

Final Project

Assignment Title:	Text Analysis and Topic Modeling		
Assignment No:		Date of Submission:	30 May 2025
Course Title:	Introduction to Data Science		
Course Code:		Section:	G
Semester:	Spring	2024-25	Course Teacher: DR. ASHRAF UDDIN

Declaration and Statement of Authorship:

1. I/we hold a copy of this Assignment/Case-Study, which can be produced if the original is lost/damaged.
2. This Assignment/Case-Study is my/our original work and no part of it has been copied from any other student's work or from any other source except where due acknowledgement is made.
3. No part of this Assignment/Case-Study has been written for me/us by any other person except where such collaboration has been authorized by the concerned teacher and is clearly acknowledged in the assignment.
4. I/we have not previously submitted or currently submitting this work for any other course/unit.
5. This work may be reproduced, communicated, compared and archived for the purpose of detecting plagiarism.
6. I/we give permission for a copy of my/our marked work to be retained by the faculty for review and comparison, including review by external examiners.
7. I/we understand that Plagiarism is the presentation of the work, idea or creation of another person as though it is your own. It is a form of cheating and is a very serious academic offence that may lead to expulsion from the University. Plagiarized material can be drawn from, and presented in, written, graphic and visual form, including electronic data, and oral presentations. Plagiarism occurs when the origin of their material used is not appropriately cited.
8. I/we also understand that enabling plagiarism is the act of assisting or allowing another person to plagiarize or to copy my/our work.

* Student(s) must complete all details except the faculty use part.

** Please submit all assignments to your course teacher or the office of the teacher concerned.

Group Name/No.: 07

No	Name	ID	Program	Contribution
1	ARSHAD ABEDIN ABIR	22-47188-1	BSc [CSE]	Text Preprocessing
2	KAMRUZZAMAN SONY	22-46797-1	BSc [CSE]	Topic Modeling
3	MRIDUL KANTI KUNDU	22-47182-1	BSc [CSE]	Exploratory Text Analysis
4	MAYSHA MARIUM	22-47197-1	BSc [CSE]	Web Scraping
5			Choose an item.	
6			Choose an item.	
7			Choose an item.	
8			Choose an item.	
9			Choose an item.	
10			Choose an item.	

Faculty use only

FACULTY COMMENTS	Marks Obtained	
	Total Marks	

1.Introduction

The purpose of this project is to analyze real-world news articles through web scraping, text preprocessing, visual analytics, and topic modeling. We have extracted news content from the Bengali news portal <https://www.banglanews24.com/>, which provides a rich source of current news articles in Bengal. It involves extracting news content from a Bengali news portal, cleaning and analyzing the text, creating meaningful visualizations, and uncovering hidden topics within the news dataset. This approach aims to gain valuable insights from the information and apply them to natural language data sourced from online media, thereby enhancing understanding of the news content.

2.Web Scraping

The first step of the project involves selecting a news portal, preferably a Bengali news website, to extract news articles. The web scraping process entails collecting news article texts from multiple pages linked on the homepage of the portal. After successfully scraping the data, the extracted content will be saved into a CSV file. This file will contain at least two essential columns: the URL of the news article and the corresponding article text. This structured data will then serve as the basis for further preprocessing and analysis.

Data Preview:

url	article_text
(character) ▼	(character) ▼
https://www.banglanews24.com/banglanews-special...	রাজধানীর প্রাণকেন্দ্রে নির্মল বাতাসে প্রাণভরে শ্বাস নেওয়ার জা...
https://www.banglanews24.com/politics/news/bd/15...	ঢাকা: প্রধান উপদেষ্টা অধ্যাপক ড. মুহাম্মদ ইউনূসের একটি সা...
https://www.banglanews24.com/climate-nature/new...	ঢাকা: স্থল গভীর নিম্নচাপ ঘীরে ঘীরে দুর্বল হলেও এর প্রভাবে এ...
https://www.banglanews24.com/national/news/bd/1...	প্রধান উপদেষ্টা অধ্যাপক ড. মুহাম্মদ ইউনূস ও জাপানের প্রধা...
https://www.banglanews24.com/environment-biodiv...	ঢাকা: গভীর নিম্নচাপ স্থলভাগে উঠে আসার পর দেশে মৌসুমে...
https://www.banglanews24.com/saradesh/news/bd/...	নোয়াখালী: বঙ্গোপসাগরে সৃষ্ট গভীর নিম্নচাপের প্রভাবে নোয়াখা...
https://www.banglanews24.com/cricket/news/bd/15...	জাতীয় ক্রীড়া পরিষদ (এনএসসি) কর্তৃক ফারুক আহমেদের প...
https://www.banglanews24.com/politics/news/bd/15...	ঢাকা: সমগ্র জাতি একটি সুস্ব্ঠ, সুন্দর, নিরপেক্ষ নির্বাচনের জ...
https://www.banglanews24.com/opinion/news/bd/15...	শহীদ প্রেসিডেন্ট জিয়াউর রহমান আমাদের জাতীয় ইতিহাসের ...
https://www.banglanews24.com/saradesh/news/bd/...	ভোলা: নিম্নচাপের প্রভাবে ভোলায় নদ-নদীর পানি অস্বাভাবিক...
https://www.banglanews24.com/saradesh/news/bd/...	দিনাজপুর: দিনাজপুরের পাবতীপুরে ট্রাক আটকে চাঁদাবাজির স...
https://www.banglanews24.com/daily-chittagong/ne...	চট্টগ্রাম: জোড়া লাগানো যমজ শিশু রিয়াশাদ ও রেনিশকে জ...
https://www.banglanews24.com/politics/news/bd/15...	ঢাকা: মহান স্বাধীনতার ঘোষক ও বি.এন.পির প্রতিষ্ঠিতা শহীদ পে...

Previewing first 50 entries.

Figure 01: CSV Preview of Scraped News Articles

3.Text Preprocessing

Text preprocessing is a fundamental step to preparing raw text data for meaningful analysis. It involves cleaning and transforming the text to reduce noise, standardizing the format, and converting it into a suitable form for modeling. The main steps include cleaning the text, removing stop words, and tokenization.

Before Text Preprocessing:

অবশেষে নাটকীয়তার অবসান ঘটলো। সাবেক অধিনায়ক আমিনুল ইসলাম বুলবুল নির্বাচিত হয়েছেন বিসিবির নতুন সভাপতি হিসেবে।

```
googletag.cmd.push(function () {
  googletag.display('div-gpt-ad-5565653-11');
});
```

বিকলে অনুষ্ঠিত জরুরি সভায় পরিচালকদের ভোট এই সিদ্ধান্ত গৃহীত হয়। বিসিবির পরবর্তী নির্বাচন পর্যন্ত তিনি এই পদ দায়িত্ব পালন করবেন। একই সভায় সিনিয়র সহ-সভাপতির দায়িত্ব পেয়েছেন অ

