off peak user hour on July 24, 2023?



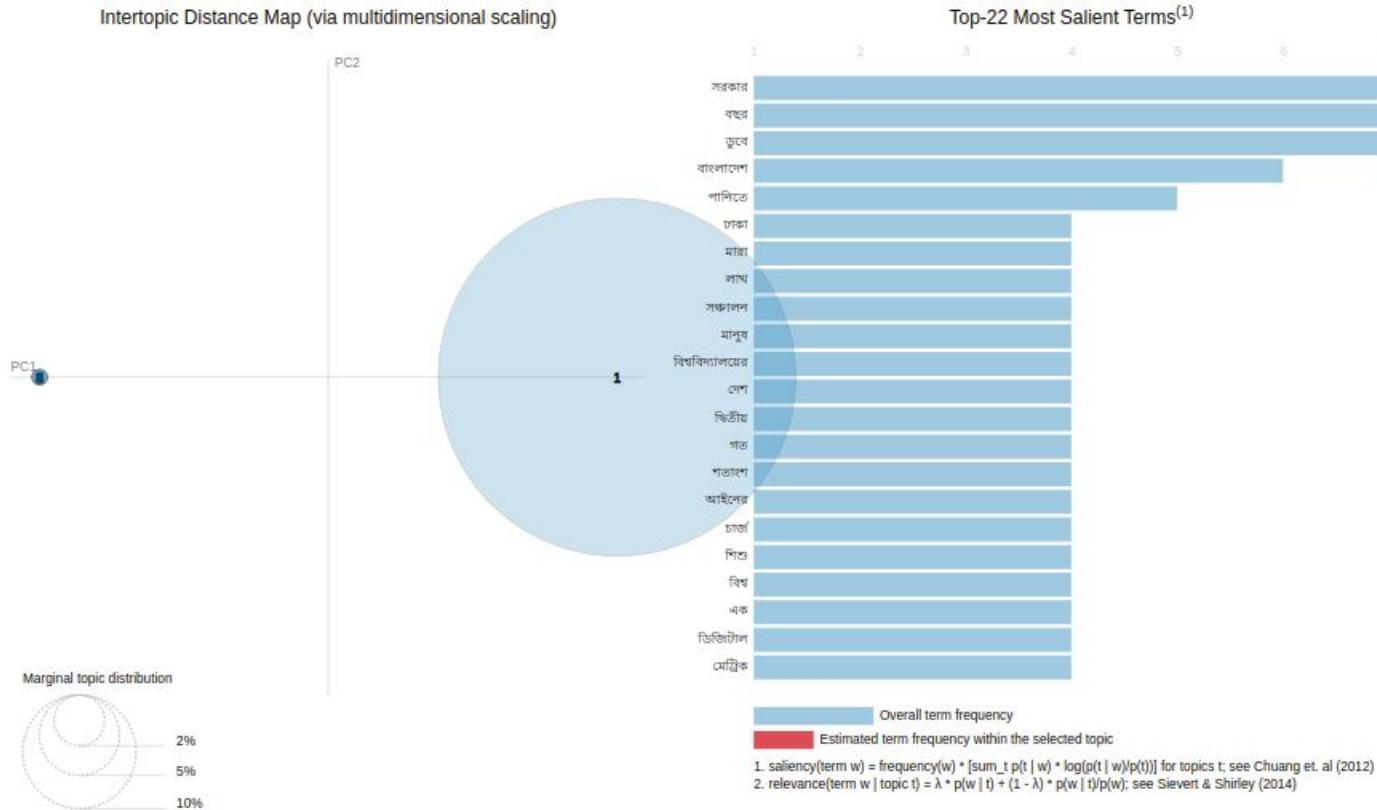
# What was the most popular topic during the off peak user hour on July 24, 2023? (Continued)



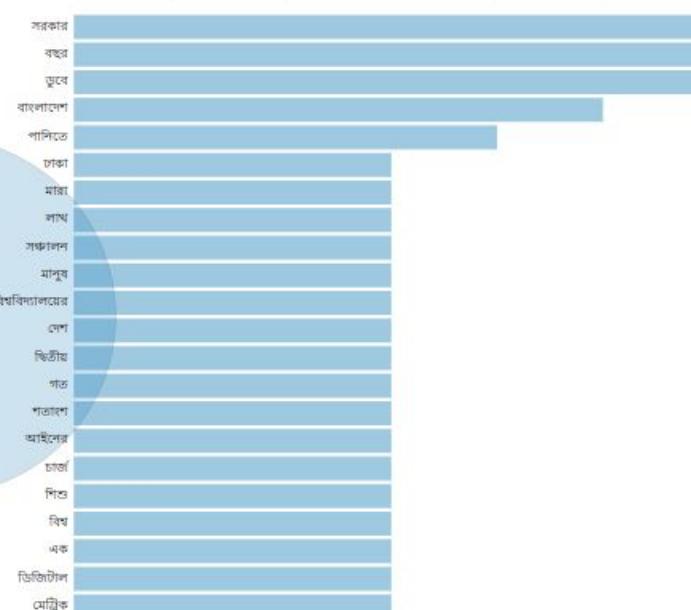
Most trending topic of the day is  
Topic: 1 with probabilities :

$$0.103 * \text{"বাংলাদেশ"} + 0.065 * \text{"বছরেন"} + \\ 0.065 * \text{"গত"} + 0.065 * \text{"তথ্য"} + 0.065 * \text{"বাংলাদেশেন"} + 0.053 * \text{"প্রধান"} + 0.053 * \text{"সরকার"} + 0.053 * \text{"সমাবেশ"} + 0.053 * \text{"বড়"} + 0.053 * \text{"আওয়ামী"}$$

# What was the most popular topic during the off peak user hour on July 25, 2023?



Top-22 Most Salient Terms<sup>(1)</sup>

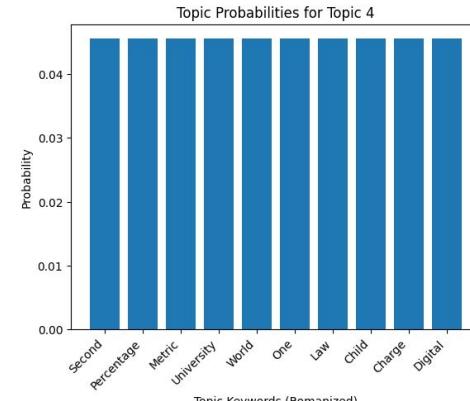
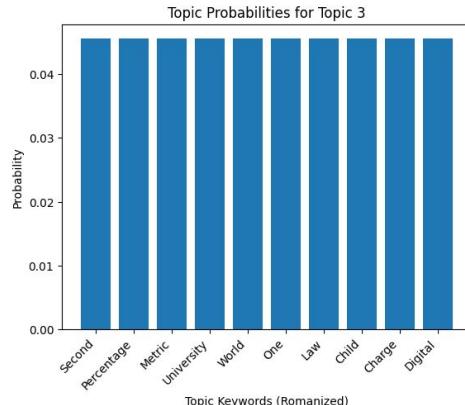
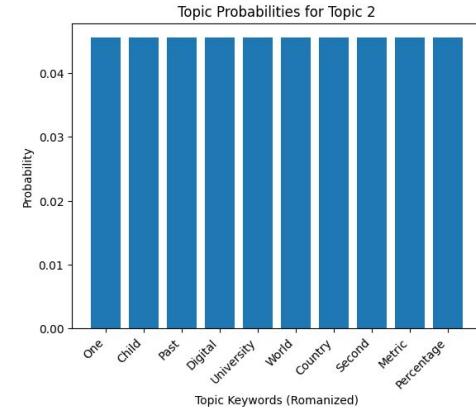
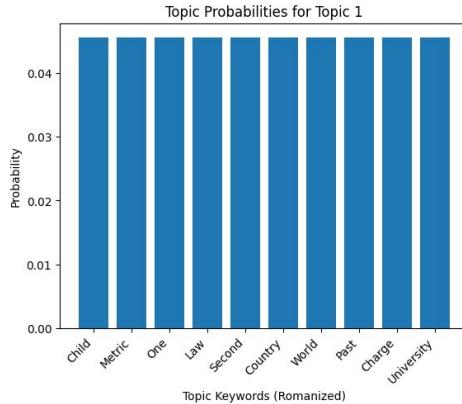
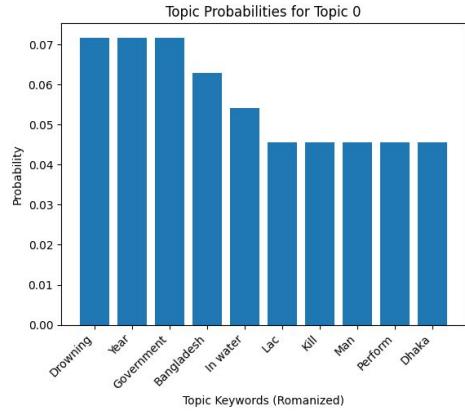


