**AMERICAN INTERNATIONAL UNIVERSITY–BANGLADESH**

**FACULTY OF COMPUTER SCIENCE & ENGINEERING**

**Course Name: INTRODUCTION TO DATA SCIENCE**

**Semester: Spring 2024-2025**

**Section: E      Group No: 11**

**Supervised By: DR. ABDUS SALAM**

**Final TERM PROJECT**

| Name | ID |
|---|---|
| MD YEASIN NEWAZ | 22-46803-1 |
| MD. RAKIBUL ISLAM | 22-47102-1 |
| MD MOSTAFIJUR RAHAMAN BIPUL | 22-46762-1 |
| MAHEDI HASAN | 22-46259-1 |

# Part-1

## Scraping NPR Newspaper by Category

## Code



```r
safe_GET <- function(url, retries = 3, delay_sec = 2) {
  headers <- add_headers(`User-Agent` = "Mozilla/5.0")
  for (i in seq_len(retries)) {
    res <- try(GET(url, headers, timeout(10)), silent = TRUE)
    if (!inherits(res, "try-error") && status_code(res) == 200) {
      return(res)
    } else {
      message(paste("Request failed:", url, "- Retry", i))
      Sys.sleep(delay_sec)
    }
  }
  return(NULL)
}
get_full_article_text <- function(article_url) {
  res <- safe_GET(article_url)
  if (is.null(res)) {
    message("❌ Failed to fetch full article:", article_url)
    return(NA_character_)
  }
  page <- read_html(res)
  selectors <- c(
    "div.article-body p",
    "div#storytext p",
    "div[data-testid='story-text'] p",
    "div.selectorgadget_suggested p"
  )
  for (sel in selectors) {
    paragraphs <- page %>% html_nodes(sel) %>% html_text(trim = TRUE)
    if (length(paragraphs) > 0 && any(nzchar(paragraphs))) {
      return(paste(paragraphs, collapse = "\n\n"))
    }
  }
  return(NA_character_)
}
get_npr_articles <- function(category, max_articles = 100) {
  base_url <- paste0("https://www.npr.org/sections/", category, "/")
  articles <- list()
  page <- 1
  while (length(articles) < max_articles) {
    url <- paste0(base_url, "?page=", page)
    cat("📄 Scraping", category, "- Page", page, "\n")
    res <- safe_GET(url)
    if (is.null(res)) break
    webpage <- read_html(res)
    article_nodes <- html_nodes(
      webpage,
      ".item-info, .bucketwrap, .internallink, .insettwocolumn, .inset2col, .bucketwrap.image.large"
    )
    if (length(article_nodes) == 0) break
    for (node in article_nodes) {
      title <- node %>% html_node("h2.title, h3.title") %>% html_text(trim = TRUE)
      link <- node %>% html_node("a") %>% html_attr("href")
      date_str <- node %>%
        html_node("time") %>%
        { if (!is.null(.)) html_attr(., "datetime") %||% html_text(., trim = TRUE) else NA_character_ }
      parsed_date <- parse_date_time(date_str, orders = c("Ymd HMS", "Ymd", "mdY", "B d, Y"), quiet = TRUE)
      if (any(is.na(c(title, link, parsed_date))) || grepl("/podcasts/", link)) {
        next
      }
      full_text <- tryCatch({
        get_full_article_text(link)
      }, error = function(e) NA_character_)
      if (is.na(full_text) || str_trim(full_text) == "") next
      articles[[length(articles) + 1]] <- tibble(
        title = title,
        description = full_text,
        date = parsed_date,
        category = category,
        link = link
      )
      if (length(articles) >= max_articles) break
      Sys.sleep(1)
    }
    page <- page + 1
    Sys.sleep(2)
  }
  bind_rows(articles)
}
categories <- c("business", "health", "politics", "technology", "world")
all_articles <- bind_rows(lapply(categories, get_npr_articles, max_articles = 100))
print(table(all_articles$category))
write_csv(all_articles, "npr.csv")
message("✅ Scraping complete. File saved: npr_csv")
```

