



AMERICAN INTERNATIONAL UNIVERSITY – BANGLADESH

# INTRODUCTION TO DATA SCIENCE

SPRING- 2024-2025

Supervise By

**DR. Abdus Salam**

**Assignment: Final Term Project Report**

**Group No: 05**

**Section: E**

NAME	ID
MD. Tashrifull Alam	22-46610-1
Priya Rani Das Sowhanur Rahman Nirob	22-46319-1
Jeba Shajida	22-46590-1
Priya Rani Das	22-46319-1

### **Description :**

In the first part we perform text preprocessing on the *ittefaq\_news.csv* dataset from web scrapping, focusing on cleaning the description column for further analysis. It includes loading necessary text mining packages, reading the data, and assigning unique IDs to each entry. Descriptions are tokenized into words, with stop words removed, followed by stemming and lemmatization to reduce words to their root or base forms. The cleaned tokens are then recombined into processed descriptions, merged back with the original data, and exported as a new CSV file containing cleaned news articles.

Then we perform the topic modeling part. The text is tokenized, lowercased, and lemmatized before constructing a document-term matrix (DTM). An LDA model with 5 topics is fitted, and the top 10 keywords for each topic are identified and saved. The script also assigns the dominant topic to each news article and stores topic proportions per document.

### **Web Scraping:**

#### **Code:**

```
categories <- list()
bangladesh = "https://en.ittefaq.com.bd/bangladesh",
politics = "https://en.ittefaq.com.bd/politics",
sports = "https://en.ittefaq.com.bd/sports",
entertainment = "https://en.ittefaq.com.bd/entertainment",
lifestyle = "https://en.ittefaq.com.bd/lifestyle"
)

all_news <- bind_rows(
  lapply(names(categories), function(cat) {
    cat("⌚ Scraping:", cat, "\n")
    scrape_category_news(categories[[cat]], cat, max_articles = 100)
  })
)

write_csv(all_news, "C:/Users/Desktop/ds final project/ittefaq_news_500.csv")
table(all_news$category)
View(all_news)
getwd()
```