Overall term frequency

Estimated term frequency within the selected topic

1. saliency(term w) = frequency(w) \* [sum\_t p(t | w) \* log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)
2. relevance(term w | topic t) =  $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

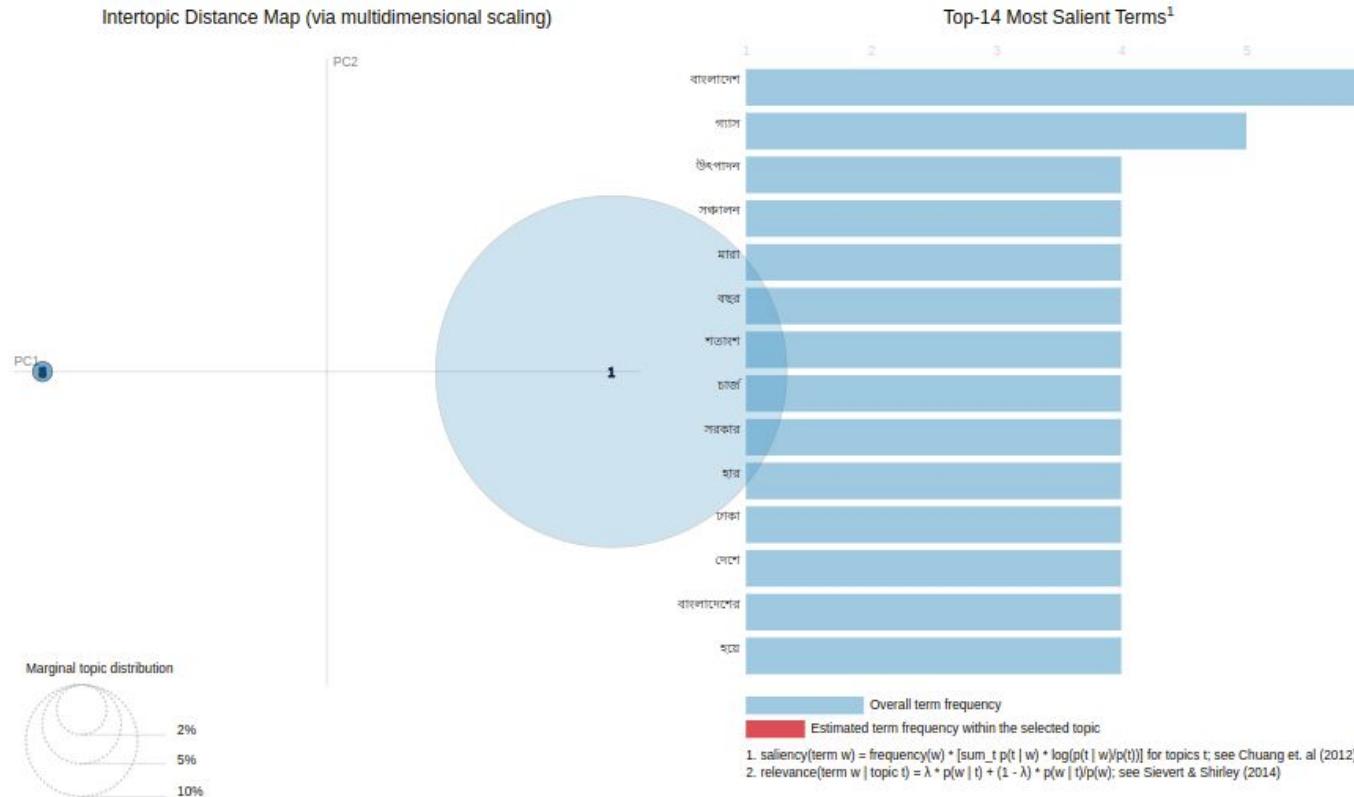
# What was the most popular topic during the off peak user hour on July 25, 2023? (Continued)



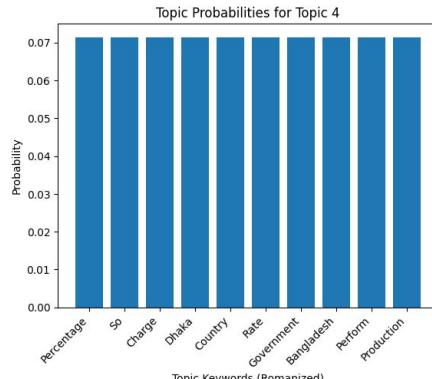
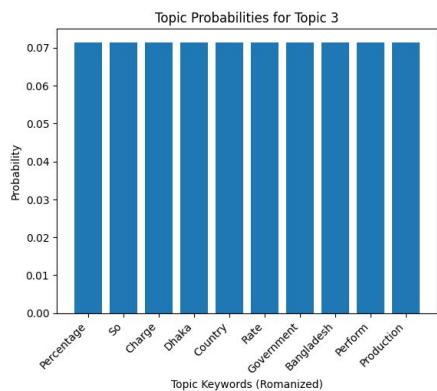
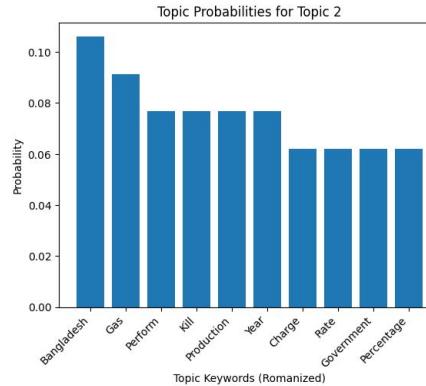
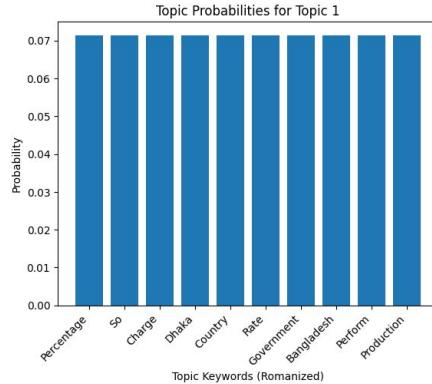
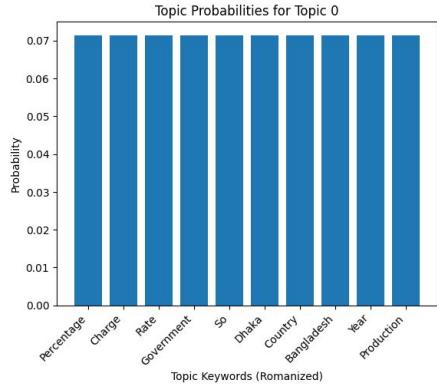
Most trending topic of the day is  
Topic: 0 with probabilities :