# Output

```
Scraping business - Page 1
Scraping business - Page 2
Scraping business - Page 3
Scraping business - Page 4
Scraping business - Page 5
Scraping business - Page 6
Scraping business - Page 7
Scraping business - Page 8
Scraping business - Page 9
Scraping business - Page 10
Scraping business - Page 11
Scraping business - Page 12
Scraping business - Page 13
Scraping business - Page 14
Scraping business - Page 15
Scraping business - Page 16
Scraping business - Page 17
Scraping business - Page 18
Scraping business - Page 19
Scraping business - Page 20
Scraping health - Page 1
Scraping health - Page 2
Scraping health - Page 3
Scraping health - Page 4
Scraping health - Page 5
Scraping health - Page 6
Scraping health - Page 7
Scraping health - Page 8
Scraping health - Page 9
Scraping health - Page 10
Scraping health - Page 11
Scraping health - Page 12
Scraping health - Page 13
Scraping health - Page 14
Scraping health - Page 15
Scraping health - Page 16
Scraping health - Page 17
Scraping health - Page 18
Scraping health - Page 19
Scraping health - Page 20
```

```
> print(table(all_articles$category))

  business    health  politics technology     world
       100       100       100       100       100
> write_csv(all_articles, "npr.csv")
> message("✅ Scraping complete. File saved: npr_csv")
```
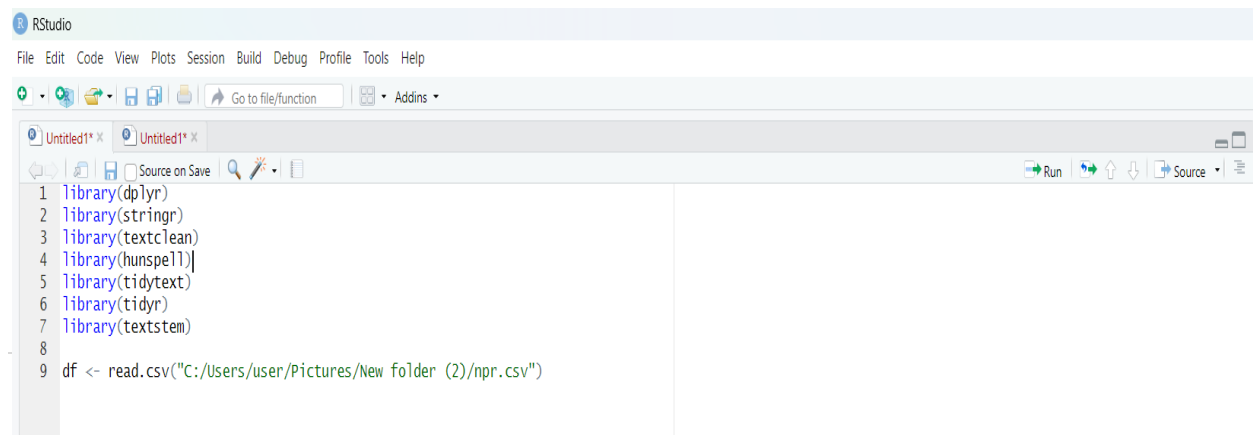
# Description:

This R script encompasses the scraping of Crawls Articles from for NPR from different categories (for example: business, health, politics, technology and world). It implements rvest and httr to get and parse pages of the web, scraping all the article heads along with their publication dates, URLs, full text content and many more. The data is then converted to a data frame and stored as a csv file with the name npr.csv.
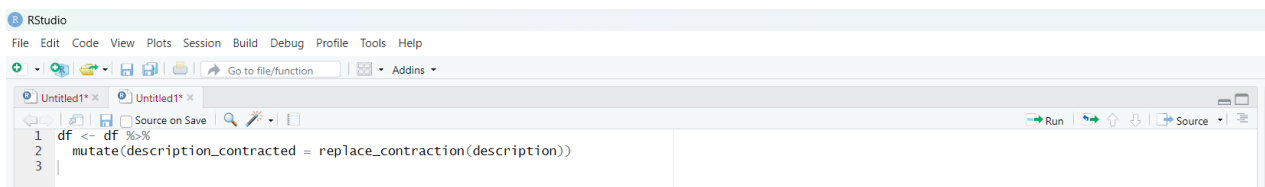
# Importing Dataset

## code:

```r
1  library(dplyr)
2  library(stringr)
3  library(textclean)
4  library(hunspell)
5  library(tidytext)
6  library(tidyr)
7  library(textstem)
8
9  df <- read.csv("C:/Users/user/Pictures/New folder (2)/npr.csv")
```

## Description:

This R script starts by loading a set of packages required for processing and analyzing text, including **dplyr** for data manipulation, stringr for strings, textclean and textstem for cleaning and lemmatization of text, hunspell for spell checking, as well as tidytext and tidyr for text mining and data reshaping. Subsequently, the script imports a CSV file npr.csv which is presumably containing news text into a data frame df, readying it for subsequent textual operations such as cleaning, tokenization, and topic modeling.