ভিন্ন ভিন্ন একটি প্রস্তাবন জারি করা হয়। প্রস্তাবন উল্লেখ করা হয়, বিসিবির গঠনতন্ত্রের ১০.২ (খ)(৪) ধারা অনুসারে জাতীয় ক্রীড়া পরিষদ কর্তৃক আমিনুলকে পরিচালক হিসেবে মনোনীত করা হয়েছে। এর আগে গতকাল রাত ১১টায় জাতীয় ক্রীড়া পরিষদ (এনএসসি) থেকে এক বিজ্ঞপ্তিতে জানানো হয়, বিসিবির পরিচালক হিসেবে ফারুক আহমেদের মনোনয়ন বাতিল করা হয়েছে। বিসিবির পরবর্তী নির্বাচন অক্টোবরের মধ্যে সম্পন্ন হওয়ার কথা রয়েছে। তার মানে, চার মাসের মতো দায়িত্ব থাকতে পারেন এই সাবেক অধিনায়ক, এরপর তিনি ফিরতে চান তার আন্তর্জাতিক ক্রিকেট সংস্থার (আইসিসি) চাকরিতে। বাংলাদেশের ক্রিকেট ইতিহাসে এক অনন্য নাম বুলবুল। ১৯৯৯ সালের বিশ্বকাপ বাংলাদেশ দলের অধিনায়ক ছিলেন তিনি। দেশের ইতিহাসের প্রথম টেস্ট সেক্সট্রিউ এসেছে তার বাট থেকে। খেলার মাঠ থেকে বিদায় নেওয়ার পর তিনি পাড়ি জমান অষ্ট্রেলিয়ায়, কোচিংয়ে লেভেল টু সম্পন্ন করে যুক্ত দ্ব দি নিউ সাউথ ওয়েলস ইউনিভার্সিটিতে কোচিং পাবেন। পরে দেশে ফিরে আরাব্বীক উপহার দেন ঢাকা প্রিমিয়ার লিগের শিরোপা। আরইউ ঢাকা: ছাত্র-জনতার অভ্যুত্থানের অন্যতম নেতা ও জাতীয় নাগরিক পার্টির (এনসিপি) আধায়ক মো. নাহিদ ইসলামের জাতীয় পরিচয়পত্র (এনআইডি) উপদেষ্টা পদে থাকারদ্বায় লক করে রেখেছিল নির্বাচন কমিশন (ইসি)। সুকণ্ঠা জানিয়েছে, এনআইডি তথা ফাঁসের অভিযোগের ভিত্তিতে ওই সিদ্ধান্ত নিয়েছিল সংস্থাটি। ২০২৪ সালে ৩৬ দিনের অভ্যুত্থান আন্দোলনের মুখে প্রধানমন্ত্রীর পদ থেকে শেখ হাসিনা পদত্যাগ করে দেশ ছেড়ে ও আগষ্ট ভারতে পালিয়ে গেলে যে অন্তর্বর্তী সরকার গঠন হয়, সেখানে তথ্য ও সম্প্রচার মন্ত্রণালয় এবং ভাঙ ও টেলিযোগাযোগ মন্ত্রণালয়ের উপদেষ্টা পদে আসীন হন নাহিদ। আর সেই পদে থাকাকালীনই তার বিরুদ্ধে একটি অনলাইন প্রাটিকর্ম এনআইডি তথ্য ফাঁস করার অভিযোগ আসে। জানা গেছে, ন্যাশনাল টেলিকমিউনিকেশন মনিটরিং সেন্টারের (এনটিএমসি) তৎকালীন মহাপরিচালক মো. নাহিদ ইসলামের বিরুদ্ধে ‘ভদ্র বারা’ নামে এক হোয়াটসঅ্যাপ গ্রুপের মাধ্যমে নাগরিকদের তথ্য সংগ্রহ করে বাইরে পাচারের অভিযোগ আনেন। সেই অভিযোগ আমলে নিয়ে তদন্ত করার সিদ্ধান্ত নেয় ইসি। তবে সহকারী প্রোগ্রামার আমিনুল ইসলামের নেতৃত্বে গঠিত দুই সদস্যের তদন্ত কমিটি প্রতিবেদন দেওয়ার আগেই নাহিদের এনআইডি লক করার সুপারিশ করেন সিস্টেম অ্যানালিস্ট মোহাম্মদ আরিফুল ইসলাম। সেই সুপারিশের ভিত্তিতে গত বছরের ১৭ সেপ্টেম্বর এনআইডি লক করার নির্দেশ দেন তৎকালীন জাতীয় পরিচয় নিবন্ধন (এনআইডি) অনুবিভাগের মহাপরিচালক মো. মাহবুব আলম তালুকদার। এদিকে তদন্ত কমিটি দুদিন পর অর্থাৎ ২০২৪ সালের ১৯ সেপ্টেম্বর প্রতিবেদন দাখিল করে। এতে উল্লেখ করা হয়, এনটিএমসি-এর মহাপরিচালক কর্তৃক যৌথিকভাবে অভিযোগ পাওয়া যায় যে, মো. নাহিদ ইসলাম ‘ভদ্র বারা’ নামক হোয়াটসঅ্যাপ গ্রুপের এডমিন হিসেবে কাজ করে নাগরিকদের তথ্যাদি বাহিরে সরবরাহ করেন। তদন্ত চলাকালীন সময়ে অনলাইনের বিভিন্ন সোশ্যাল মিডিয়ায় গ্রুপে যোগে রাষ্ট্রের বিরুদ্ধে তথ্যাদি সরবরাহ করেন এমন অভিযোগটি সভা নয় এবং ভাঙা সরবরাহ করার বিষয়ে তার কোনো সম্পৃক্ততা পাওয়া যায়নি। প্রতিবেদন বলা হয়, যেহেতু ওই ভোটার কর্তৃক ভাঙা সরবরাহ করার বিষয়ে কোনো সম্পৃক্ততা পাওয়া যায়নি এবং অভিযোগটি মিথ্যা প্রমাণিত হয়েছে, সেহেতু মো. নাহিদ ইসলামের এনআইডি আনলক করার জন্য সুপারিশ করে কমিটি। তাদের প্রতিবেদনের ভিত্তিতে নথি ইশ্তারাদ করা হলে ২২ সেপ্টেম্বর এনআইডি মহাপরিচালক তখন এনআইডিটি আনলক করার সিদ্ধান্ত দেন। এভাবেই পাঁচ দিনের জন্য লক থাকে সাবেক উপদেষ্টা নাহিদ ইসলামের এনআইডি। ইসি কর্মকর্তারা জানিয়েছেন, এনটিএমসি থেকে এনআইডি নম্বর পাঠিয়ে সংশ্লিষ্ট ভোটারের বিরুদ্ধে তথ্য পাচারের অভিযোগ তোলা হয়। আর সেই এনআইডি যে উপদেষ্টা নাহিদ ইসলামের সংশ্লিষ্টরা সেটা বুঝতে পারেননি। বরং লক করার পর রোধমান হয়ে যে এটা কার এনআইডি। পরবর্তীতে ফ্রন্টই আবার সেটা আনলক করা হয়। এই কাজের সঙ্গে যুক্ত দায়িত্বশীল একজন কর্মকর্তা বলেন, বিষয়টি সঠিক যে একজন ছাত্র উপদেষ্টার এনআইডি তখন লক করা হয়েছিল। তবে খুব দ্রুত সময়ের ব্যবধানে আবার আনলক করা হয়। ৫ আগষ্ট আগুয়ানী লীগ সরকার পড়নের চার মাস আগে ১ এপ্রিল বেসামরিক বিমান চলাচল কর্তৃপক্ষের সদস্য অভিযুক্ত সচিব মাহবুব আলম তালুকদারকে এনআইডি মহাপরিচালক হিসেবে পদায়ন করা হয়। রাজনৈতিক পট পরিবর্তনের পর ৬ নভেম্বর বিশেষ ভারপ্রাপ্ত কর্মকর্তা হিসেবে তাকে সরিয়ে নেওয়া হলেও বর্তমানে তিনি চট্টগ্রাম পোর্ট অথরিটির সদস্য হিসেবে দায়িত্ব পালন করছেন। নাহিদ ইসলামের এনআইডি লক করার বিষয়ে তিনি বাংলাদেশিউজকে বলেন, যখন কোনো অভিযোগ আসে তখন তদন্তকালীন এনআইডি লক করা হয়। এটা একটা প্রসিদ্ধি। যদি তদন্তে অভিযোগ প্রমাণিত না হয়, তখন আবার আনলক করা হয়। এটা প্রক্রিয়া। তো সেই ছাত্রলতার ক্ষেত্রেও এমনটি হয়েছিল। হয়তো কোনো এক জেসি থেকে অভিযোগ এসেছিল। তবে সেটা প্রমাণিত হয়নি। তাই কার্যকমিনের জন্য এনআইডি লক ছিল। ইসি কর্মকর্তারা জানিয়েছেন, সাধারণত কারও এনআইডি লক করা হলে কমিশনের সিদ্ধান্ত নেওয়া হয়। সরকার পরিবর্তনের ঠিক এক মাসের মাথায় গত বছরের ৫ সেপ্টেম্বর কাজী হাবিবুল আউয়ালের নেতৃত্বাধীন কমিশন পদত্যাগ করেন। এক্ষেত্রে কমিশন নূরু পরিচিহিত এনআইডি মহাপরিচালক ওই সিদ্ধান্ত নেন। বর্তমানে এনআইডি মহাপরিচালক এসএম ছমামুন কবীর বলেন, সে সময় আমি ছিলাম না। তাই সেটা আমার বিবেচনার বিষয় নয়। আর পুরোনো বিষয় যেটার সঙ্গে আমার কোনো সম্পৃক্ততা নেই, সেটার মতামতও দিতে চাই না। এনআইডি লক করা হলে সংশ্লিষ্ট ব্যক্তির আর নাগরিক সেবা মেলে না। সম্প্রতি শেখ হাসিনা ও পরিবারের ১০ সদস্য এবং সাবেক এনআইডি মহাপরিচালক সুলতানুজ্জামান মো. সালেহ উদ্দিনের এনআইডিও লক করেছে ইসি। অতীতেও অনেকের এনআইডি নানা প্রেক্ষাপটে লক করে সংস্থাটি। ফার্সিাদবিরাগী আন্দোলনের অবসান ঘটা সশস্ত্রকর্মের নিয়ে নতুন দল গঠনের জন্য গত ২৫ ফেব্রুয়ারি সরকারের উপদেষ্টা পরিষদ থেকে পদত্যাগ করেন নাহিদ ইসলাম। এরপর ২৮ ফেব্রুয়ারি এনসিপিরা আত্মপ্রকাশে আসে নতুন তাকে আধায়ক ঘোষণা করা হয়। ইউটিউ/এইচএ/ঢাকা: দেশের ছয়টি বিভাগের ওপর ছড়িয়ে পড়েছে দক্ষিণ-পশ্চিম নেদুশি বায়ু, অর্থাৎ বর্ষা। এর প্রভাবে এসব এলাকায় হাত পায় অতিভারী বৃষ্টি। শুক্রবার (৩০ মে) এমন পূর্বাভাস দিয়েছে আবহাওয়া অধিদপ্তর। আবহাওয়াবিদ ড. মুহাম্মদ আবুল কালাম মল্লিক জানান, দক্ষিণ-পশ্চিম নেদুশি বায়ু বরিশাল, চট্টগ্রাম, সিলেট, ময়মনসিংহ, ঢাকা ও রংপুর বিভাগে পড়তে অসমর্থ রয়েছে। শনিবার (৩০ মে) সন্ধ্যা পর্যন্ত ময়মনসিংহ, ঢাকা, খুলনা, বরিশাল, চট্টগ্রাম ও সিলেট বিভাগের অধিকাংশ জায়গায় এবং রংপুর ও রাজশাহী বিভাগের অনেক জায়গায় অস্বাভাবিকভাবে দমকা হাওয়াসহ বিদ্যুৎ চমকানো হালকা থেকে মাঝারি ধরনের বৃষ্টি/বজ্রসহ বৃষ্টি হতে পারে। একইভাবে ময়মনসিংহ, ঢাকা, খুলনা, বরিশাল, চট্টগ্রাম ও সিলেট বিভাগের কোথাও কোথাও ভারী থেকে অতিভারী বর্ষণ হতে পারে। সারাদেশে দিন ও রাতের তাপমাত্রা প্রায় অপরিবর্তিত থাকতে পারে। রোববার (১ জুন) সন্ধ্যা থেকে পরবর্তী ২৪ ঘণ্টার পূর্বাভাসে বলা হয়েছে, বরিশাল, চট্টগ্রাম ও সিলেট বিভাগের অধিকাংশ জায়গায়, ময়মনসিংহ ও ঢাকা বিভাগের অনেক জায়গায় এবং রংপুর, রাজশাহী ও খুলনা বিভাগের কিছু কিছু জায়গায় আ স্বাভাবিক দমকা হাওয়াসহ বিদ্যুৎ চমকানো হালকা থেকে মাঝারি ধরনের বৃষ্টি/বজ্রসহ বৃষ্টি হতে পারে। সর্বোচ্চ বরিশাল, চট্টগ্রাম ও সিলেট বিভাগে মাঝারি থেকে ভারী বর্ষণ হতে পারে। সারাদেশে দিনের তাপমাত্রা ১০-৩০ ডিগ্রি সেলসিয়াস