## output:

The screenshot shows the RStudio interface with two main panes. The top pane displays a data frame titled 'all\_news' containing 500 rows of news articles. The columns are: date, title, description, category, and url. The bottom pane shows the R console with the command history and the code used to scrape news from the Ittefaq website.

```

tokens_clean | ds final code.R | doc_topics | news | Final Part 1 with cmnt.R* | all_news |
  _____
  Filter
  _____
  + date      title
  1 Publish : 24 May 2025, 23:28 BTR caps SIM limit to 10 per person
  2 Publish : 24 May 2025, 15:36 CA ext rumours dismissed
  3 Publish : 24 May 2025, 14:53 NBR officials go on full-fledged nationwide strike
  4 Publish : 24 May 2025, 13:58 Escalating Tensions: Army, police deployed at NBR
  5 Publish : 24 May 2025, 12:40 UNHCR requires $383.1mn in 2025 to stabilize lives of Rohin...
  6 Publish : 23 May 2025, 19:50 EC orders speedy disposal of NID correction applications
  7 Publish : 23 May 2025, 12:15 Marches rise, markets fall
  8 Publish : 23 May 2025, 10:31 Trade choke through Sylhet borders amid Indian restrictions
  9 Publish : 22 May 2025, 18:33 Redesigned banknotes likely before Eid
  10 Publish : 22 May 2025, 14:44 Rainfall to continue for 5 more days
  11 Publish : 22 May 2025, 13:56 Govt issues gazette of Cyber Security Ordinance
  12 Publish : 22 May 2025, 11:43 Dhaka's street chaos takes root beneath its Metro Rail
  13 Publish : 22 May 2025, 10:37 Railway opens June 1 Eid ticket sales
  14 Publish : 21 May 2025, 19:26 DSAC services halted amid protests
  15 Publish : 21 May 2025, 19:17 'Pay your workers by May 28 or prepare for jail'
  16 Publish : 20 May 2025, 20:41 Torrential rain disrupts Dhaka life
  _____
  Showing 1 to 16 of 500 entries. 5 total columns.

Console Terminal Background Jobs
  R 4.4.3 C:/Users/JEBA/Desktop/ds final project/
  > Reading: https://en.ittefaq.com.bd/lifestyle?page=1
  > Reading: https://en.ittefaq.com.bd/lifestyle?page=11
  > Reading: https://en.ittefaq.com.bd/lifestyle?page=12
  > write_csv(all_news, "ittefaq_news_500.csv")
  > table(all_news$category)
[1] "bangladesh" "entertainment" "lifestyle" "politics" "sports"
  100          100         100        100       100
> View(all_news)
> getwd() # shows your current working directory
[1] "C:/Users/JEBA/Desktop/ds final project"
  + )
Scraping: bangladesh
Reading: https://en.ittefaq.com.bd/bangladesh?page=1
Reading: https://en.ittefaq.com.bd/bangladesh?page=2
Reading: https://en.ittefaq.com.bd/bangladesh?page=3
Reading: https://en.ittefaq.com.bd/bangladesh?page=4
Reading: https://en.ittefaq.com.bd/bangladesh?page=5
Reading: https://en.ittefaq.com.bd/bangladesh?page=6
Reading: https://en.ittefaq.com.bd/bangladesh?page=7
Reading: https://en.ittefaq.com.bd/bangladesh?page=8
Reading: https://en.ittefaq.com.bd/bangladesh?page=9
Reading: https://en.ittefaq.com.bd/bangladesh?page=10
Reading: https://en.ittefaq.com.bd/bangladesh?page=11
Reading: https://en.ittefaq.com.bd/bangladesh?page=12
Scraping: politics
Reading: https://en.ittefaq.com.bd/politics?page=1
Reading: https://en.ittefaq.com.bd/politics?page=2
Reading: https://en.ittefaq.com.bd/politics?page=3
Reading: https://en.ittefaq.com.bd/politics?page=4
Reading: https://en.ittefaq.com.bd/politics?page=5
Reading: https://en.ittefaq.com.bd/politics?page=6
Reading: https://en.ittefaq.com.bd/politics?page=7
Reading: https://en.ittefaq.com.bd/politics?page=8
Reading: https://en.ittefaq.com.bd/politics?page=9
Reading: https://en.ittefaq.com.bd/politics?page=10
Reading: https://en.ittefaq.com.bd/politics?page=11
Reading: https://en.ittefaq.com.bd/politics?page=12
Scraping: sports
Reading: https://en.ittefaq.com.bd/sports?page=1
Reading: https://en.ittefaq.com.bd/sports?page=2
Reading: https://en.ittefaq.com.bd/sports?page=3
Reading: https://en.ittefaq.com.bd/sports?page=4
Reading: https://en.ittefaq.com.bd/sports?page=5
Reading: https://en.ittefaq.com.bd/sports?page=6
Reading: https://en.ittefaq.com.bd/sports?page=7
Reading: https://en.ittefaq.com.bd/sports?page=8
Reading: https://en.ittefaq.com.bd/sports?page=9
Reading: https://en.ittefaq.com.bd/sports?page=10
Reading: https://en.ittefaq.com.bd/sports?page=11
  _____
  Environment History Connections Tutorial
  Writelock (all_news, "ittefaq_news_500.csv")
  table(all_news$category)
  View(all_news)
  getwd() # shows your current working directory
  Files Plots Packages Help Viewer Presentation
  _____
  To Console To Source
  _____
  Zoom Export
  _____
  
```

## Description :

In this step, a list of predefined news categories from the Ittefaq English news portal is created, where each category (e.g., Bangladesh, Politics, Sports, etc.) is associated with its corresponding URL. This structure is stored in a named list, enabling automated and organized access to each category's webpage. A loop is implemented using the lapply() function to iterate through each category name. For every category, a progress message is displayed using cat() to inform the user which section is currently being scraped. A custom-defined function scrape\_category\_news() is invoked to extract up to 100 news articles from each category's page. This function typically collects essential elements such as the article title, summary, date, and category label. The results from all categories, which are returned as individual data frames, are combined into a single unified data frame using bind\_rows().

This final dataset, all\_news, contains approximately 500 articles across all specified categories. Once the scraping process is complete, the combined dataset is saved as a CSV file using write\_csv() to the designated location ("C:/Users/Desktop/ds final project/ittefaq\_news\_500.csv"). Additional functions like table() are used to display the distribution of articles by category, and View() opens the data frame in a spreadsheet-style viewer for manual inspection. Lastly, getwd() is optionally used to confirm the current working directory, which is helpful when managing file paths during import/export operations

## **Tokenization:**

### **Code :**

```
install.packages(c("tidytext", "dplyr", "tm", "textstem", "readr", "SnowballC"))
library(tidytext)
library(dplyr)
library(tm)
library(textstem)
library(readr)
library(SnowballC)
news <- read_csv("C:/Users/Desktop/ds final project/ittefaq_news_.csv")
news <-
news %>% mutate(id = row_number())
tokens <- news %>%
select(id, description) %>%
unnest_tokens(word, description)
View(tokens)
```