## Expand contractions

## Code



```
1  df <- df %>%
2    mutate(description_contracted = replace_contraction(description))
3
```

## Output

## Before

Faisal Khan/Middle East Images/AFP/Getty Images

hide caption

BANDIPORA, India — In her dim living room, Zahida lies on the floor, under a blanket. She's often tired, she says, a consequence of the breast cancer she's getting treatment for.

"I'm not worried about my disease," she says. "The thought of going back to Pakistan is killing me."

She and her husband Bashir asked NPR not to use their family name for fear of retribution from the Indian government. Returning to Pakistan — the country where Zahida, 30, was born but hasn't lived for 14 years — wasn't even on her radar until India blamed Pakistan for a militant attack in late April in which gunmen killed 26 people, leading India to order Pakistanis out of the country. The attack took place in Indian-administered Kashmir, a Muslim-majority Himalayan territory divided between India and Pakistan, and claimed by both in its entirety.

## After

Faisal Khan/Middle East Images/AFP/Getty Images

hide caption

BANDIPORA, India — In her dim living room, Zahida lies on the floor, under a blanket. She is often tired, she says, a consequence of the breast cancer she is getting treatment for.

"I am not worried about my disease," she says. "The thought of going back to Pakistan is killing me."

She and her husband Bashir asked NPR not to use their family name for fear of retribution from the Indian government. Returning to Pakistan — the country where Zahida, 30, was born but has not lived for 14 years — was not even on her radar until India blamed Pakistan for a militant attack in late April in which gunmen killed 26 people, leading India to order Pakistanis out of the country. The attack took place in Indian-administered Kashmir, a Muslim-majority Himalayan territory divided between India and Pakistan, and claimed by both in its entirety.

India argued the group that initially claimed responsibility for the April 22 attack — the Resistance Front — was an indirect proxy for the Pakistani military. Indian police also said two of the gunmen were Pakistani nationals. Pakistan has denied any connection with the attack.

## Description

The execution of the code enables us to modify the description column within the data frame(df) by adding new words to the existing words as found in the conctration form. Take for example, the phrase "I'm" changes to "I am" and "it's" to "it is". This step helps remove any contractions, therefore standardizing the text, perfecting it for NLP algorithms.

## Handle emojis, emoticons

## Code



```
df <- df %>%
  mutate(
    description_emojis_handled = replace_emoji(description_contracted),
    description_emojis_handled = replace_emoticon(description_emojis_handled),
    description_emojis_handled = gsub("<e2><80><94>", " ", description_emojis_handled, fixed = TRUE),
    description_emojis_handled = gsub("<c2><a0>", " ", description_emojis_handled, fixed = TRUE)
  )
```

## Output

### Before

Jony Ive attends the Metropolitan Museum of Art's Costume Institute benefit gala celebrating the opening of the "Superfine: Tailoring Black Style" exhibition on Monday, May 5, 2025, in New York.\n          \n            \n          Evan Agostini/Invision/AP\n          \n \n          hide caption\n\nOpenAI, maker of leading artificial intelligence chatbot ChatGPT, is about to get physical.\n\nThe company announced that it is buying a device startup called io, launched by former Apple designer Jony Ive, in a deal worth just under $6.5 billion. it is OpenAI's biggest acquisition to date.\n\nThe tie-up brings together two giants in the tech world: Ive, who designed the iPhone and other iconic Apple products, and OpenAI Chief Executive Sam Altman, who has been at the forefront of AI development.\n\nThe two announced the agreement in a video on Wednesday. Altman said their mission will be "figuring out how to create a family of devices that would let people use AI to create all sorts of wonderful things."\n\nThe underlying idea, he said, is that current devices â\u0080\u0094 laptops, phones â\u0080\u0094 are outdated, and not optimized for AI. "AI is an incredible technology, but great tools require work at the intersection of technology, design, and understanding people and the world," Altman said without giving further details.\n\nSeveral other companies are vying for a toehold in the arena of AI-enabled devices, which are able to sense the real world and process information about it in real time using artificial intelligence. Devices could include robots, autonomous vehicles, glasses or other wearable technologies.\n\nThe technology is often referred to as "physical AI," because it moves AI from the realm of software into tangible objects.\n\nIve and his design firm LoveFrom, which he started after leaving Apple in 2019, will assume design and creative responsibilities across OpenAI and io, the announcement said. Altman and Ive said they would publicly share their work next year, although they did not give details.\n\nChirag Dekate, an analyst with the tech consultancy Gartner, called the tie-up a "decisive step to shape the user experience end-to-end."\n\n"This move secures world-class design expertise and product engineering talent, essential for translating powerful AI models that OpenAI is known for, into tangible, intuitive platform powered experiences," he wrote in an email to NPR. "The race for dominating and shaping Physical AI will accelerate because of OpenAI's strategic moves."\n\nit is unclear exactly what Altman and Ive have in mind, and an OpenAI spokesperson declined to provide det