Figure 02: Raw Bengali Text Before Preprocessing

3.1 Clean the Text

The cleaning step removes unnecessary characters and standardizes the text. The following actions are performed:

- **Before Clean Text:**

Before clean text:

ঢাকা: জ্বালানি তেল বিক্রির কমিশন ১৫ কর্মদিবসে সমাধানের আশ্বাস দিয়েছে বাংলাদেশ পেট্রোলিয়াম কর্পোরেশন (বিপিসি)। ফল পেট্রোলপাম্প ও ট্যাংক লরি মালিক ঐক্য পরিষদ সারা দেশে তাদের কর্মবিরতি কর্মসূচি স্থগিত করেছে। রোববার (২৫ মে) দুপুরে রাজধানীর কারওয়ানবাজারে বিপিসির কার্যালয়ে বিপিসির চেয়ারম্যান আমিন উল আহসানের সঙ্গে আলোচনা করে পেট্রোলপাম্প ও ট্যাংক লরি মালিক ঐক্য পরিষদ ও সিদ্ধান্ত নেন। সাত দফা দাবিতে আজ সকাল ৬টা থেকে ধর্মঘট শুরু করে বাংলাদেশ পেট্রোলপাম্প ও ট্যাংক লরি মালিক ঐক্য পরিষদ। দুপুর ২টা পর্যন্ত পেট্রোলপাম্প বন্ধের পাশাপাশি জিন্দা থেকে তেল উত্তোলন ও পরিবহন বন্ধ রাখার ঘোষণা দিয়েছিল ঐক্য পরিষদ। এরপর ধর্মঘট প্রত্যাহার বারমাসীদের সঙ্গে বৈঠক করেন বাংলাদেশ পেট্রোলিয়াম কর্পোরেশন (বিপিসি)। কমিশন বৃদ্ধিসহ পেট্রোলপাম্প মালিকদের সাত দফার মধ্যে আছে, সড়ক অধিদপ্তরের ইজারা মাসুল আগের মতো বহাল করা, পাম্পের সংযোগ সড়কের ইজারা নবায়নের সময় পে-অর্ডারকে নবায়ন বাল প্রণয় করা, বিএসটিআই কর্তৃক আভ্যন্তরীণ ট্যাংক কালিকেশন, ডিপ রত পরীক্ষণ ফিস এবং নিবন্ধন প্রথা বাতিল করা, পেট্রোলপাম্পের ক্ষেত্রে পরিবেশ, বিইআরসি, কলকারখানা পরিদপ্তর, ফায়ার সার্ভিসের লাইসেন্স গ্রহণ প্রথা বাতিল করা, বিপসন কোম্পানি থেকে ভিলায়শিপ ছাড়া সরাসরি তেল বিক্রয় বন্ধ করা, ট্যাংক লরি চালকদের লাইসেন্স নবায়ন এবং নতুন লাইসেন্স বাধা বিপণিত ছাফা ইস্যু, সব ট্যাংক লরি জন্য আন্তঃজেলা রুট পারমিট ইস্যু করা, বিভিন্ন স্থানে অবনুমানিত এবং অবৈধভাবে ঘরের মধ্যে থোলা স্থানে যন্ত্রপাতি মেশিন নিয়ে জ্বালানি তেল বিক্রি বন্ধ করা। আরেকবার/৪৪টি

Figure 03: Bengali Text Before Cleaning Text

- **Remove non-Bengali characters:** Any character that is not a Bengali letter or whitespace is removed.

After removing non-Bengali characters:

ঢাকা বাংলাদেশ পেট্রোলপাম্প ও ট্যাংকলরি মালিক ঐক্য পরিষদের তাক ১০ দফা দাবিতে আজ রোববার ২৫ মে অর্ধদিবস কর্মবিরতি পালন চলাছে পূর্বাঘোষিত দাবি পূরণ না হওয়ায় রবিবার সকাল ৬টা থেকে এই কর্মবিরতি শুরু হ ছেছে এদিকে সকালে বাংলাদেশ পেট্রোলিয়াম কর্পোরেশন একটি বৈঠক ডেকেছে এই বৈঠকের ফলাফলের ওপর ভিত্তি করে তারা পরবর্তী সিদ্ধান্ত জানাবেন বলেও জানান তিনি ১১ মে এক সংবাদ সম্মেলনে পরিষদের পক্ষে জানানো হা ছিল ১৪ নের মধ্যে দাবি আদায় না হলে তারা ২৫ মে প্রতীকী কর্মসূচি পালন করা হবে পরিষদের ঘোষণা অনুযায়ী কর্মবিরতিতে জ্বালানি তেলের উত্তোলন পরিবহন ও বিপসন বন্ধ থাকবে তবে হজ ফ্লাইট ও আন্তর্জাতিক ফ্লাইট স চল রাখার জন্য উত্তোভাজকরা তেল পরিবহন চালু থাকবে শুধু আত্মসল ফায়ার সার্ভিস ও হেসে পেট্রোলপাম্পের সড়ক পুলিশের পাড়িতে জ্বালানি সরবরাহের চুক্তি আছে কেবল তাইই পুলিশের পাড়িতে জ্বালানি সরবরাহ করতে পারবে পরিষদের দাবিগুলোর মধ্যে রয়েছে পরিষদের দাবিগুলোর মধ্যে রয়েছে তেল বিক্রির কমিশন নূরুতথ্য ও শতাংশ করা সড়ক ও জনপথ অধিদপ্তরের ইজারা মাসুল আগের মতো বহাল রাখা পাম্প সংযোগ সড়কের ইজারা নবায়ন আবেদনবন্ধের সঙ্গে নির্ধারিত পে অর্ডার জমা দিলেই তা নবায়ন বিবেচিত করা বিএসটিআই শুধু ডিপার্মেন্ট ইউনিট ট্রাশিং ও পরিমাপ যাচাই করবে আরও রয়েছে আভ্যন্তরীণ ট্যাংক কালিকেশন ডিপ রত পরীক্ষণ ফিস ও নিবন্ধনপ্রথা বাতিল করা পরিবেশ কলকারখানা ও ফায়ার সার্ভিস লাইসেন্সের বিধান বাতিল করা ঘরের মধ্যে বা থোলা স্থানে অবৈধভাবে মেশিন বসিয়ে জ্বালানি বিক্রি বন্ধ করা এবং ভিলায়শিপ ছাড়া বিপসন কোম্পানির সরাসরি তেল বিক্রি বন্ধ করা ট্যাংকলরির চালকদের লাইসেন্স নবায়ন ও নতুন লাইসেন্স ইস্যু সহজ করা রাজায় যথোনে সেখানে ট্যাংকলরি থামিয়ে কাগজপত্র পরীক্ষা না করে তা তেলের জিন্দা গেটেই সম্পন্ন করা সব ট্যাংকলরির জন্য আন্তঃজেলা রুট পারমিট ইস্যু করা আরএইচ

Figure 04: Bengali Text After Remove Non-Bengali Characters

- [illegible]

- **Remove Bengali digits:** Bengali numerals (০ to ৯) are removed.

- After removing Bengali digits:
- | | | | |
|--|-----------------------------------|--|---|
| চাকা বাহোশাল পোতাশাল ও টাকালবির মানিক চক্রা পরিষদের ডাক ঘোষা দারিত আজ রোবাবর | যে অর্থমন্ত্রী কবিরিজি পালন চলাছে | পূর্ববাংলায় দাবি পূরণ না হয়য়া রবারির সকাল | ঐ থাকে এই কবিরিজি চক্রা হ |
| যেহে এফিক সচকা বাহোশাল পোতাশালয় ককরাপোনে এফিক বৈরিক কোকেছ | এই বৈরিকের লফাফোর এনর জিতিত | করা পরবর্তী শিয়ার জাফানাব পালন জিনি | যে ক সহকারে সফলতায় পরিষদের পক্ষ জানানো হ |
| হেফেল | যেহা মধ্যে দাবি আদায় না হলে ডাক | যে প্রকীরী অনুষ্ঠি পালন করা | যেহা পরিষদের যোগো অধ্যাপী |
| | | কবিরিজি জিয়ারা জেনের উত্থাপন | পরিষদে ও বিপনন পক্ষ |
| সফল রাখার জন্য উচ্চজাজেরাজে জেলে পরিষদে ডাক থাকবে | অনু আদায়ের ফায়ার সাডিস ও | যেহে পোতাশালপার সচকা | পুলিশের গাফিল জালানি |
| | | সবকারের জালানি | সবকারের জালানি |
| যে পরিষদের দাবিসেলার মধ্যে রয়েছে | পরিষদের পরিচালনার মধ্যে রয়েছে | যেহে বিক্রির কমিশন নুনাম | নতালন পক্ষ |
| | | সহকারে ডাক | সহকারে ও কলম্ব |
| | | অধিবাস্তুরের জালানা | জুমি ইকারা |
| জালানা বাহোশাল আদায়েরপার সচকা | নির্ধারিত যে অর্থ জালানি দা | নতায়ন জালানি | রিডেসটিআই |
| ফি ও নিকশপত্রী জালানি | সচকা পরিষদ | কলকাতাওয়া ও ফায়ার সাডিস | লাইসেনসের বিধান |
| সরাসরি জেলে বিক্রি বন্ধ করা | টাকালবির চালকদের | লাইসেনস নবায়ন ও | নতুন লাইসেনস |
| | | ইয়া সাহকার ডাক | রাষ্ট্রায় যোগো |
| কন্যা আজেকল পিটি পামিউ টি | মধ্য করা | আরউর | |

- **Remove short words:** Words with lengths between 1 to 3 characters are removed to reduce noise.