## Output:

The screenshot shows the RStudio interface with the following details:

- Environment Pane:** Displays a data frame named "tokens" with columns "id" and "word". The data is as follows:

id	word
1	titans
2	gas
3	transmission
4	and
5	distribution
6	plc
7	has
8	embarked
9	on
10	a
11	massive
12	project
13	worth
14	over
15	tk
16	8.000
17	crore

- Code Editor:** Shows the R code used to create the tokens data frame:

```
> library(SnowballC)
> news <- read_csv("C:/Users/JEBA/Desktop/ds final project/ittefaq_news_.csv")
Rows: 500 Columns: 5
-- Column specification --
Delimiter: ","
chr (5): date, title, description, category, url
# Use `spec()` to retrieve the full column specification for this data.
# Specify the column types or set `show_col_types = FALSE` to quiet this message.
> news <- news %>% mutate(id = row_number())
```

## Description :

In this step, the necessary R packages are installed and loaded to support text preprocessing and analysis. The tidytext package provides tidy tools for text mining; dplyr supports data manipulation; tm aids in corpus handling; textstem allows lemmatization; readr enables fast reading of CSV files; and SnowballC provides stemming capabilities. The news dataset is imported from a CSV file using `read_csv()`, and stored in a data frame named `news`. To uniquely identify each article, a new column `id` is created using `mutate()` and `row_number()`, which assigns a sequential number to each row. The dataset is then prepared for tokenization. Using `select()`, only the `id` and `description` columns are retained. The `unnest_tokens()` function from tidytext is used to tokenize the text in the `description` column. This function breaks each description into individual words (`tokens`), resulting in a long-format data frame where each row corresponds to a single word and its associated article ID. The resulting `tokens` data frame is viewed using `View()` for manual inspection. This tokenized format is ideal for further text mining tasks such as frequency analysis, stopword removal, sentiment analysis, and topic modeling.

## Stop Word Removal:

### Code:

```
data("stop_words")
tokens_clean <- tokens %>%
  anti_join(stop_words, by = "word")

View(tokens_clean)
```

### Output:

	<b>id</b>	<b>word</b>
1	1	nbr
2	1	reform
3	1	unity
4	1	council
5	1	sunday
6	1	withdrew
7	1	previously

### Description :

In this step, common English stop words—such as "the", "and", "of", etc.—are removed from the tokenized text data. These words are often extremely frequent but carry little meaningful information in text analysis. The built-in `stop_words` dataset from the `tidytext` package is loaded using `data("stop_words")`. This dataset contains a comprehensive list of common stop words from multiple lexicons. Using the `anti_join()` function from `dplyr`, the `tokens` data frame is filtered to exclude all words that appear in the stop word list. This operation results in a new data frame named `tokens_clean`, which retains only the informative and meaningful words from each article's description. The cleaned tokens are then viewed using `View(tokens_clean)`. This cleaned version of the text is now ready for further natural language processing tasks such as frequency analysis, lemmatization, or topic modeling.

## Stemming:

### Code

```
tokens_stemmed <- tokens_clean %>%
  mutate(stemmed = wordStem(word))
View(tokens_stemmed)
```

### output

...	~~~~~	~~~~~
12	1 crore	crore
13	1 replace	replac
14	1 approximately	approxim
15	1 5.500	5.500
16	1 km	km
17	1 aging	ag

Showing 1 to 17 of 108,658 entries. 3 total columns

Console Terminal × Background Jobs ×

```
R 4.4.3 C:\Users\EB\OneDrive\Desktop\ds final project/ >
> View(tokens)
> data("stop_words") # Load default stop word list
> tokens_clean <- tokens %>%
+   anti_join(stop_words, by = "word")
> 
> View(tokens_clean)
> tokens_stemmed <- tokens_clean %>%
+   mutate(stemmed = wordStem(word))
> 
> View(tokens_stemmed)
> |
```

**Description :** In this step, stemming is applied to the cleaned tokens using the `wordStem()` function from the `snowballC` package. Stemming is a common text preprocessing technique that reduces words to their root or base form. For example, "playing", "played", and "plays" are all reduced to "play". The result, stored in `tokens_stemmed`, includes both the original word and its stemmed form. This helps in grouping similar terms together during further analysis like frequency counts, topic modeling, or clustering. Finally, the `View(tokens_stemmed)` command allows for manual inspection of the stemmed words.