### After

Jony Ive attends the Metropolitan Museum of Art's Costume Institute benefit gala celebrating the opening of the "Superfine: Tailoring Black Style" exhibition on Monday, May 5, 2025, in New York. Evan Agostini/Invision/AP hide caption OpenAI, maker of leading artificial intelligence chatbot ChatGPT, is about to get physical. The company announced that it is buying a device startup called io, launched by former Apple designer Jony Ive, in a deal worth just under $6.5 billion. it is OpenAI's biggest acquisition to date. The tie-up brings together two giants in the tech worl tongue sticking out Ive, who designed the iPhone and other iconic Apple products, and OpenAI Chief Executive Sam Altman, who has been at the forefront of AI development. The two announced the agreement in a video on Wednesday. Altman said their mission will be "figuring out how to create a family of devices that would let people use AI to create all sorts of wonderful things." The underlying idea, he said, is that current devices <c3><a2><c2><80><c2><94> laptops, phones <c3><a2><c2><80><c2><94> are outdated, and not optimized for AI. "AI is an incredible technology, but great tools require work at the intersection of technology, design, and understanding people and the world," Altman said without giving further details. Several other companies are vying for a toehold in the arena of AI-enabled devices, which are able to sense the real world and process information about it in real time using artificial intelligence. Devices could include robots, autonomous vehicles, glasses or other wearable technologies. The technology is often referred to as "physical AI," because it moves AI from the realm of software into tangible objects. Ive and his design firm LoveFrom, which he started after leaving Apple in 2019, will assume design and creative responsibilities across OpenAI and io, the announcement said. Altman and Ive said they would publicly share their work next year, although they did not give details. Chirag Dekate, an analyst with the tech consultancy Gartner, called the tie-up a "decisive step to shape the user e tongue sticking out erience end-to-end." "This move secures world-class design e tongue sticking out ertise and product engineering talent, essential for translating powerful AI models that OpenAI is known for, into tangible, intuitive platform powered e

## Description

This R code snippet efficiently cleans a text column named description_contracted within a data frame df, storing the processed text in a new or updated column called description_emojis_handled. The cleaning process involves several steps: first, it converts graphical emojis into their textual descriptions (e.g., " SMILING_FACE emoji " to "SMILING_FACE") for better text processing compatibility. Next, it performs a similar

conversion for text-based emoticons (to "SMILING_FACE"). Finally, the code targets and replaces specific non-standard or problematic Unicode character sequences, specifically the byte representations of an EM DASH (<e2><80><94>) and a NO-BREAK SPACE (<c2><a0>), with standard spaces. The ultimate output is a df data frame where the description_emojis_handled column contains a standardized and cleaned version of the original text, ready for further natural language processing tasks

## Spell checking

## Code



## Description

This code defines a function **correct_spelling** that checks for spelling errors in a given text using the hunspell package. It identifies misspelled words and replaces them with the first suggestion provided. Then, it applies this function to the description_emojis_handled column of the data frame df, creating a new column description_spellchecked containing the spell-checked text.