- [illegible]

- **Squishing whitespace:** Multiple spaces are replaced by a single space, and leading/trailing spaces are trimmed.

- After squishing whitespace:
- চাকা বোলানো পেছোঁপালো পুটানোর মালিক ঐক্য পরিষদের ডাক মারিত্ত রোবোর অর্থদেবস কমবিরিত্তি পালন চলছে পুর্ন্যুথিত্তি দাবি পুর্ন হুজুম রবিবার সকাল থেকে কমবিরিত্তি শুরু হয়েছে এদিকে সকাল বাৎসাব্দে পেছোঁপালিয় কর পোষোঁপন চাকি তৈরিক তৈরিক ফমফালন ভিত্তি তার পরবর্তী আন্তজ্ঞ জানাবন বালো জানান ভিত্তি সংবাদ সমবেদন পরিষদের পক্ষ জানাবনা রোয়ামি দায়িত্ব দাবি দায়িত্ব তার প্রতীকী কর্মবিত্তি পরিষদের ঘোষণা অত্যাচারী চাকি পরিষদের জ্ঞানালি থেকে উদ্ভাবন ফমফালন বিপদন বাক ধাবার মুখী আন্তজ্ঞ জ্ঞানালি রাখা চান উদ্ভাবনজ্ঞানোর বিপদন বাক ধাবার বাক আত্মলোক রোয়ামি দায়িত্ব থেকে পেছোঁপালিয়পালি সঙ্গ পালিশন পালিত্ত জ্ঞানালি সরকারের মুক্তি কলন বাক পুর্নালি পালিত্ত জ্ঞানালি সরকার বাক ধাবার পরিষদ মালিশ্রালার মাধ্যমে পরিষদের মালিশ্রালার মাধ্যমে রোয়ামি পরিষদের মালিশ্রালার মাধ্যমে রিক্রি কমিশন নুননয় জননয় জননয় অধিদপ্তর জ্ঞানার ভূমির ইগানো মালিশ্রালার বালন রাওয়াল মালিশ্রালার সফর জ্ঞানালি রোয়ামি আবেদনপত্রের সাজে নিবিত্তিত্তি রোয়ামি বিলেনে বিবিত্তিত্তি বিসম্মতীয় বাক জিসমপালি উদ্ভিত্তি জ্ঞানালি পরিষদ মাধ্যম জ্ঞানার রোয়ামি জ্ঞানার আভাব্যভাব্য উঠাং জ্ঞানালিগনন পলীকাল নিবন্ধনপত্র জ্ঞানালি পরিষদন কলারাবানো যোয়ামি সাজিত্তি লালেশ্রালার বিধান বাবিল ফালন মাধ্যম থোলা থোলা আবেদনপত্রের সঙ্গ পরিষদ বিবিত্তি বাক জিলারপালি জ্ঞানার জ্ঞানালি বাক জিলারপালি হুজা বিপদন কোম্পানির সারসরি বিবিত্তি বাক ট্যাকলরির চাকলরির লাইসেন্স নয়দায় নতুন লাইসেন্স কলার রাওয়াল ফোয়াম মেয়াল ট্যাকলরির জ্ঞানালি কাসপসন পলীকাল জেনের জিনো গেটেই স্পেন্সি চাকলরির জ্ঞান আন্তজ্ঞনা জ্ঞানালি হুজা আবেদন

Data Preview:	
url	cleaned_text
(character)	(character)
https://www.banglanews24.com/banglanews-special...	রাজধানীর প্রাণকেন্দ্রে নির্মল বাতাসে প্রাণভরে শ্বাস নেওয়ার জা...
https://www.banglanews24.com/politics/news/bd/15...	ঢাকা প্রধান উপদেষ্টা অধ্যাপক মুহাম্মদ ইউনুসের সাক্ষাৎকারের...
https://www.banglanews24.com/climate-nature/new...	ঢাকা হ্রল গভীর নিম্নচাপ ঘীরে ঘীরে দুর্বল হলেও প্রভাবে এখনো ...
https://www.banglanews24.com/national/news/bd/1...	প্রধান উপদেষ্টা অধ্যাপক মুহাম্মদ ইউনুস জাপানের প্রধানমন্ত্রী ...
https://www.banglanews24.com/environment-biodiv...	ঢাকা গভীর নিম্নচাপ হ্রলভাগে আসার দেশে মৌসুমের রেকর্ড বৃ...
https://www.banglanews24.com/saradesh/news/bd/...	নোয়াখালী বঙ্গোপসাগরে সৃষ্ট গভীর নিম্নচাপের প্রভাবে নোয়াখা...
https://www.banglanews24.com/cricket/news/bd/15...	জাতীয় ক্রীড়া পরিষদ এন.এস.সি কর্তৃক ফারুক আহমেদের পরি...
https://www.banglanews24.com/politics/news/bd/15...	ঢাকা সমগ্র জাতি সুশ্রু সুন্দর নিরপেক্ষ নির্বাচনের অপেক্ষমাণ ...
https://www.banglanews24.com/opinion/news/bd/15...	শহীদ প্রেসিডেন্ট জিয়াউর জাতীয় ইতিহাসের অবিস্মরণীয় জাতি...
https://www.banglanews24.com/saradesh/news/bd/...	ভোলা নিম্নচাপের প্রভাবে ভোলায় নদীর পানি অস্বাভাবিকভাবে ...
https://www.banglanews24.com/saradesh/news/bd/...	দিনাজপুর দিনাজপুরের পার্বতীপুরে ট্রাক আটকে চাঁদাবাজির তা...
https://www.banglanews24.com/daily-chittagong/ne...	চট্টগ্রাম জোড়া লাগানো শিশু রিয়াশাদ রেনিশকে জন্মের ঘণ্টার ...
https://www.banglanews24.com/politics/news/bd/15...	ঢাকা মহান স্বাধীনতার ঘোষক বিএনপির প্রতিষ্ঠিতা শহীদ প্রেসি...

This process helps reduce irrelevant content and prepares the text for deeper analysis.

3.2 Remove Stop Words

A custom list of commonly occurring but non-informative Bengali stop words is defined. These include words such as

```

bangla_stopwords <- c(
  "অবশ্য", "অবশ্যই", "অল্প", "অপর", "অথচ", "অথবা", "অনেক", "অন্য", "অন্যান্য", "অন্যান্যদের", "অনেকে", "অনেকেই",
  "অতএব", "অতি", "অতিরিক্ত", "আগে", "আগেও", "আগে থেকেই", "আগেই", "আগুন", "আছে", "আজ", "আবার", "আসলে", "ই",
  "ইতোমধ্যে", "উনি", "উপরে", "উপস্কা", "ও", "ওদিকে", "ওরা", "কত", "কখনও", "কখনো", "কখনোই",
  "কবলা", "কব", "কবেই", "কবি", "কম", "কম্পিউটার", "কয়েক", "কয়েকজন", "কয়েকটি", "কয়েকটা", "কয়েকজনই", "কয়েকজনের",
  "কয়েকদিন", "কয়েকটা", "কারণ", "কারণে", "কারো", "কারোই", "কারোকে", "কারোকেও", "কামীর", "কাজেই", "কাজেই",
  "কাজে", "কাজে লাগানো", "কাজের", "কাজের জন্য", "কাজের সময়", "কাজের কাজে", "কাজেই", "কাজে লাগানো", "কাজের জন্য",
  "কাজে লাগানো", "কাজেই", "কাজের জন্য", "কাজের কাজ", "কাজের কাজেই", "কাজের জন্য", "কাজের জন্যই", "কাজের জন্যে", "কাজের সময়",
  "কারও", "কারো", "কারোই", "কারোকে", "কারোকেও", "কারণ", "কারণে", "কখনও", "কখনো", "কখনোই", "কখনোও", "কখনোওই", "অবশেষে",
  "বাংলাদেশের", "নিয়", "পর", "বর্তমান", "করবে", "যে", "কোনো", "হয়", "দেশের", "আমাদের", "যে", "পারে", "মাঝে", "হিরে", "টাকা",
  "এর", "এবং", "না", "থেকে", "তিনি", "করা", "এ", "হয়", "তার", "বলেন", "কিন্তু", "তবে", "হয়েছে", "ছিল", "যা", "সে",
  "এই", "হয়", "জনা", "না", "ওরে", "তুমি", "আমি", "তারা", "আমরা", "আপনি", "তাহলে", "করেন", "হয়েছে", "ড"
)

```

Figure 10: Custom Bengali Stop Words List

among others. Removing these stop words from the cleaned text helps to focus the analysis on meaningful terms that contribute more significantly to the text's content.