$$0.072 * \text{"ডুবে"} + 0.072 * \text{"বছর"} + 0.072 * \text{"সরকার"} + 0.063 * \text{"বাংলাদেশ"} + 0.054 * \text{"পানিতে"} + 0.045 * \text{"লাথ"} + 0.045 * \text{"মারা"} + 0.045 * \text{"মানুষ"} + 0.045 * \text{"সঞ্চালন"} + 0.045 * \text{"ঢাকা"}$$

# What was the most popular topic during the off peak user hour on July 26, 2023?



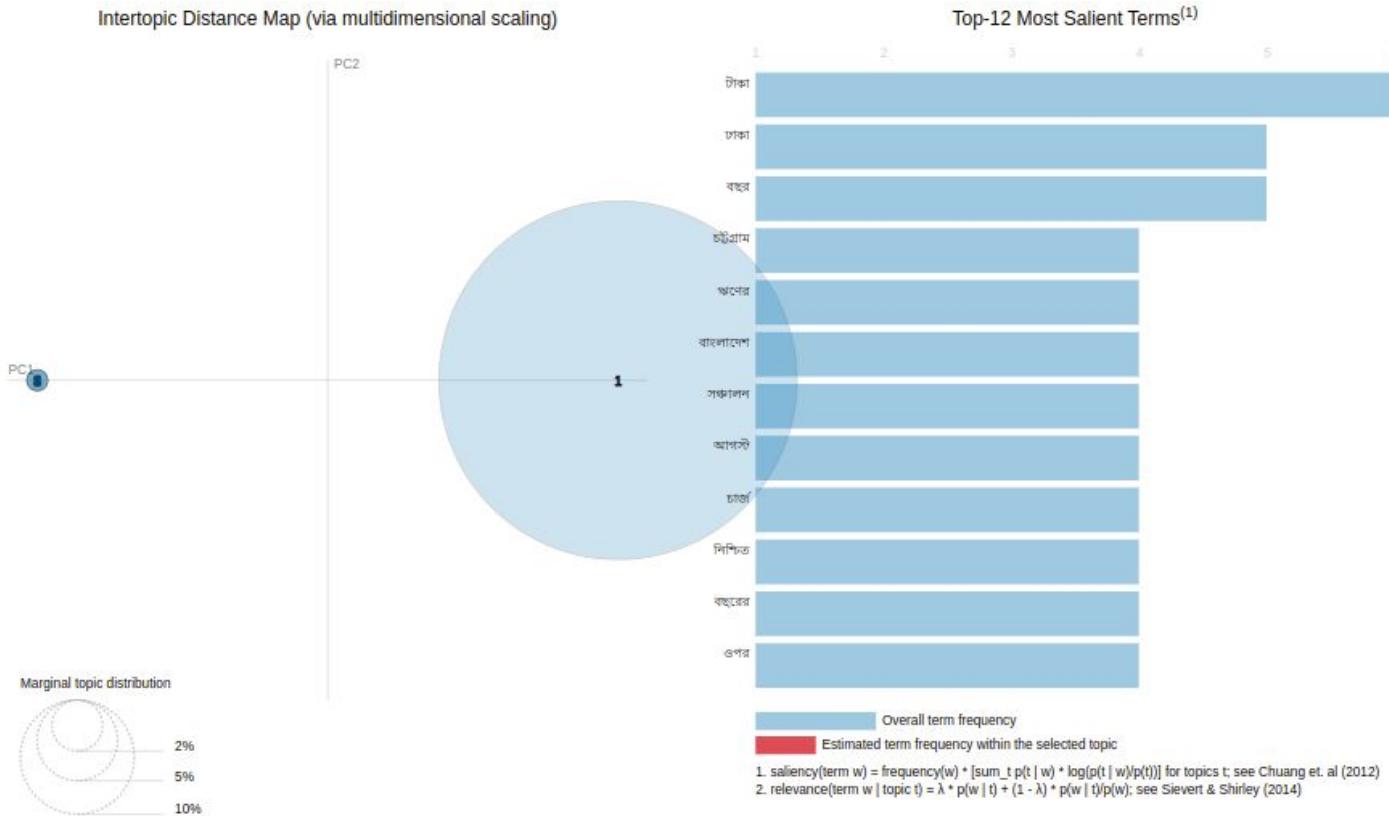
# What was the most popular topic during the off peak user hour on July 26, 2023? (Continued)



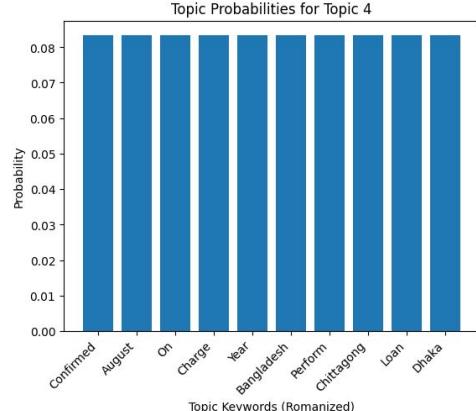
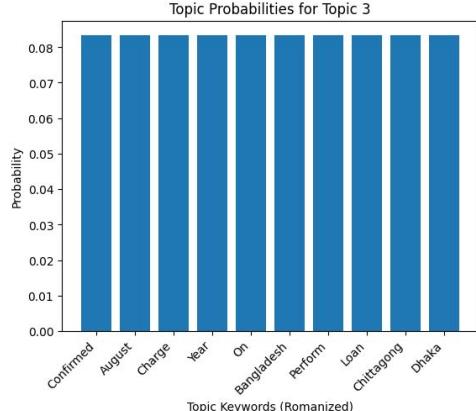
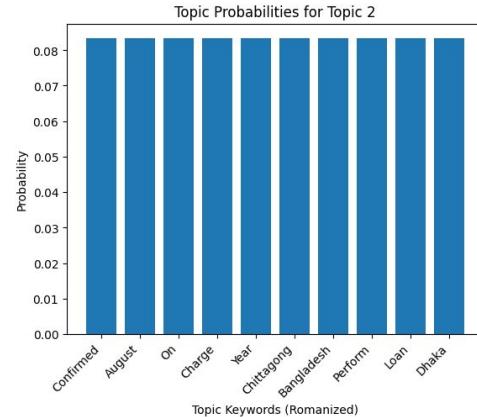
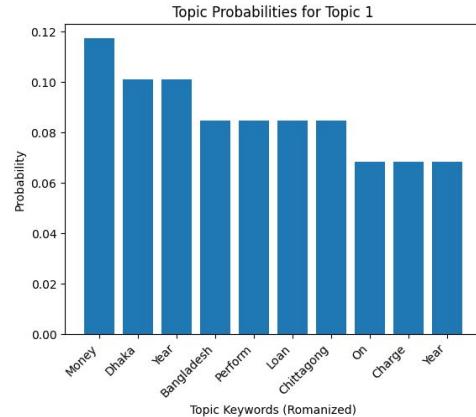
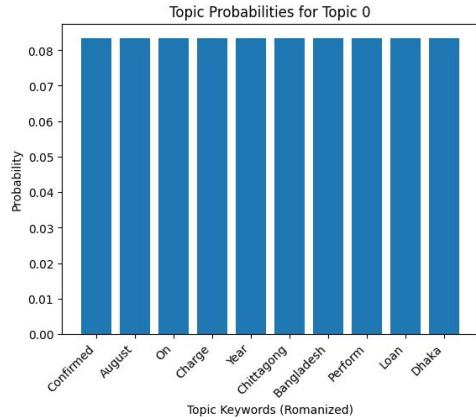
Most trending topic of the day is  
Topic: 2 with probabilities :