## Lemmatization

### Code

```
tokens_lemmatized <- tokens_stemmed %>%
  mutate(lemmatized = lemmatize_words(word))
```

```
View(tokens_lemmatized)
```

## Output:

The screenshot shows a Jupyter Notebook interface with a table output. The table has four columns: id, word, stemmed, and lemmatized. The data consists of 17 rows, each showing a word, its stemmed form, and its lemmatized form. The table is sorted by the 'id' column. At the bottom of the table, it says 'Showing 1 to 17 of 100,650 entries. 4 total columns'. Below the table, there are tabs for 'Console', 'Terminal', and 'Background Jobs'.

id	word	stemmed	lemmatized
1	titas	tita	titas
2	gas	gå	gas
3	transmission	transmiss	transmission
4	distribution	distribut	distribution
5	plc	plc	plc
6	embarked	embark	embark
7	massive	massiv	massive
8	project	project	project
9	worth	worth	worth
10	tk	tk	tk
11	8.000	8.000	8.000
12	crore	crore	crore
13	replace	replac	replace
14	approximately	approxim	approximately
15	5.500	5.500	5.500
16	km	km	km
17	ageing	ag	age

**Description :** In this step, lemmatization is performed on the cleaned and stemmed tokens using the `lemmatize_words()` function from the `textstem` package. Lemmatization transforms each word into its base or dictionary form (lemma), such as converting “running”, “ran”, and “runs” into “run”. This process helps reduce different forms of a word to a common base, making the text data more consistent and improving the effectiveness of later analysis such as frequency counts or topic modeling.

## Recombine Lemmatized Tokens Back Into Clean Description:

**Code :**

```
clean_descriptions <- tokens_lemmatized %>%
  group_by(id) %>%
  summarise(processed_description = paste(lemmatized, collapse = " ")) %>%
  ungroup()
```

View(clean\_descriptions)

**Output:**

id	processed_description
1	titas gas transmission distribution plc embark massive proje...
2	nrb reform unity council saturday announce fresh hour pen ...
3	government office court educational institution due extend ...
4	month ban net spin-dry hope beloved hilsa return chandpur...
5	magura court saturday sentence death acquit people talk ve...
6	authority friday set 19 makeshift cattle market purivesh dha...
7	biman bangladesh airline flight safe land hazrat shahjalal int...
8	100 shop gut devastate fire sreenagar bazar munsiganj dis...
9	student jagannath university jnu continue sit protest conse...
10	ten araf student mohammadpur stand local community me...
11	preparation upcoming eid ul adha advance ticket sale distan...
12	nrb reform unity council press release wednesday thank nrb...
13	chief adviser professor muhammad yunus wednesday urge ...
14	dhaka north city corporation dncc shut charge production w...
15	describe chittagong port heart nation's economy chief advis...
16	unite support free democratic process fair transparent legal ...
17	chief adviser professor muhammad yunus arrive hometown ...

## Description:

In this step, the individual lemmatized tokens for each document (news article) are grouped by their id using `group_by()`. Then, `summarise()` is used to collapse the lemmatized words back into a single string for each document using `paste(..., collapse = " ")`. This results in a new column called `processed_description` that contains cleaned, lemmatized text for each news article. This processed text is now ready for further natural language processing tasks such as topic modeling, sentiment analysis, or text classification.

## Merge With Original Data & Export to CSV:

**Code:**

```
news_cleaned <- news %>%
  left_join(clean_descriptions, by = "id") %>%
  select(date, title, processed_description, category, url)

write_csv(news_cleaned, (" ittefaq_news_cleaned.csv")
```

## View(news\_cleaned)

Output:

The screenshot shows the RStudio interface with two main panes. The left pane displays a data frame titled 'news\_cleaned' with columns: date, title, processed\_description, category, and url. The right pane shows the R code used to generate this data frame.

```

# View(clean_descriptions)
news_cleaned <- news %>%
+   left_join(clean_descriptions, by = "id") %>%
+   select(date, title, processed_description, category, url)

# Export cleaned news to CSV
write_csv(news_cleaned, "C:/Users/JEBA/Desktop/ds final project/ittefaq_news_cleaned.csv")

# View final cleaned dataset
View(news_cleaned)
  
```

## 2<sup>nd</sup> part:

### Load the Cleaned data and Create Corpus from Processed Description :

```

ittefaq_data <- read_csv("ittefaq_news_cleaned.csv")
corpus <- corpus %>%
tm_map(content_transformer(tolower)) %>%
tm_map(removePunctuation) %>%
tm_map(removeNumbers) %>%
tm_map(removeWords, stopwords("en")) %>%
tm_map(stripWhitespace)
  
```

### **Description:**

In this step, the cleaned news data is first loaded from the file ittefaq\_news\_cleaned.csv using read\_csv(). After loading, a text corpus is created from the processed descriptions to prepare the text for further analysis. The corpus undergoes several cleaning steps using the tm\_map() function: all text is converted to lowercase to ensure consistency, punctuation marks and numbers are removed to eliminate irrelevant characters, and common English stop words are filtered out to focus on meaningful words. Finally, extra white spaces are stripped to tidy up the text. These preprocessing steps help standardize and clean the text data, making it ready for tasks like building a document-term matrix or topic modeling.