[illegible]

Figure 11: Bengali Text After Removing Bengali Stop Words

3.3 Tokenization

Tokenization is the process of splitting each cleaned document into individual words, called tokens. This allows the text to be analyzed at the word level. For example, using R's `tokenize_words()` function, each article's text is broken down into a list of words. These tokens are the basic units used for tasks such as frequency analysis, creating document-term matrices, and topic modeling.

[1]	"রাজধানীর"	"প্রাণকেন্দ্রে"	"নির্মল"	"বাতাসে"	"প্রাণভরে"
[6]	"শ্বাস"	"নেওয়ার"	"জায়গার"	"অভাব"	"অভাবের"
[11]	"জায়গা"	"কিছুটা"	"হলও"	"পুরণ"	"করছিল"
[16]	"পানি"	"সবুজে"	"মেশা"	"হাতিরঝিল"	"রাস্তার"
[21]	"দুধারে"	"সবুজ"	"গাছের"	"ছায়ায়"	"দুদণ্ড"
[26]	"শান্তি"	"খুঁজতে"	"অনেকেই"	"ছুটে"	"আসেন"
[31]	"হাতির"	"বিলে"	"হুটরি"	"দিনে"	"ঘুরতে"
[36]	"আসেন"	"পরিবারসহ"	"প্রতিদিন"	"সন্ধ্যার"	"হাতিরঝিলের"
[41]	"কৃত্রিম"	"আলোর"	"ঝলকানিও"	"নগরের"	"মানুষের"
[46]	"কাছে"	"অন্যতম"	"আকর্ষণের"	"কেন্দ্র"	"এখানে"
[51]	"দর্শনার্থী"	"পথচারী"	"কিংবা"	"আশপাশের"	"বাসিন্দা"
[56]	"সবারই"	"চলাচল"	"চেপে"	"আয়তনের"	"হাতিরঝিল"
[61]	"এফডিসি"	"গুপ্তশান"	"পর্যন্ত"	"সংযুক্ত"	"জলাশয়ের"
[66]	"পানি"	"দুর্গন্ধমুক্ত"	"রাখতে"	"কর্তৃপক্ষের"	"উদ্যোগ"
[71]	"থাকলেও"	"দায়িত্ব"	"পাওয়া"	"ঠিকাদারি"	"প্রতিষ্ঠানের"
[76]	"কারণে"	"সেটি"	"দৃশ্যমান"	"বিলের"	"পানি"
[81]	"পরিষ্কার"	"রাখতে"	"বহুরের"	"ঠিকাদারি"	"প্রতিষ্ঠান"
[86]	"মেহজাবিন"	"এন্টারপ্রাইজকে"	"কোটি"	"বরাদ্দ"	"দেওয়া"
[91]	"অগ্রগতির"	"তেমন"	"প্রমাণ"	"মধুবাগ"	"এলাকার"
[96]	"স্থায়ী"	"বাসিন্দা"	"শাওন"	"আহমেদ"	"চলতি"
[101]	"হাতিরঝিলে"	"প্রাণভরে"	"শ্বাস"	"নেওয়ার"	"উপায়"
[106]	"সকালে"	"বাচ্চাদের"	"দিয়ে"	"হাটতে"	"হতাম"
[111]	"এমনই"	"অবস্থা"	"নিজের"	"ঠিকানা"	"এলাকা"
[116]	"ছেড়ে"	"যেতাম"	"প্রকট"	"দুর্গন্ধে"	"হাতিরঝিল"
[121]	"এলাকার"	"মানুষের"	"স্বাস্থ্যঝুঁকি"	"বাড়িছে"	"এখানে"
[126]	"বিশুদ্ধ"	"বাতাস"	"নির্মল"	"পরিবেশ"	"স্থানীয়"
[131]	"রফিকুল"	"নামে"	"আরেকজনের"	"হাতিরঝিল"	"প্রকল্পের"
[136]	"শুরুর"	"দিকে"	"ওয়াটার"	"ট্যান্ডিতে"	"গন্তব্যে"
[141]	"ফিরতাম"	"অন্যরকম"	"একটা"	"সতেজতা"	"অনুভব"
[146]	"চলাই"	"অস্বস্তিকর"	"উঠেছে"	"বাচ্চারা"	"সারা"
[151]	"লেকের"	"পাশে"	"বইসা"	"খেলত"	"গন্ধ"
[156]	"পানি"	"পোলাপান"	"আসতেই"	"হাতিরঝিল"	"এলাকার"
[161]	"বাসিন্দা"	"গৃহিণী"	"হাসনা"	"বিলের"	"ধারাই"
[166]	"ঝালমুড়ি"	"বিক্রি"	"হাসান"	"মিয়া"	"আগের"
[171]	"চেয়ে"	"মানুষ"	"হাতিরঝিলে"	"দূষণ"	"দুর্গন্ধে"
[176]	"এলাকাবাসীর"	"পাশাপাশি"	"আমরাও"	"অস্বস্তিতে"	"সালে"
[181]	"জনসাধারণের"	"উন্মুক্ত"	"হাতিরঝিল"	"প্রকল্প"	"পানি"
[186]	"পরিষ্কার"	"পরিশোধনের"	"পর্যাপ্ত"	"ব্যবস্থা"	"থাকায়"
[191]	"পানির"	"গুণগত"	"খারাপ"	"প্রকট"	"আকার"
[196]	"ধারণ"	"করেছে"	"পয়োবর্জ্য"	"ময়লা"	"আবর্জনা"
[201]	"ড্রেনের"	"পানি"	"ঢুকে"	"বিষাক্ত"	"উঠেছে"

[206]	"বিলের"	"পানি"	"বাতাসে"	"ভাসছে"	"উৎকট"
[211]	"গন্ধ"	"আশপাশ"	"ঘুরে"	"এমনটি"	"দেখছেন"
[216]	"প্রতিবেদক"	"হাতিরঝিল"	"ঘুরে"	"দেখা"	"গেছে"
[221]	"বিলের"	"প্রায়"	"থেকেই"	"পানির"	"দুর্গন্ধ"
[226]	"ভেসে"	"আসছে"	"পানিতে"	"নানা"	"বর্জ্য"
[231]	"পলিথিন"	"প্লাস্টিকের"	"পাইপও"	"ভাসতে"	"দেখা"
[236]	"গেছে"	"কারওয়ান"	"বাজারের"	"প্যান"	"প্যাসিফিক"
[241]	"সোনারগাঁও"	"হাটেলের"	"পেছনে"	"বর্জ্যের"	"মাত্রা"
[246]	"বেশি"	"কাজী"	"নজরুল"	"ইসলাম"	"অ্যাভিনিউ"
[251]	"থেকেও"	"দুর্গন্ধ"	"পাওয়া"	"বিলে"	"মাছও"
[256]	"ভাসতে"	"দেখা"	"পানির"	"ওপরে"	"বর্জ্যের"
[261]	"পুরু"	"স্তর"	"তৈরি"	"ওয়াটার"	"ট্যান্ডিতে"
[266]	"চলার"	"সময়"	"স্পষ্টভাবে"	"দেখা"	"যাচ্ছিল"
[271]	"হাতির"	"বিলের"	"পানি"	"জীববৈচিত্র্য"	"পরিবেশের"
[276]	"হুমকি"	"বলছেন"	"পরিবেশবিদরা"	"রাজউকের"	"কর্মকর্তারা"
[281]	"বলছেন"	"প্রায়"	"নর্দমার"	"পানি"	"হাতিরঝিলের"
[286]	"পানিতে"	"মিশে"	"নিষেধাজ্ঞা"	"থাকলেও"	"মানুষ"
[291]	"আবর্জনা"	"পলিথিন"	"এমনকি"	"শিল্প"	"বর্জ্য"
[296]	"ফেলছে"	"বিলের"	"পানিতে"	"পানি"	"দূষিত"
[301]	"হচ্ছে"	"ময়লা"	"আবর্জনা"	"গ্যাস"	"গ্যাস"
[306]	"সৃষ্টি"	"হচ্ছে"	"কারণেই"	"দুর্গন্ধ"	"ছিড়িয়ে"
[311]	"পড়ছে"	"রাজউকের"	"তত্ত্বাবধায়ক"	"প্রকৌশলী"	"যান্ত্রিক"
[316]	"সাবির"	"তাহের"	"বাংলানিউজকে"	"হাতিরঝিলের"	"চারপাশে"
[321]	"প্রায়"	"নর্দমা"	"গর্ত"	"রয়েছে"	"প্রতি"
[326]	"মাসে"	"একবার"	"সেগুলো"	"পরিষ্কার"	"কঠিন"
[331]	"বর্জ্য"	"গর্তের"	"পেছনে"	"অরপরে"	"লেকের"
[336]	"পানিতে"	"ভেসে"	"রামপুরা"	"কাঁঠালবাগানে"	"দুটি"
[341]	"সলুইস"	"রয়েছে"	"বর্ষাকালে"	"এগুলো"	"খুলতে"
[346]	"কঠিন"	"বর্জ্য"	"প্রবেশ"	"প্রতি"	"রক্ষণাবেক্ষণের"
[351]	"বরাদ্দের"	"অংশই"	"বিদ্যুৎ"	"খাতে"	"যায়"
[356]	"পানি"	"পরিশোধনের"	"বরাদ্দ"	"রাসায়নিক"	"প্রায়"
[361]	"হয়ে"	"সালের"	"আগের"	"অবশিষ্ট"	"রাসায়নিক"
[366]	"অল্প"	"পরিমাণে"	"ব্যবহার"	"হচ্ছে"	"বছরে"
[371]	"মাত্রা"	"একবার"	"সালে"	"হাতিরঝিল"	"সংগৃহীত"
[376]	"নমুনায়"	"ক্ষতিকর"	"রাসায়নিক"	"পিএফওএ"	"পারফ্লুরোঅকটোনোয়িক"
[381]	"অ্যাসিড"	"পিএফওএস"	"পারফ্লুরোঅকটোনোসালফোনিক"	"অ্যাসিড"	"উভয়ই"
[386]	"দীর্ঘমেয়াদি"	"বিশাক্ততার"	"দায়ী"	"পিএফওএসের"	"স্তর"
[391]	"পরামর্শমূলক"	"স্তরের"	"চেয়ে"	"বেশি"	"মাত্রায়"
[396]	"পাওয়া"	"গেছে"	"বিভিন্ন"	"গবেষণায়"	"হাতিরঝিলে"
[401]	"পানি"	"দূষণের"	"বিষয়টি"	"এসেছে"	"এরপরও"
[406]	"পানি"	"দূষণ"	"যাচ্ছে"	"বাংলানিউজের"	"বলেছে"