0.106 \* "বাংলাদেশ" + 0.091 \* "গ্যাস" +  
0.077 \* "সঞ্চালন" + 0.077 \* "মারা" +  
0.077 \* "উৎপাদন" + 0.077 \* "বছর" +  
0.062 \* "চার্জ" + 0.062 \* "হার" + 0.062 \* "  
সরকার" + 0.062 \* "শতাংশ"

# What was the most popular topic during the off peak user hour on July 27, 2023?



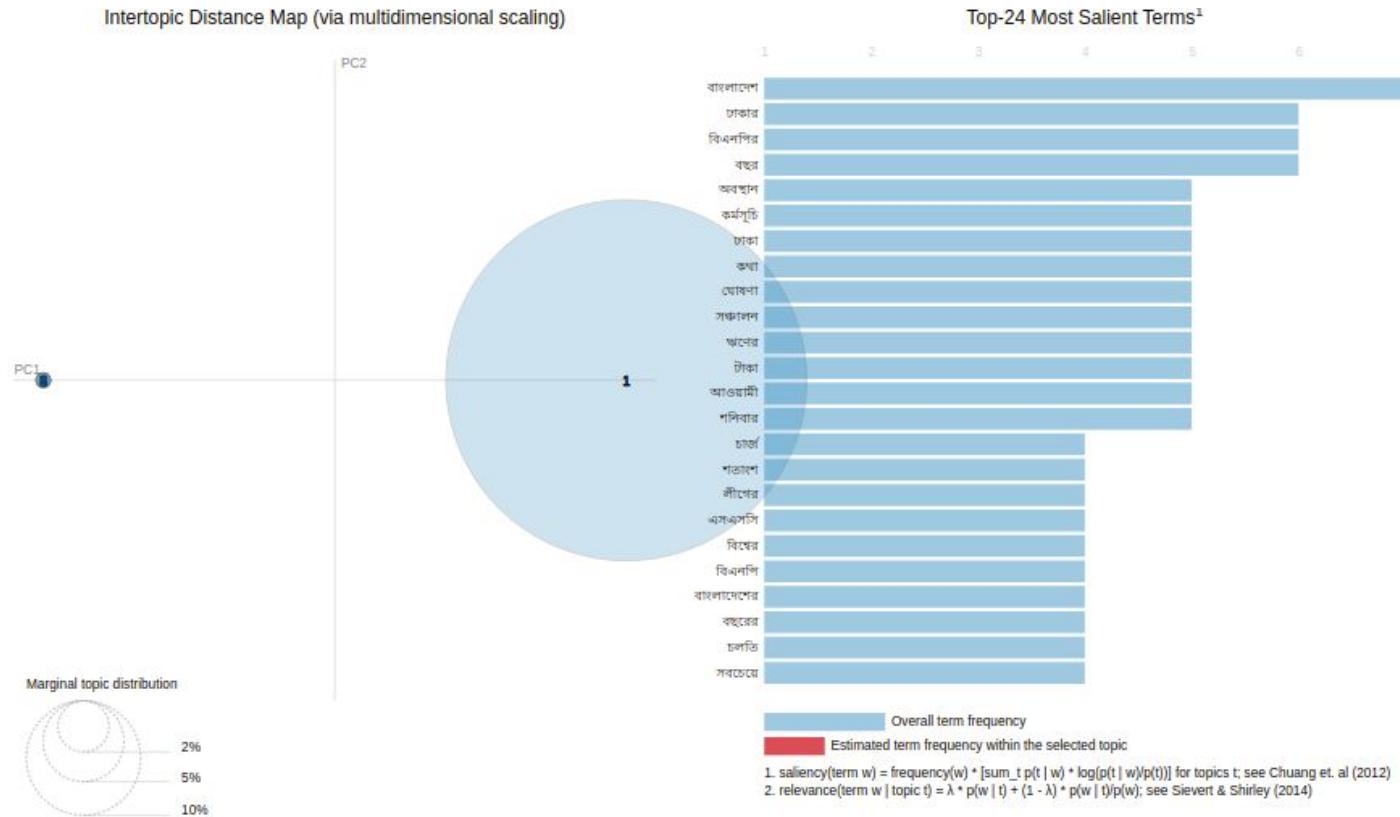
# What was the most popular topic during the off peak user hour on July 27, 2023? (Continued)



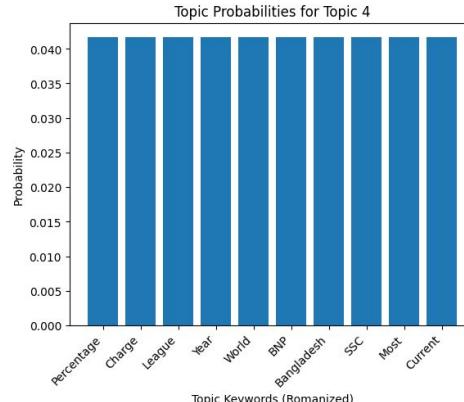
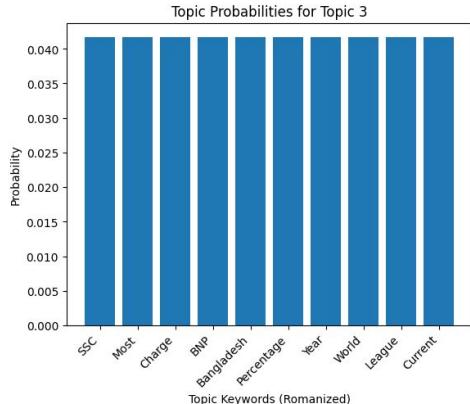
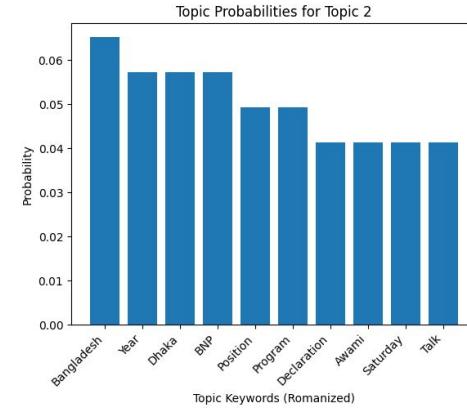
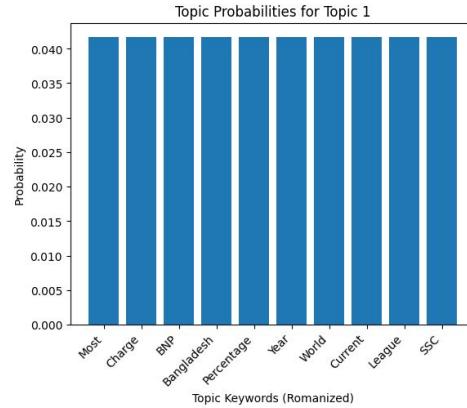
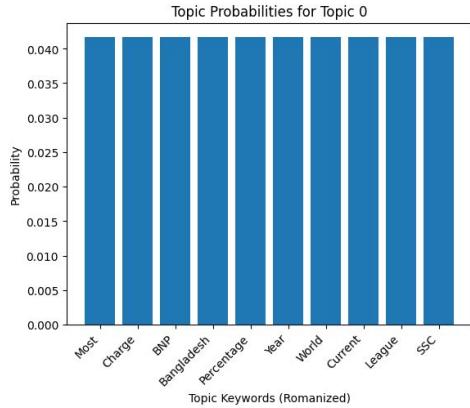
Most trending topic of the day is  
Topic: 1 with probabilities :