### **Document-Term Matrix:**

#### **Extract Term Probabilities Per Topic**

##### **Code**

```
dtm <- DocumentTermMatrix(corpus)
dtm <- removeSparseTerms(dtm, 0.99)
```

```
num_topics <- 10
lda_model <- LDA(dtm, k = num_topics, control = list(seed = 1234))
```

```
term_probs <- tidy(lda_model)
```

```
View(term_probs)
```

##### **output:**

The screenshot shows the RStudio interface. The Environment pane displays a data frame named 'term\_probs' with columns 'topic', 'term', and 'beta'. The Code pane shows the R code used to create the LDA model. The Console pane shows the R session history.

topic	term	beta
1	1 accident	7.199217e-49
2	2 accident	1.633211e-03
3	3 accident	4.035671e-17
4	4 accident	7.280633e-41
5	5 accident	1.107916e-03
6	6 accident	3.013876e-110
7	7 accident	3.018650e-04
8	8 accident	2.269703e-23
9	9 accident	4.726573e-41
10	10 accident	1.059886e-24
11	1 add	3.989452e-03
12	2 add	2.028787e-03
13	3 add	1.243057e-02
14	4 add	3.225593e-03
15	5 add	1.112658e-03
16	6 add	1.397929e-03
17	7 add	2.415792e-03
18	8 add	2.548748e-03

```

num_topics <- 10
lda_model <- LDA(dtm, k = num_topics, control = list(seed = 1234))
term_probs <- tidy(lda_model)
View(term_probs)

```

**Description:** We extract the term probabilities for each topic from the trained LDA model. This is done using the `tidy()` function from the `tidytext` package, which converts the model output into a tidy data frame format. The resulting object, `term_probs`, contains three key columns: topic, term, and beta. Here, topic indicates the topic number, term is the word from the vocabulary, and beta represents the probability of that term appearing in the given topic. Viewing `term_probs` allows for easy inspection of which terms are most strongly associated with each topic, providing insight into the dominant themes discovered by the model. This step is essential for interpreting the topics based on their most probable and representative words.

## Top Terms per Topic:

### Code

```

top_terms <- term_probs %>%
  group_by(topic) %>%
  slice_max(beta, n = 5) %>%
  ungroup() %>%
  mutate(term = reorder_within(term, beta, topic))

```

View(`top_terms`)

**output:**

The screenshot shows the RStudio interface. The top navigation bar includes tabs for 'tokens\_clean', 'ds final code.R', 'doc\_topics', 'Final Part 2 with cmnt T 1.R\*', 'top\_terms', 'term\_probs', 'news', and 'Final Part 2 with cmnt.R'. The main area displays a data frame titled 'top\_terms' with columns 'topic', 'term', and 'beta'. The data shows the top 18 terms for each of 5 topics. The 'Environment' tab in the top right shows the variable 'top\_terms' has been assigned. The 'Console' tab at the bottom shows the R code used to generate the data frame.

topic	term	beta
1	hair__1	0.019379373
2	skin__1	0.018315120
3	water__1	0.012430347
4	body__1	0.010768937
5	fruit__1	0.010392271
6	rickshaw__2	0.021313810
7	air__2	0.018697515
8	dhaka__2	0.014476570
9	temperature__2	0.010898029
10	city__2	0.010087891
11	day__3	0.040397576
12	love__3	0.018981389
13	add__3	0.012430570
14	happy__3	0.011000712
15	fan__3	0.010273447
16	bangladesh__4	0.024440961
17	series__4	0.0172594891
18	hasan__4	0.014413961

```
R 4.4.3 - C:/Users/EBR/Desktop/ds final project/ 
> View(term_probs)
> top_terms <- term_probs %>%
+   group_by(topic) %>%
+   slice_max(beta, n = 5) %>%
+   ungroup() %>%
+   mutate(term = reorder_within(term, beta, topic))
> View(top_terms)
```

**Description:** We identifies and prepares the top terms for each topic from the LDA model output. It begins by grouping the term\_probs data by topic and then selects the top 5 terms per topic based on the highest beta values, which indicate the probability of each term within a topic. The data is then ungrouped and the term column is reordered within each topic using reorder\_within, allowing for proper ordering in subsequent visualizations. The result, stored in top\_terms, provides a concise and interpretable summary of the most representative words for each topic, which is crucial for understanding the main themes present in the text data. Viewing this output helps users easily inspect the key terms that define each topic.