[411]	"রাজউক"	"কর্তৃপক্ষ"	"লেকের"	"নির্বাহী"	"তত্ত্বাবধায়ক"
[416]	"প্রকৌশলী"	"যান্ত্রিক"	"সাবির"	"তাহের"	"জানান"
[421]	"হাতিরঝিল"	"লেক"	"ময়লা"	"এরইমধ্যে"	"তিনটি"
[426]	"বিলে"	"মেহজাবিন"	"এন্টারপ্রাইজকে"	"পরিশোধ"	"পরিশোধ"
[431]	"হলেও"	"ঠিকাদার"	"কোম্পানি"	"মেহজাবিনের"	"দেখা"
[436]	"যায়নি"	"সরেজমিনে"	"হাতিরঝিলের"	"পানি"	"পরিষ্কার"
[441]	"পরিচ্ছন্ন"	"রাখতে"	"বসানো"	"কয়েকটি"	"স্পেশাল"
[446]	"সুয়ারেজ"	"ডাইভারশন"	"স্ট্রাকচার"	"এসএসডিএস"	"স্ক্রিনিং"
[451]	"মেশিন"	"মেশিন"	"অপারেটর"	"হিসেবে"	"থাকার"
[456]	"করছেন"	"ছয়জন"	"তাদের"	"বেতন"	"দেওয়া"
[461]	"হাজার"	"এসএসডিএস"	"হোটেল"	"সোনারগাঁও"	"এলাকায়"
[466]	"জাহিদ"	"হাসান"	"গোলাপ"	"এসএসডিএস"	"মগবাজার"
[471]	"সংলগ্ন"	"এলাকায়"	"সেলিম"	"এসএসডিএস"	"মধুবাগ"
[476]	"এলাকায়"	"মামুন"	"এসএসডিএস"	"নিকেতন"	"এলাকায়"
[481]	"ওহিদ"	"শহিদ"	"ছাড়া"	"এসএসডিএস"	"ম্যানুয়াল"
[486]	"সাতটি"	"মেশিন"	"এগুলোর"	"কাজের"	"চারজন"
[491]	"শ্রমিক"	"থাকার"	"যারা"	"ম্যানুয়ালি"	"করবেন"
[496]	"সরেজমিনে"	"একজনকেও"	"পাওয়া"	"যায়নি"	"এমনকি"
[501]	"অন্য"	"অপারেটরদের"	"কাছে"	"জানতে"	"চাইলে"
[506]	"জানান"	"ম্যানুয়াল"	"মেশিনে"	"কর্মচারীই"	"এসএসডিএসের"
[511]	"বর্জ্য"	"সরানোর"	"শ্রমিক"	"গাড়ি"	"লেবার"
[516]	"সপ্তাহে"	"একদিন"	"কারওয়ান"	"বাজার"	"গাড়ি"
[521]	"ভাড়া"	"তিনজন"	"শ্রমিক"	"ভাড়া"	"কাজটি"
[526]	"জানিয়েছেন"	"প্রকাশে"	"অনিচ্ছুক"	"কর্মচারী"	"শুরুতে"
[531]	"জলাশয়"	"পরিষ্কার"	"করার"	"চারজন"	"ডুবুরি"
[536]	"ছিলেন"	"যারা"	"বিষয়"	"জানতে"	"চাইলে"
[541]	"রাজউক"	"প্রকৌশলী"	"সাবির"	"তাহের"	"তেমন"
[546]	"কিছু"	"জানাতে"	"পারেননি"	"এমনকি"	"সরেজমিনে"
[551]	"পাওয়া"	"কর্মচারীর"	"সংখ্যার"	"দেওয়া"	"সংখ্যার"
[556]	"তথ্যেরও"	"অমিল"	"পাওয়া"	"ঠিকাদারি"	"প্রতিষ্ঠান"
[561]	"মেহজাবিন"	"এন্টারপ্রাইজের"	"মালিক"	"জাকির"	"হোসেনের"
[566]	"নিজের"	"অধীনে"	"থাকা"	"কর্মীর"	"দেখে"
[571]	"নামই"	"বলতে"	"পারেননি"	"পর্যায়"	"সুপারভাইজারকে"
[576]	"দিয়ে"	"জানতে"	"কাজের"	"চিত্র"	"সম্পর্কে"
[581]	"কাছে"	"জানতে"	"চাইলেও"	"সদুত্তর"	"দিতে"
[586]	"পারেননি"	"থেকেই"	"স্পষ্ট"	"নিলেও"	"লেকের"
[591]	"ময়লা"	"পরিষ্কারে"	"অগ্রগতি"	"বরাদ্দের"	"কোথায়"
[596]	"জানতে"	"চাইলে"	"প্রকাশে"	"অনিচ্ছুক"	"ঠিকাদার"
[601]	"কোম্পানির"	"কর্মচারী"	"বিভিন্ন"	"নামে"	"তোলা"
[606]	"হলেও"	"নিজের"	"সহকর্মীদের"	"এখনো"	"দেখেননি"
[611]	"অর্থাৎ"	"বিভিন্ন"	"কর্মচারীর"	"দিয়ে"	"তৈরি"
[616]	"তোলা"	"ঠিকই"	"নিয়োগই"	"রাজউক"	"চেয়ার"
[621]	"ম্যানু"	"রিয়াজুল"	"ইসলাম"	"বাংলানিউজকে"	"জানান"
[626]	"দুর্নীতি"	"অনিয়ম"	"তদন্তে"	"কমিটি"	"তদন্ত"
[631]	"প্রতিবেদনের"	"প্রতিবেদন"	"অনুযায়ী"	"দোষীদের"	"বিরুদ্ধে"
[636]	"ব্যবস্থা"	"নেওয়া"	"দৃষ্টিতে"	"ক্ষতি"	"হচ্ছে"
[641]	"পরিবেশবিদ"	"আহমদ"	"কামরুজ্জামান"	"মজুমদার"	"এলিভেটেড"
[646]	"এক্সপ্রেসওয়ার"	"পিলারই"	"সোনারগাঁও"	"হোটেলের"	"পেছনে"
[651]	"হাতিরঝিল"	"সংলগ্ন"	"এলাকা"	"সুয়ারেজ"	"সিস্টেম"
[656]	"ব্যাহত"	"হচ্ছে"	"মাটি"	"দিয়ে"	"ভরাতের"
[661]	"কারগেও"	"পরিবেশের"	"মারাত্মক"	"ক্ষতি"	"হচ্ছে"
[666]	"পানিতে"	"অক্সিজেনের"	"পরিমাণ"	"গিয়ে"	"ভেসে"
[671]	"উঠছে"	"ছাড়া"	"বিলের"	"সামগ্রিক"	"জীববৈচিত্র্যে"
[676]	"মারাত্মক"	"প্রভাব"	"পড়ছে"	"শহরের"	"প্রতিবেশ"
[681]	"ব্যবস্থায়"	"বিরূপ"	"প্রভাব"	"ফেলছে"	"কামরুজ্জামান"
[686]	"বিলে"	"আগের"	"ব্যাঙ"	"বিলীন"	"যাচ্ছে"
[691]	"উপকারী"	"ক্ষুদ্র"	"পোকামাকড়ও"	"নগরের"	"মানুষের"
[696]	"স্বাস্থ্য"	"ঝুঁকি"	"বাড়াচ্ছে"	"পানি"	"পরিষ্কার"
[701]	"পরিচ্ছন্ন"	"রাখতে"	"ঘটুকু"	"বরাদ্দ"	"রাখা"
[706]	"সেটুকুও"	"ঠিকঠাক"	"কার্যকর"	"হাতিরঝিলের"	"পানির"
[711]	"দুর্দশা"	"নগরের"	"মানুষও"	"স্বস্তি"	"দুর্নীতিরও"
[716]	"কারপে"	"পরিবেশে"	"বিরূপ"	"প্রভাব"	"পড়ছে"
[721]	"রাজধানীর"	"মানুষের"	"স্বাস্থ্যকে"	"ঝুঁকিতে"	"ফেলছে"
[726]	"কামরুজ্জামান"	"আরএইচ"			

Figure 12: Tokenized Bengali Words

4.Exploratory Text Analysis

A Term Document Matrix (TDM) was created to analyze term frequencies across documents. Word frequencies were calculated and the most common words identified. Visualizations such as word clouds and bar plots of the top 20 words were generated to aid intuitive exploration of the data.