0.117 \* "টাকা" + 0.101 \* "ঢাকা" + 0.101 \* "বছর" + 0.085 \* "বাংলাদেশ" + 0.085 \* "সফ্লাইন" + 0.085 \* "ধণের" + 0.085 \* "চট্টগ্রাম" + 0.068 \* "ওপর" + 0.068 \* "চার্জ" + 0.068 \* "বছরেন"

# What was the most popular topic during the off peak user hour on July 28, 2023?



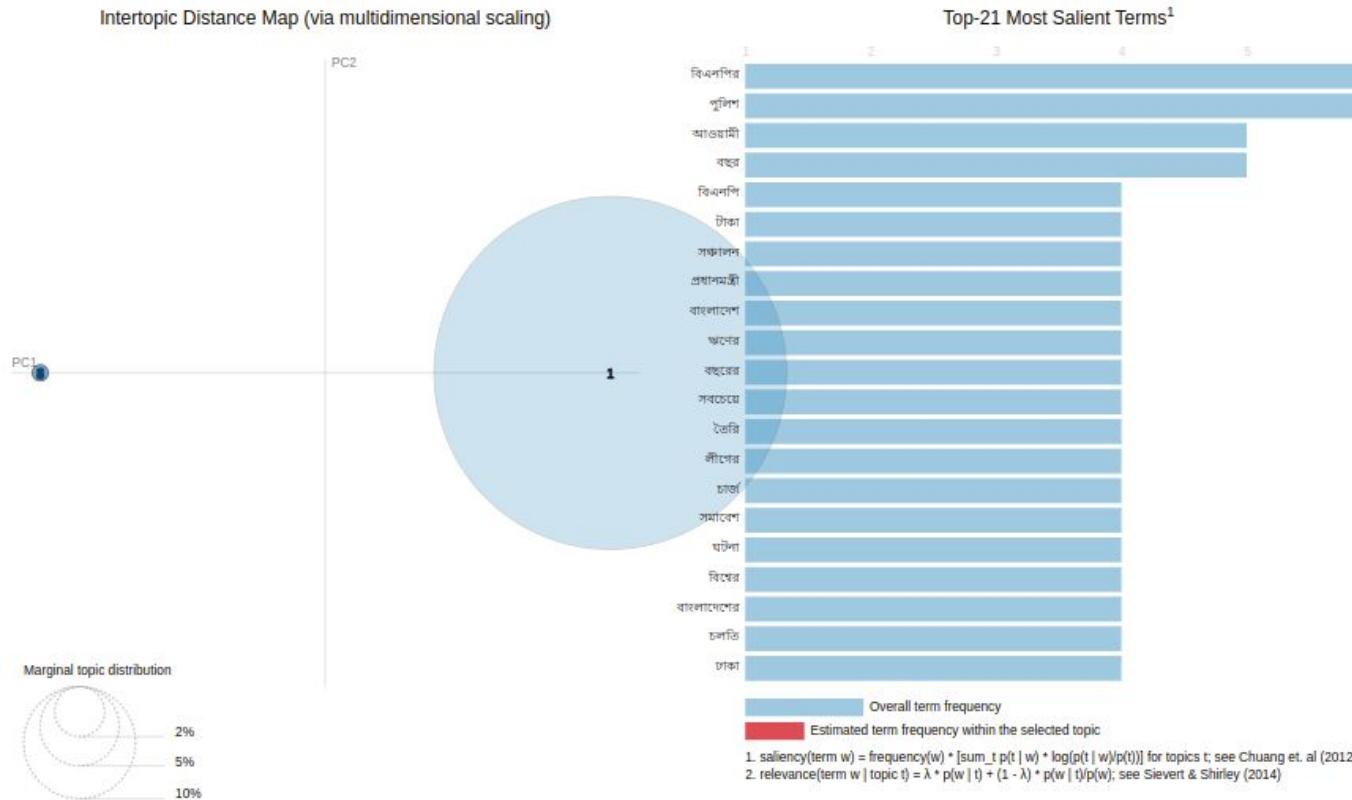
# What was the most popular topic during the off peak user hour on July 28, 2023? (Continued)



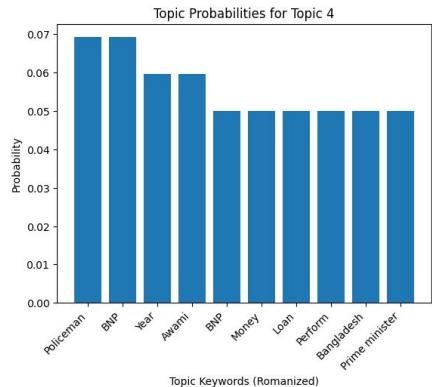
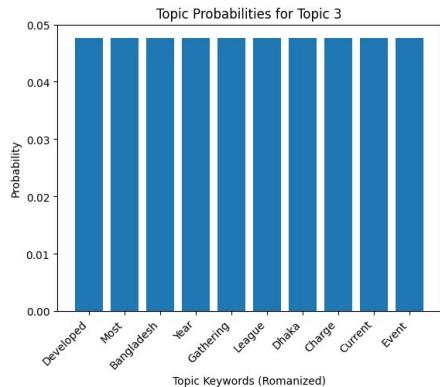
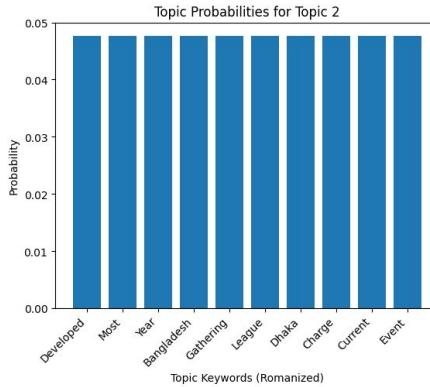
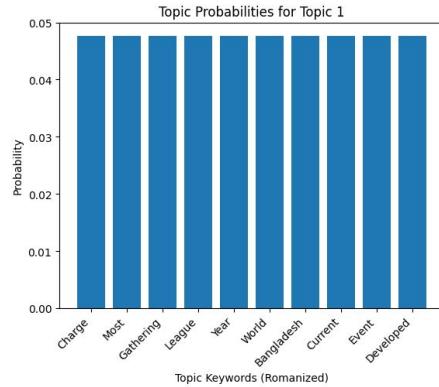
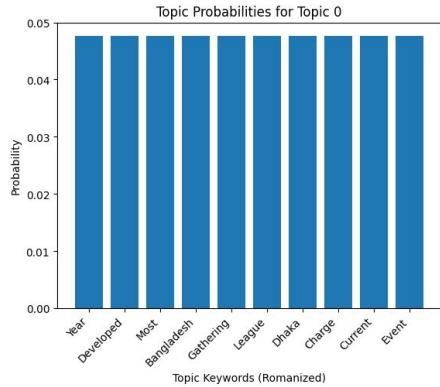
Most trending topic of the day is  
Topic: 2 with probabilities :

$$0.065 * \text{"বাংলাদেশ"} + 0.057 * \text{"বছর"} + \\ 0.057 * \text{"ঢাকাৰ"} + 0.057 * \text{"বিএনপিৱ"} + \\ 0.049 * \text{"অবস্থান"} + 0.049 * \text{"কৰ্মসূচি"} + \\ 0.041 * \text{"ঘোষণা"} + 0.041 * \text{"আওয়ামী"} + \\ 0.041 * \text{"শনিবাৰ"} + 0.041 * \text{"কথা"}$$