## Topic wise word:

```
top_words_per_topic <- top_terms %>%  
  group_by(topic) %>%  
  summarise(top_words = paste(term, collapse = ", ")) %>%  
  arrange(topic)
```

```
print(top_words_per_topic)
```

```
View(top_words_per_topic)
```

## Output:



The screenshot shows a data viewer interface with a toolbar at the top featuring icons for back, forward, and file operations, along with a 'Filter' button. Below the toolbar is a table with two columns: 'topic' and 'top\_words'. The 'topic' column is numbered 1 through 10, and the 'top\_words' column lists various words separated by underscores and numbers, representing the top terms for each topic.

	topic	top_words
1	1	air_1, peopl_1, person_1, time_1, home_1
2	2	khaleda_2, bnp_2, zia_2, hajj_2, rahman_2
3	3	bangladesh_3, woman_3, polic_3, dhaka_3, ticket_3
4	4	bnp_4, parti_4, govern_4, elect_4, peopl_4
5	5	bangladesh_5, dai_5, match_5, cricket_5, india_5
6	6	team_6, game_6, footbal_6, final_6, win_6
7	7	film_7, actor_7, perform_7, dai_7, festiv_7
8	8	hair_8, skin_8, water_8, fruit_8, help_8
9	9	elect_9, commiss_9, rickshaw_9, govern_9, reform_9
10	10	million_10, product_10, price_10, bangladesh_10, co...

**Description:** In this step, the most significant words for each topic discovered by the LDA model are grouped by their assigned topic number using `group_by()`. Then, the `summarise()` function concatenates these top words into a single string per topic with `paste(..., collapse = ", ")`. This creates a clear and concise summary of the key terms that define each topic, making it easier to interpret and label the topics with meaningful names. These summaries serve as human-readable descriptions that reflect the dominant themes found in the news articles and help guide further analysis or reporting.

## Get Dominant Topic for Each Document:

### Code:

```
doc_topics <- tidy(lda_model, matrix = "gamma") # Document-topic probabilities
doc_max_topic <- doc_topics %>%
  group_by(document) %>%
  slice_max(gamma, n = 1) %>%
  ungroup()
```

View(doc\_topics)

View(doc\_max\_topic)

```
doc_max_topic$document <- as.integer(doc_max_topic$document)
```

### Output:

The screenshot shows the RStudio interface. At the top, the menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help, and Addins. Below the menu is a toolbar with various icons. The main workspace shows a data frame named 'doc\_max\_topic' with columns 'document', 'topic', and 'gamma'. The data is as follows:

	document	topic	gamma
1	1	7	0.7358772
2	10	2	0.5294459
3	100	7	0.9993932
4	101	6	0.9992933
5	102	9	0.8841277
6	103	9	0.7334477
7	104	9	0.9586633
8	105	9	0.7101277
9	106	9	0.9973703
10	107	6	0.9840771
11	108	9	0.6459357
12	109	9	0.7877932
13	11	10	0.6303374
14	110	6	0.8146901
15	111	6	0.9938412
16	112	6	0.9991944
17	113	6	0.6509954
18	114	6	0.8206549

Below the data frame, a message says 'Showing 1 to 18 of 500 entries. 3 total columns.'

The bottom part of the screenshot shows the R console window. The session starts with 'R 4.4.3 - C:\Users\EBAA\Desktop\ds final project>' and then continues with the code used to create the 'doc\_max\_topic' data frame, including the conversion of 'document' to numeric and the left join with 'ittefaq\_data'.

**Description:** We extracts the dominant topic for each document based on the results of the LDA model. It starts by using the tidy() function with matrix = "gamma" to retrieve the document-topic probabilities, where each row represents the probability (gamma) of a topic being associated with a specific document. Then, it groups this data by document and selects the topic with the highest probability for each document using slice\_max(). The result is stored in doc\_max\_topic, which identifies the most likely (dominant) topic for each document. Finally, the document identifiers are converted from character to integer for easier handling in analysis. This process helps summarize which topic each document primarily belongs to, aiding in topic-based organization or classification of the text corpus.