Data Preview:	
word (character) ▾	freq (double) ▾
জিয়াউর	189
হিসেবে	159
ঢাকা	153
ঘোষণা	141
তাদের	134
জাতীয়	131
জিয়া	123
সালের	122
দিয়ে	121
সরকার	117
শহীদ	116
শুরু	115
বিভিন্ন	111
সালে	111
করেছেন	111
প্রধান	111
করার	109
ছিলেন	109
রহমানের	109
হচ্ছে	107
পানি	101
করেছে	96
মাধ্যমে	93
রয়েছে	92
কিছু	91
সরকারের	91
নির্বাচন	90
চট্টগ্রাম	90
স্বাধীনতার	89
রাজনৈতিক	88
কারণে	87
পর্যন্ত	86
বিএনপির	85
বৃহস্পতিবার	81
প্রথম	80
দলের	79
বিএনপি	79
মেজর	79
বৃষ্টি	78
গেছে	77
জানান	76
দেখা	76
নতুন	76
মার্চ	73
সাবেক	71
তাকে	70
কাছে	69
দেওয়া	69
দাবি	69
অনেক	66

Figure 13: CSV Preview of Bengali Words by Frequency

4.1 Word Cloud with Most Frequent Word



Figure 14: Word Cloud of Most Frequent Bengali Words

This word cloud provides a visual summary of word frequencies found in a collection of Bengali text documents, most likely related to news or political topics.

1. Most Frequent Words:

- জিয়াউর (Ziaur): This is the largest word, showing it is the most commonly mentioned term in the text. It probably refers to “Ziaur Rahman” or the Zia political faction, indicating a political focus.
- সরকার (Government): Indicates discussions involving government policies, activities, or official statements.
- ঢাকা (Dhaka): The capital of Bangladesh, frequently mentioned in political news and national events.
- ঘোষণা (Announcement): Points to frequent mentions of official declarations or statements.

2. Political Context:

- Words such as নির্বাচন (Election), বিএনপি (BNP - Bangladesh Nationalist Party), সদস্য (Member), দলের (Party's) suggest that the text mainly covers political events, party actions, and elections.
- Terms like স্বাধীনতা (Independence), শ্রদ্ধা (Respect/Honor), শহীদ (Martyr) reflect themes of nationalism and historical significance within political narratives.

3. Geographical and Social References:

- The prominence of ঢাকা (Dhaka) highlights the central role of the capital in the news or discussion topics.
- Other terms like পানি (Water) and বৃষ্টি (Rain) hint at social or environmental issues alongside political reporting.

4. Other Frequent Terms:

- Names and titles like প্রধান (Chief/Head), মেয়র (Mayor), আহমেদ (Ahmed) suggest mentions of specific leaders or officials.
- Words such as কারণ (Reason), বিরুদ্ধে (Against), প্রথম (First) indicate the presence of narrative or argumentative elements typical in news stories.

This word cloud reflects a text dataset rich in political and governmental themes, focusing on prominent figures likely connected to the Zia family or faction, the government, and national matters centered mainly on Dhaka. It also includes related topics such as elections, official announcements, social/environmental concerns, and expressions of national pride and historical remembrance.

4.2 Bar Chart of Top 20 Words

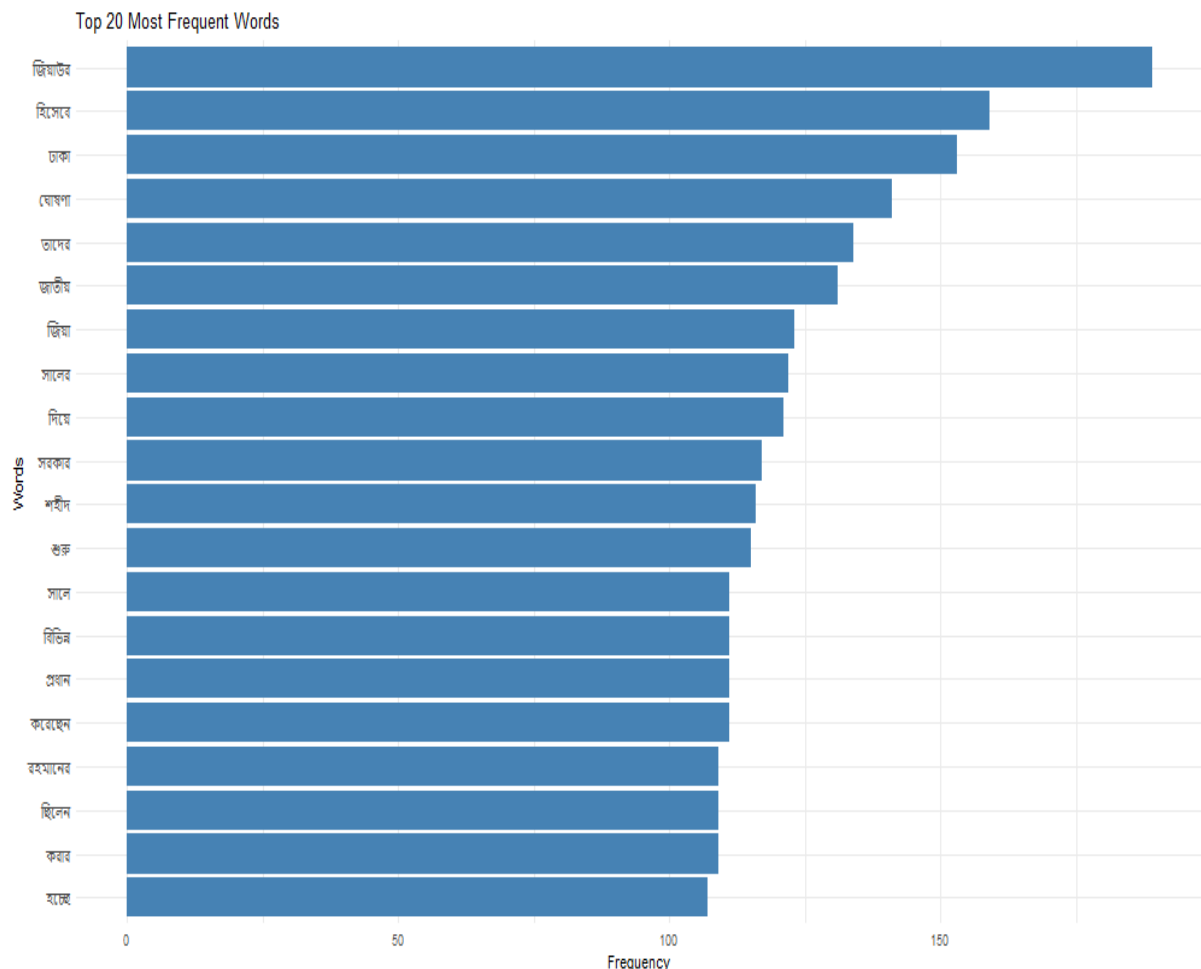


Figure 15: Bar Chart of Top 20 Most Frequent Bengali Words

This bar chart presents the top 20 most frequent Bengali words from a text corpus, likely centered on political and social topics. The word “জিয়াউর” (referring to Ziaur Rahman or his political faction) appears most frequently, followed by terms like “তারা” (they), “যেখানে” (where), and “ঢাকা” (the capital city), indicating a strong focus on political figures, places, and government-related discussions. Words such as “সরকার” (government), “নির্বাচন” (election), and “বিএনপি” (BNP) further highlight themes of political events, party activities, and electoral processes. Additionally, nationalistic terms like “স্বাধীনতা” (independence) and “শ্রদ্ধা” (respect) appear alongside social and environmental references such as “পানি” (water) and “বৃষ্টি” (rain), reflecting a broader context of news coverage. The presence of official titles and narrative words like “প্রধান” (chief), “মেয়র” (mayor), “কারণ” (reason), and “বিরুদ্ধে” (against) suggests detailed reporting and argumentation within the texts. Overall, the chart offers a clear quantitative visualization that complements the word cloud, emphasizing the political and social nature of the dataset with a strong geographic focus on Dhaka.

5. Topic Modeling

Topic modeling is a technique used to automatically discover the main themes or topics present in a large collection of text documents. It helps to organize and summarize the text by grouping words that frequently appear together into meaningful topics.

5.1 Document-Term Matrix (DTM)

This step involves creating a structured matrix where rows represent documents and columns represent terms (words). Each cell in the matrix contains the frequency of a particular term in a specific document. This matrix serves as the foundational input for topic modeling algorithms, allowing them to analyze the distribution of words across documents systematically.

5.2 Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) is a probabilistic model used to uncover hidden thematic structures in a large collection of text documents. By applying LDA to the Document-Term Matrix, the algorithm identifies 3 to 5 main topics. Each topic is represented as a distribution over words that frequently occur together, revealing underlying themes or subjects within the dataset.

5.3 Top Words per Topic

After the topics are identified, bar plots are generated to display the top significant words for each topic. These visualizations make it easier to interpret the content of each topic by highlighting the most influential words, aiding in understanding the essence of the themes discovered.