# What was the most popular topic during the off peak user hour on July 29, 2023?



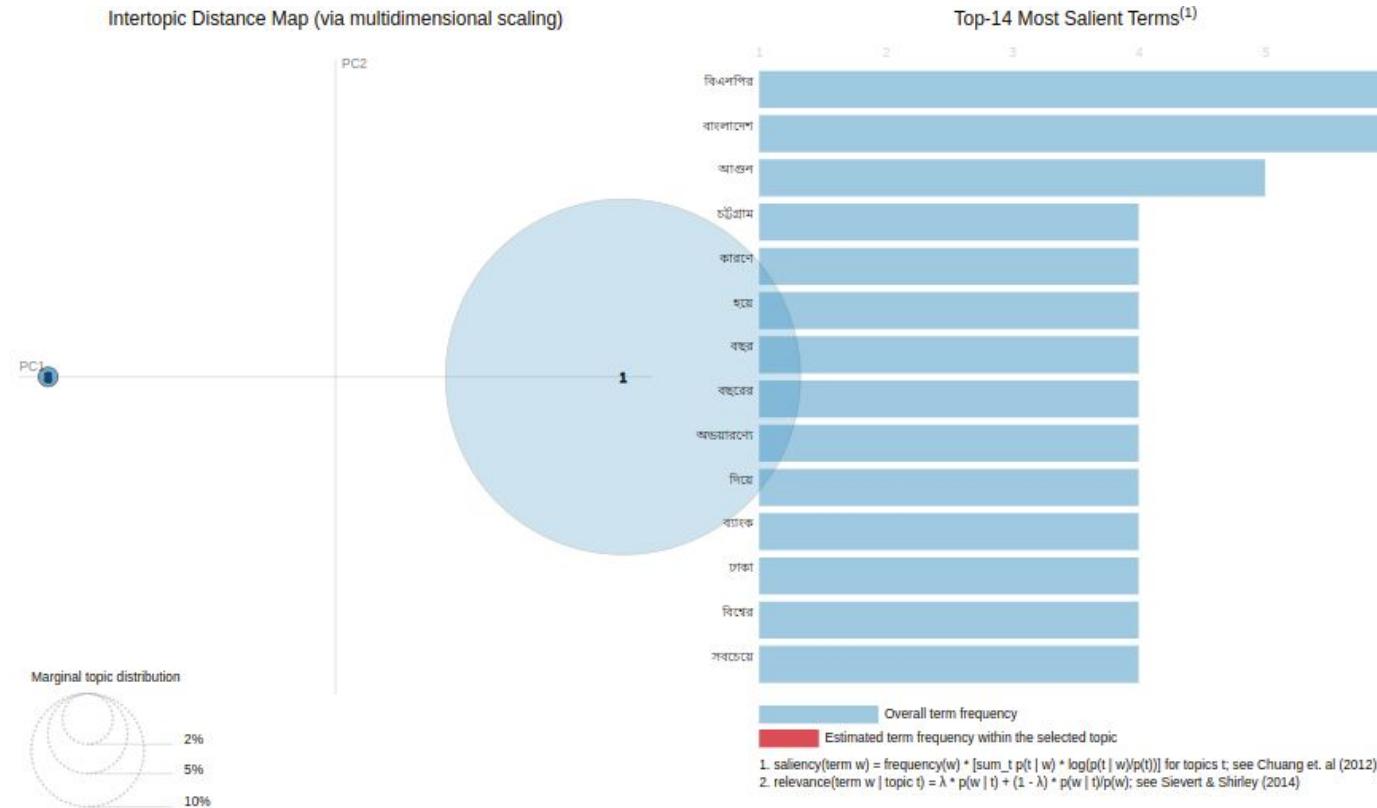
# What was the most popular topic during the off peak user hour on July 29, 2023? (Continued)



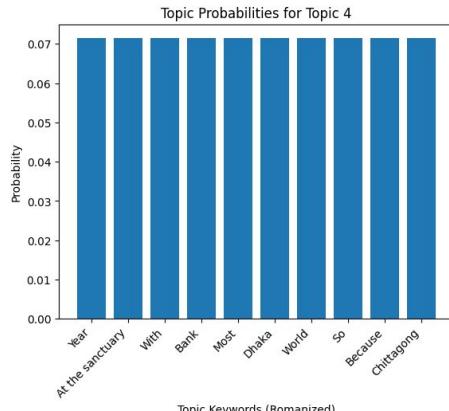
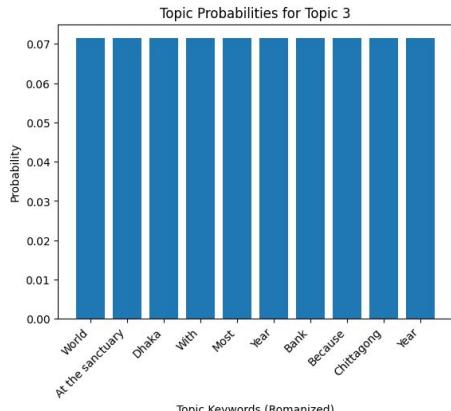
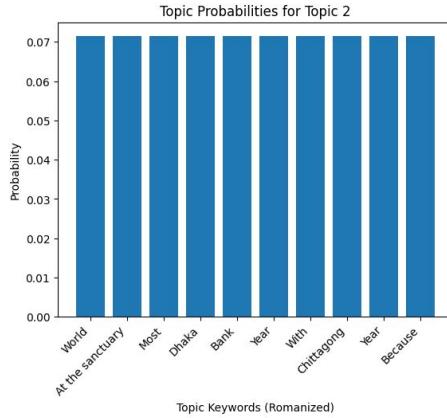
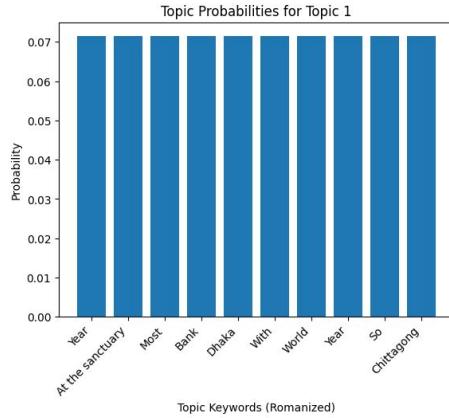
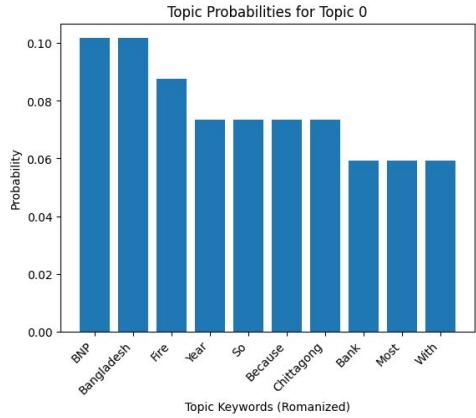
Most trending topic of the day is  
Topic: 4 with probabilities :

$$0.069 * \text{"পুলিশ"} + 0.069 * \text{"বিএনপির"} + \\ 0.060 * \text{"বছর"} + 0.060 * \text{"আওয়ামী"} + \\ 0.050 * \text{"বিএনপি"} + 0.050 * \text{"টাকা"} + \\ 0.050 * \text{"ঝগড়"} + 0.050 * \text{"সঞ্চালন"} + \\ 0.050 * \text{"বাংলাদেশ"} + 0.050 * \text{"প্রধানমন্ত্রী"}$$