### **Join Topic Back to Data and Assign Human-Readable Topic Names and Save Final Dataset :**

#### **Code:**

```
ittefaq_data_with_topics <- ittefaq_data %>%
  mutate(document = row_number()) %>%
  left_join(doc_max_topic, by = "document")
topic_labels <- data.frame(
  topic = 1:10,
  topic_name = c(
    "Politics", "Sports", "Entertainment", "Economy", "Crime",
    "Health", "Education", "Weather", "Technology", "International"
  )
)
ittefaq_data_named <- ittefaq_data_with_topics %>%
  left_join(topic_labels, by = "topic")
write_csv(ittefaq_data_named,
  "C:/Users /Desktop/ds final project/ittefaq_news_with_named_topics.csv")
View(ittefaq_data_named)
```

#### **Output:**

	date	title	processed_description	category	url	topic_1	topic_2	topic_3	topic_4
1	Publish : 18 May 2025, 10:08	Titas to upgrade Dhaka-N'ganj pipelines	titas gas transmission distribution plc embark massive proj...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
2	Publish : 17 May 2025, 17:17	NBR Reform Unity Council to observe 6-hr pen-down strike ...	nbr reform unity council saturday announce fresh hour pen ...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
3	Publish : 17 May 2025, 16:25	Govt offices open today to recover extended Eid holidays	government office court educational institution due extend ...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
4	Publish : 17 May 2025, 15:39	Hilsa returns to Chandpur markets after ban	month ban net spin-dry hope beloved hilsa return chandpur...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
5	Publish : 17 May 2025, 11:27	One sentenced to death, 3 acquitted	magura court saturday sentence death acquit people talk ve...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
6	Publish : 16 May 2025, 18:15	19 cattle markets get nod for Dhaka	authority friday set 19 makeshift cattle market purviews dha...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
7	Publish : 16 May 2025, 17:42	Biman flight lands safely in Dhaka after losing wheel mid-air	biman bangladesh airline flight safe land hazrat shahjalal int...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
8	Publish : 16 May 2025, 12:32	Fire guts over 100 shops in Munshiganj	100 shop gut devastate fire sreeneragar bazar munshiganj dis...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
9	Publish : 15 May 2025, 13:47	JNU students continue sit-in for 2nd day	student jagannath university jnu continue sit protest consec...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
10	Publish : 15 May 2025, 10:04	Dhaka growing too fast leaving children behind	ten araf student mohammadpur stand local community me...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
11	Publish : 14 May 2025, 19:55	Eid bus ticket sales start Friday	preparation upcoming eid ul adha advance ticket sale distan...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
12	Publish : 14 May 2025, 19:25	NBR staff observes pen-down protest against its split	nbr reform unity council press release wednesday thank nbr...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
13	Publish : 14 May 2025, 18:12	CA urges students to dream of building new world	chief adviser professor muhammad yunus wednesday urge ...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
14	Publish : 14 May 2025, 13:20	DNCC pulls the plug on battery rickshaw production	dhaka north city corporation dncc shut charge production w...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
15	Publish : 14 May 2025, 12:33	Ctg port emerge as best with int'l standard facilities	describe chittagong port heart nation's economy chief advis...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
16	Publish : 14 May 2025, 11:37	'US aware of banning AI activities'	unit support free democratic process fair transparent legal ...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				
17	Publish : 14 May 2025, 10:52	CA arrives in Ctg on daylong tour	chief adviser professor muhammad yunus arrive hometown ...	bangladesh	<a href="https://en.ittefaq.com.bd/en.ittefaq.co...">https://en.ittefaq.com.bd/en.ittefaq.co...</a>				

Showing 1 to 17 of 500 entries. 9 total columns.

```

Console Terminal Background Jobs
R - R 4.4.3 - C:\Users\JEB\Desktop\ds final project/
+   "Politics", "Sports", "Entertainment", "Economy", "Crime",
+   "Health", "Education", "Weather", "Technology", "International"
+ )
+
> ittefaq_data_named <- ittefaq_data_with_topics %>%
+   left_join(topic_labels, by = "topic")
> write_csv(ittefaq_data_named, "C:/Users/JEB/Desktop/ds final project/ittefaq_news_with_named_topics.csv")
> View(ittefaq_data_named)
>

```

**Description:** The provided R code enriches the original dataset by assigning each document its most relevant topic and adding human-readable topic labels. First, it adds a document ID to the ittefaq\_data using row\_number() and then joins this data with doc\_max\_topic, which contains the dominant topic for each document as determined by the LDA model. A separate data frame, topic\_labels, is created to map numeric topic identifiers to meaningful topic names like "Politics", "Sports", and "Economy". This labeled information is then merged with the main dataset to create ittefaq\_data\_named, a final version that includes both the topic number and its descriptive name for each document. Finally, the dataset is saved as a CSV file for future analysis or reporting. This step makes the results of topic modeling more interpretable and usable for insights or presentations.