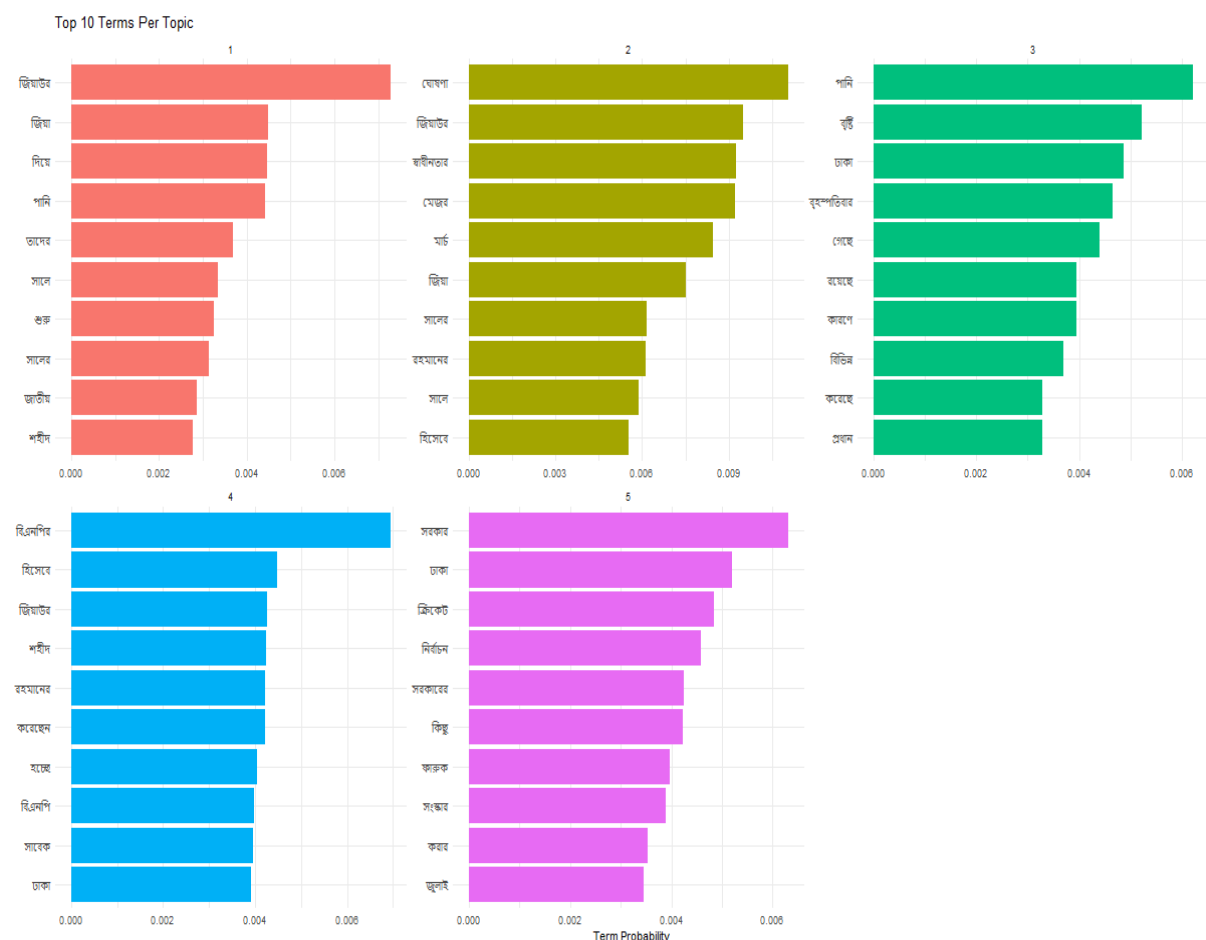


Figure 16: Bar Charts of Top Words per Topic

Topic 1 (Red Bar) Word Analysis:

- জিয়াউর, জিয়া: These words refer to prominent political figures in Bangladesh, specifically related to Ziaur Rahman and his political legacy, indicating a political theme.
- দিয়ে: A verb indicating an action, likely used in narratives describing events or processes.
- পানি: Refers to water, suggesting environmental or natural resource topics.
- তাদের, সালে, শুরু: Words related to time and groups of people, indicating the start or timeline of certain events.
- সালেন, জাতীয়, শহীদ: Words that relate to national events or commemorations, such as martyrdom or national days.

Topic 2 (Olive Bar) Word Analysis:

- ঘোষণা: Means "announcement," indicating official statements or declarations.
- জিয়াউর, জিয়া, রহমানের: Names related to political leaders, linking this topic to political history or current politics.
- স্বাধীনতার: Refers to "independence," indicating discussions about national freedom or historical struggles.
- মেজর, মার্চ, সালের: Military rank and time indicators, likely referencing specific dates or periods.
- হিসেবে: Means "as" or "in the capacity of," showing roles or positions.

Topic 3 (Green Bar) Word Analysis:

- শনি, বৃষ্টি: Related to weather, specifically rain and days of the week (Saturday).
- ঢাকা: The capital city of Bangladesh, a central location for political, social, and environmental news.
- বৃহস্পতিবার, গেছেকারণে: Time references and causal words, indicating when and why something happened.
- বিভিন্ন, করেছে, প্রধান: Describing various events or actions with importance.

Topic 4 (Blue Bar) Word Analysis:

- বিএনপি: Refers to the Bangladesh Nationalist Party, a major political party.
- হিসেবে: Denotes roles or perspectives.
- জিয়াউর, রহমানের: Political leader names, indicating political discussions.
- শহীদ: Refers to martyrs, highlighting historical or commemorative themes.
- করেছেন, হচ্ছে: Action verbs indicating ongoing or completed activities.
- সারেক, ঢাকা: Names of places or organizations, indicating local context.

Topic 5 (Purple Bar) Word Analysis:

- সরকার, সরকারের: Refers to government and related administrative matters.
- ঢাকা: Political and administrative hub.
- ক্রিকেট: Refers to the sport, indicating social or entertainment-related topics.
- নির্বাচন: Refers to elections, pointing to democratic processes and political events.
- কিছু, ফারুক, সংস্কার, করার, জুলাই: Various terms relating to actions, individuals, reforms, and time (July).

Each topic's key words reflect important aspects of Bangladesh's political landscape, social environment, and daily life. Political figures and parties such as জিয়াউর and বিএনপি, as well as terms like সরকার and নির্বাচন, highlight the political discourse. Environmental and social themes emerge from words like পানি, বৃষ্টি, and ঢাকা, while entertainment and cultural life appear through words like ক্রিকেট. Time-related words such as মার্চ, শনি, and জুলাই emphasize temporal contexts within the data. Overall, the dataset provides a rich representation of the multifaceted realities of Bangladesh's society, politics, and environment.

5.4 Document-Topic Distribution

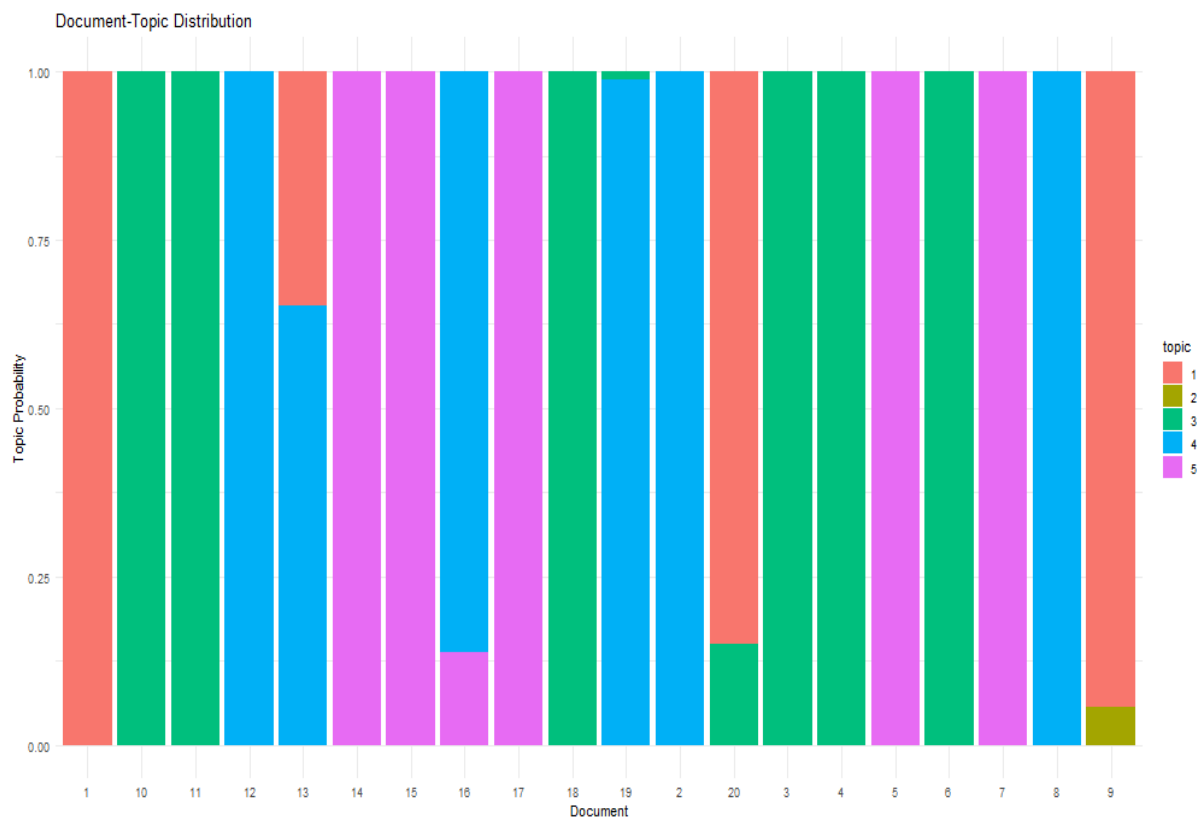


Figure 17: Document–Topic Distribution

To understand how topics are spread across individual documents, charts like stacked bar plots or pie charts are created. These visualizations show the proportion or probability of each topic within each document, illustrating the thematic composition of the documents and highlighting dominant topics.