# What was the most popular topic during the off peak user hour on July 30, 2023?



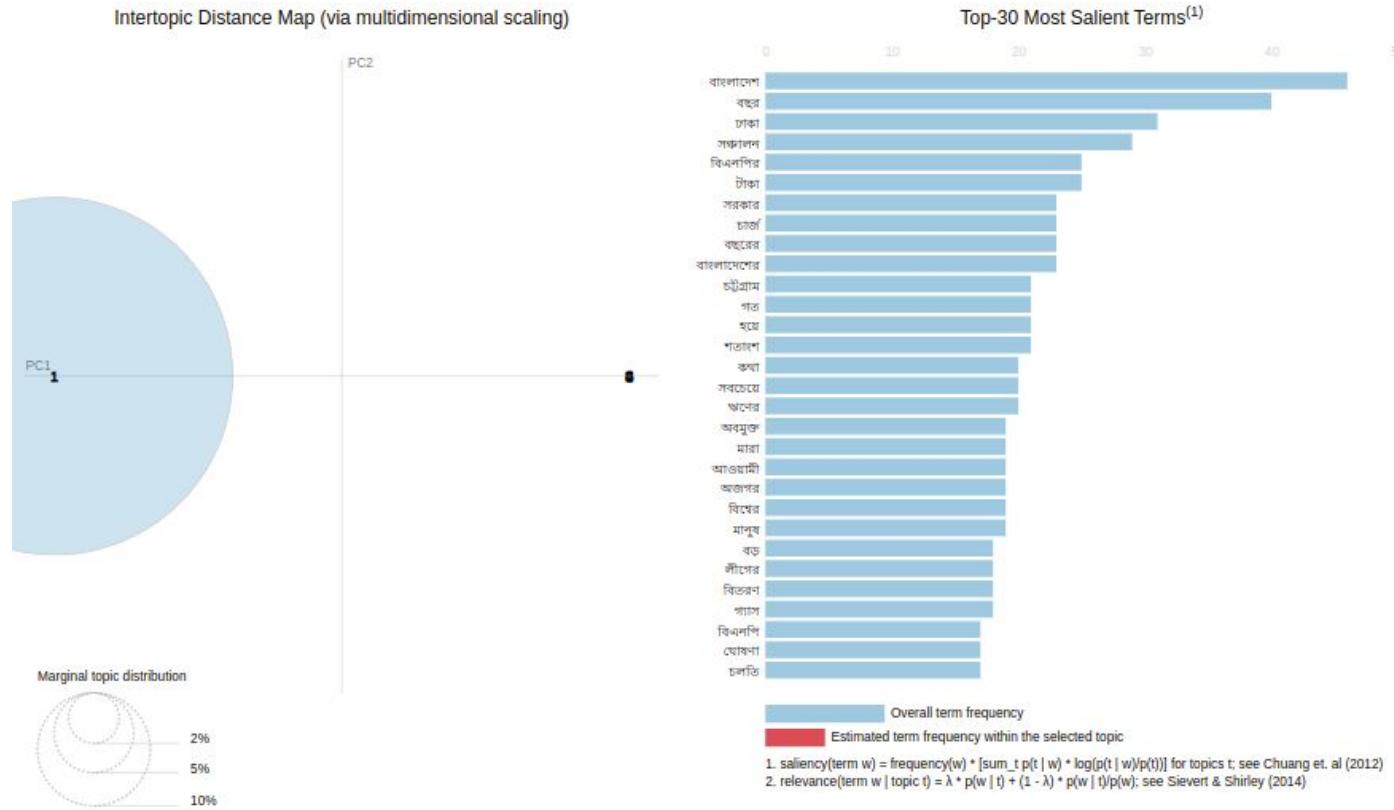
# What was the most popular topic during the off peak user hour on July 30, 2023? (Continued)



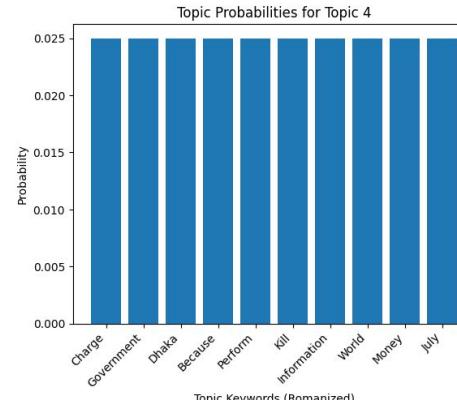
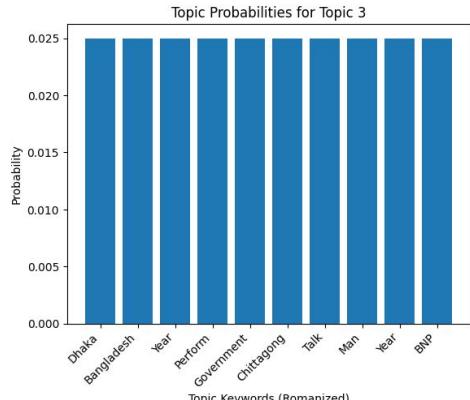
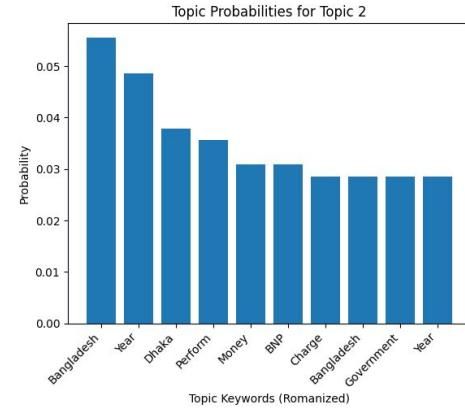
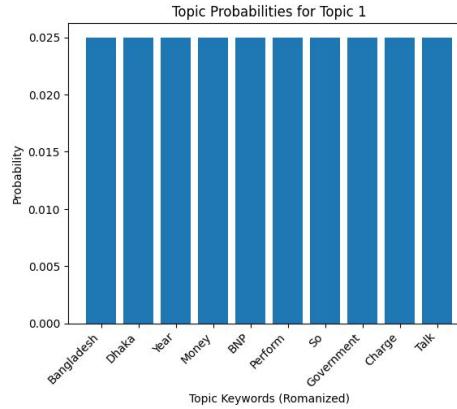
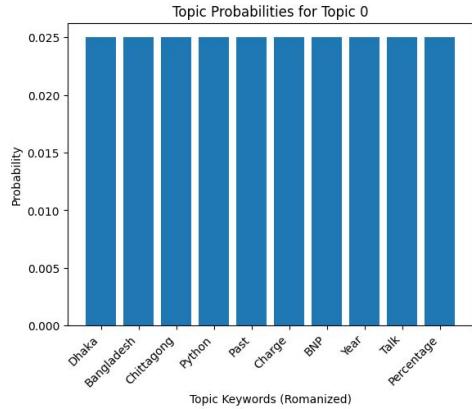
Most trending topic of the day is  
Topic: 0 with probabilities :

$$0.102 * \text{"বিএনপির"} + 0.102 * \text{"বাংলাদেশ"} + \\ 0.088 * \text{"আওলন"} + 0.073 * \text{"বছর"} + \boxed{0.073} \\ 0.073 * \text{"হয়ে"} + 0.073 * \text{"কারণে"} + 0.073 * \text{"চট্টগ্রাম"} + 0.059 * \text{"ব্যাংক"} + 0.059 * \text{"সবচেয়ে"} + 0.059 * \text{"দিয়ে"}$$