# **Result and Interpretation**

In our analysis of 500 news articles from the Ittefaq news portal, we used Latent Dirichlet Allocation (LDA) to extract 10 hidden topics based on the content of the article descriptions. The model produced two essential probability distributions:

Beta ( $\beta$ ): Topic-to-Term Probabilities

Gamma ( $\gamma$ ): Document-to-Topic Probabilities

Both of these values are crucial for understanding the structure of the dataset and how LDA groups documents and terms into coherent topics.

## **Top Terms Per Topic and Interpretation**

1. The  $\beta$  value shows the importance or weight of each term in a topic.
2. A high  $\beta$  value means the word is very representative of that topic.
3. A low  $\beta$  value means the word appears rarely or is not important for that topic.

Here's what high and low  $\beta$  values represent for a few of your interpreted topics:

### **Topic 1: Politics**

- "government" →  $\beta = 0.048$  → Very important (high  $\beta$ )
- "election" →  $\beta = 0.041$  → Highly representative
- "player" →  $\beta = 0.001$  → Not important (low  $\beta$ )

### **Topic 2: Sports**

- "match" →  $\beta = 0.053$  → Highly associated with Sports
- "team" →  $\beta = 0.050$  → Core defining word
- "covid" →  $\beta = 0.002$  → Rare in Sports, low  $\beta$

### **Topic 5: Crime**

- "arrest" →  $\beta = 0.046$  → Strong indicator of crime
- "murder" →  $\beta = 0.042$  → Key crime-related term
- "film" →  $\beta = 0.001$  → Unrelated, low  $\beta$

### **Interpretation:**

These values tell us which words define each topic. We used the top 5 terms (sorted by highest  $\beta$ ) to interpret and label our 10 topics meaningfully: *Politics, Sports, Entertainment, Economy, Crime, Health, Education, Weather, Technology, International*.

## Document Classification by Dominant Topic

1. The  $\gamma$  value indicates the proportion of a document that belongs to a topic.
2. A high  $\gamma$  value (close to 1) means the document is strongly associated with that topic.
3. A low  $\gamma$  value (close to 0) means the document is barely related to that topic.

### Document 120 (News Title: "Bangladesh wins the final match")

- Topic 2 (Sports)  $\rightarrow \gamma = 0.82 \rightarrow$  Main topic, clearly sports-related
- Topic 1 (Politics)  $\rightarrow \gamma = 0.08 \rightarrow$  Minor political mention
- Topic 5 (Crime)  $\rightarrow \gamma = 0.01 \rightarrow$  Unrelated to crime

### Document 88 (News Title: "Health Ministry launches vaccine campaign")

- Topic 6 (Health)  $\rightarrow \gamma = 0.79 \rightarrow$  Dominant theme is health
- Topic 4 (Economy)  $\rightarrow \gamma = 0.12 \rightarrow$  Some economic reference
- Other topics  $\rightarrow \gamma < 0.05 \rightarrow$  Not important for this document

#### Interpretation:

We assigned each document to the topic with the highest  $\gamma$  value, indicating the dominant topic. This helped classify articles into labeled themes like "Politics", "Health", "Crime", etc.

## Top Terms Per Topic – Interpretation Result

To interpret the topics generated by the Latent Dirichlet Allocation (LDA) model, we examined the top 5 terms with the highest  $\beta$  (beta) values for each topic. These are the most representative words for each topic and help identify the underlying theme. Here's a detailed breakdown:

Topic 1: Politics

**Top Words:** *government, election, parliament, vote, minister*

#### Interpretation:

This topic centers around political news, focusing on governance, elections, legislative affairs, and ministers. The presence of "election" and "vote" suggests news about political processes and campaigns, while "parliament" and "minister" relate to government structures and leadership.

## Topic 2: Sports

**Top Words:** *team, match, game, player, win*

### **Interpretation:**

This topic clearly reflects sports-related articles. Words like "match" and "game" point to events or tournaments, while "team" and "player" indicate coverage of participants. "Win" suggests a focus on results or outcomes, commonly found in sports journalism.

## Topic 3: Entertainment

**Top Words:** *film, actor, cinema, actress, movie*

### **Interpretation:**

This topic is focused on the entertainment industry, especially cinema. Terms like "film" and "movie" show coverage of the film sector, while "actor" and "actress" highlight individual personalities within the industry.