# What is the most popular topic during the off peak user hour of this week?



# What is the most popular topic during the off peak user hour of this week? (continued)



Most trending topic of the day is  
Topic: 2 with probabilities :

$$0.056 * \text{"বাংলাদেশ"} + 0.049 * \text{"বছর"} + \\ 0.038 * \text{"ঢাকা"} + 0.036 * \text{"সঞ্চালন"} + \\ 0.031 * \text{"টাকা"} + 0.031 * \text{"বিএনপির"} + \\ 0.029 * \text{"চার্জ"} + 0.029 * \text{"বাংলাদেশের"} + \\ 0.029 * \text{"সরকার"} + 0.029 * \text{"বছরের"}$$

# Descriptions

Kazi Shadman Sakib  
4th Year Undergraduate Student

Department of Computer Science & Engineering  
University of Dhaka



# Dataset Description

The datasets were extracted from The Daily Star **Bangla** website using Python scripts.

**Daily Datasets (Each Day's News):** These datasets contain news from different topics for each day, ranging from **26 to 31** news items. The news topics vary from day to day. This variation in topics and the number of news items per day contributes to the **heterogeneity** of these datasets. Each dataset exhibited a size variation spanning from **10 kB to 26 kB**.

**Weekly Dataset (7 Days' News):** This dataset combines news from all 7 days, encompassing a broader range of topics and potentially a larger number of news items. The mixture of news topics across the week further adds to the **heterogeneity** of this dataset. Each dataset exhibited a size variation spanning from **76 kB to 175 kB**.

# Model Description

## Data Collection and Preprocessing :

- The Bengali stemmer function ('**bengali\_stem(word)**') is defined to remove common suffixes from words for stemming.
- The '**normalize(text)**' function is defined to preprocess and normalize the text data, converting it to lowercase and tokenizing it while removing stopwords and numerical values.
- The code reads the Bengali dataset from a file and applies text normalization to the entire text content.

## Topic Modeling with LDA:

- The normalized text is tokenized, and a dictionary is created from the tokens using the **corpora.Dictionary** class.
- A bag of words representation (corpus) is generated using the created dictionary for the LDA model input.
- The **perform\_lda\_topic\_modeling** function is defined to execute LDA topic modeling using the **Gensim** library. The function takes the corpus, dictionary, number of topics, and number of passes as parameters.
- Inside the **perform\_lda\_topic\_modeling** function, the LDA model is trained on the provided corpus and dictionary using the **gensim.models.LdaModel** class. The top keywords with probabilities for each topic are printed using **Lda\_model.show\_topics**.

# Model Description (Continued)

After extensive experimentation and parameter tuning, I have refined the following predefined variables to optimize the performance of the Latent Dirichlet Allocation (LDA) model for topic modeling in the context of Bengali news articles:

- **num\_topics = 5**: This predefined variable represents the number of distinct topics that the LDA model aims to identify within the dataset. Through careful adjustment, I have settled on five topics, ensuring a balance between granularity and coherence in the extracted themes.
- **passes = 80**: The predefined variable passes is indicative of the number of iterations employed during the convergence process of the LDA model. By setting it to 80, I have established a convergence threshold that enables the model to sufficiently capture the underlying topic patterns while avoiding overfitting.
- **num\_words = 10**: This predefined variable controls the quantity of top keywords associated with each identified topic that are displayed. After iterative experimentation, I have chosen to display 10 keywords per topic. This value facilitates meaningful insights while maintaining concise topic representations.

# Model Performance

The LDA model has successfully processed the dataset of Bengali news articles from The Daily Star Bangla and identified **five distinct topics** based on the analysis of the text content. Each topic is represented by a set of top keywords along with their corresponding probabilities. The associated probabilities indicate the relevance of each keyword within its respective topic.

The model's performance in conducting topic extraction from Bengali news articles **showcases its efficacy in identifying distinct subject clusters**. Furthermore, it has the capability to retrieve relevant and frequently discussed topics along with their corresponding keywords for each dataset. This applies equally to distinct datasets encompassing daily news, as well as separate datasets focusing on peak user hour news, off-peak user hour news, weekly news, peak user hour weekly news, and off-peak user hour weekly news.

# Model Performance (Continued)

Evaluation of Model Performance across Three Distinct Datasets :

	<b>Weekly News</b>	<b>Peak User Hour Weekly News</b>	<b>Off Peak User Hour Weekly News</b>
<b>Best Topic</b>	4	3	2
<b>Upper bound of the most relevant topics based on keyword probabilities</b>	0.131, 0.131, 0.092, 0.085, 0.082, 0.082, 0.075, 0.075, 0.064, 0.061	0.060, 0.047, 0.039, 0.037, 0.034, 0.033, 0.031, 0.031, 0.030, 0.029	0.056, 0.049, 0.038, 0.036, 0.031, 0.031, 0.029, 0.029, 0.029, 0.029
<b>Lower bound of the most relevant topics based on keyword probabilities</b>	0.083, 0.083, 0.083, 0.083, 0.083, 0.083, 0.083, 0.083, 0.083, 0.083	0.027, 0.027, 0.027, 0.027, 0.027, 0.027, 0.027, 0.027, 0.027, 0.027	0.025, 0.025, 0.025, 0.025, 0.025, 0.025, 0.025, 0.025, 0.025, 0.025

# Findings

1. The findings suggest a pattern such as specific topics in the weekly dataset related to Awami League and BNP, discussed in the last 3-4 days, continue to have a notable impact on the overall dataset of the past 7 days. This influence is reminiscent of the final results, indicating a recurrent focus on discussions concerning Awami League and BNP over the span of the preceding week.
2. Notably, certain topics consistently experience heightened discussion and engagement throughout 7 days duration. These prominently recurring topics, accompanied by their associated keywords, are subsequently integrated into the final results (for all the datasets). This strategic inclusion of recurrently emphasized subjects sheds light on their enduring significance within the dataset.
3. The meticulously fine-tuned values of the predefined variables (num\_topics, passes and num\_words) enhance the model's ability to discern relevant and coherent topics from Bengali news articles, thereby contributing to an informed understanding of the dataset's content.
4. The emergence of redundant topics is a consequence of setting the predefined variable "num\_topics" to a value exceeding 5. In cases where the predefined variable "passes" is set significantly below 80, the model convergence may not reach an optimal state. The presence of redundant keywords arises when the assigned value to "num\_words" surpasses 10. This phenomenon can be attributed to the inherent data heterogeneity, which impacts the efficacy of topic extraction and the model's convergence behavior.

# Thank You

Any Questions?