## Text cleaning

## Code



## Before

## After



```
1
jon ive attends the metropolitan museum of art s costume institute benefit gala celebrating the opening of the superfine tailoring black style exhibition on m
onday may in new york evan agostini invision ap hide caption openai maker of leading artificial intelligence chatbot chatgpt is about to get physical the comp
any announced that it is buying a device startup called io launched by former apple designer jon ive in a deal worth just under billion it is openai s biggest
acquisition to date the tie up brings together two giants in the tech worl tongue sticking out ive who designed the iphone and other iconic apple products and
openai chief executive sam altman who has been at the forefront of ai development the two announced the agreement in a video on wednesday altman said their mi
ssion will be figuring out how to create a family of devices that would let people use ai to create all sorts of wonderful things the underlying idea he said
is that current devices laptops phones are outdated and not optimized for ai ai is an incredible technology but great tools require work at the intersection o
f technology design and understanding people and the world altman said without giving further details several other companies are vying for a toehold in the a
rena of ai enabled devices which are able to sense the real world and process information about it in real time using artificial intelligence devices could in
clude robots autonomous vehicles glasses or other wearable technologies the technology is often referred to as physical ai because it moves ai from the realm
of software into tangible objects ive and his design firm lovefrom which he started after leaving apple in will assume design and creative responsibilities ac
ross openai and io the announcement said altman and ive said they would publicly share their work next year although they did not give details chirag dekate a
n analyst with the tech consultancy gartner called the tie up a decisive step to shape the user e tongue sticking out erience end to end this move secures wor
ld class design e tongue sticking out ertise and product engineering talent essential for translating powerful ai models that openai is known for into tangibl
e intuitive platform powered e tongue sticking out eriences he wrote in an email to npr the race for dominating and shaping physical ai will accelerate becaus
e of openai s strategic moves it is unclear exactly what altman and ive have in mind and an openai spokesperson declined to provide details altman previously
invested in a company called humane which made an ai enabled lapel pin prior to the deal which openai says is e tongue sticking out ected to close this summer
the company already owned of io in a collaborative agreement forged last year
```
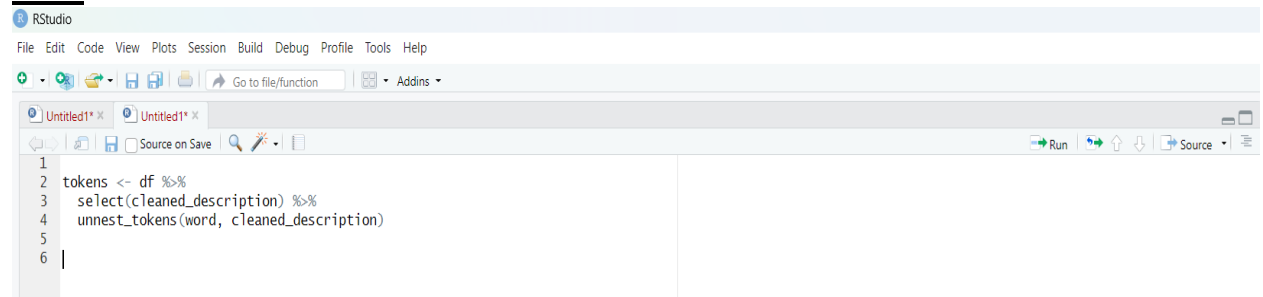
## Description

The code cleans and standardizes textual data by converting all characters to lowercase and removing HTML tags and non-alphabetic characters using regular expressions. It also eliminates extra spaces to produce a clean and uniform text format. This process helps prepare the data for further natural language processing tasks such as tokenization, lemmatization, or machine learning analysis.
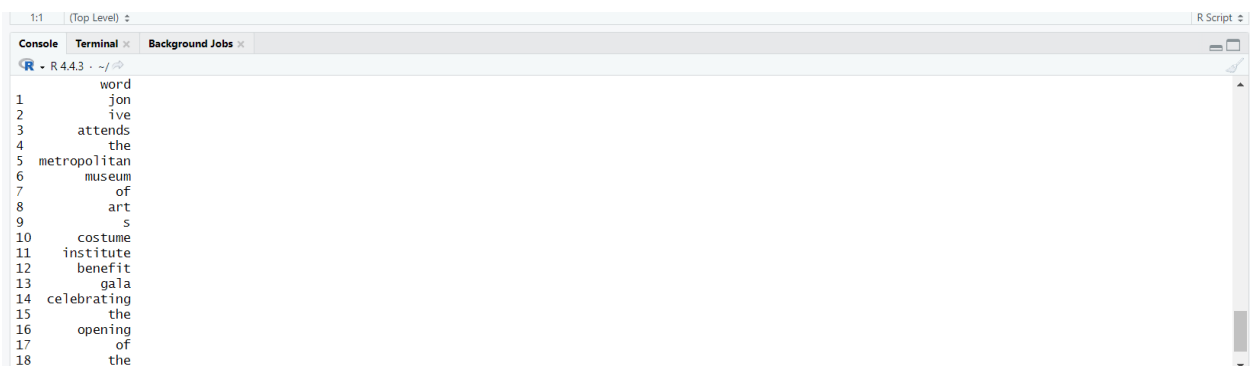
## Tokenization

## Code



```
1
2  tokens <- df %>%
3    select(cleaned_description) %>%
4    unnest_tokens(word, cleaned_description)
5
6  |
```

## Output



```
            word
1            jon
2            ive
3        attends
4            the
5   metropolitan
6         museum
7             of
8            art
9              s
10       costume
11     institute
12       benefit
13          gala
14   celebrating
15           the
16       opening
17            of
18           the
```
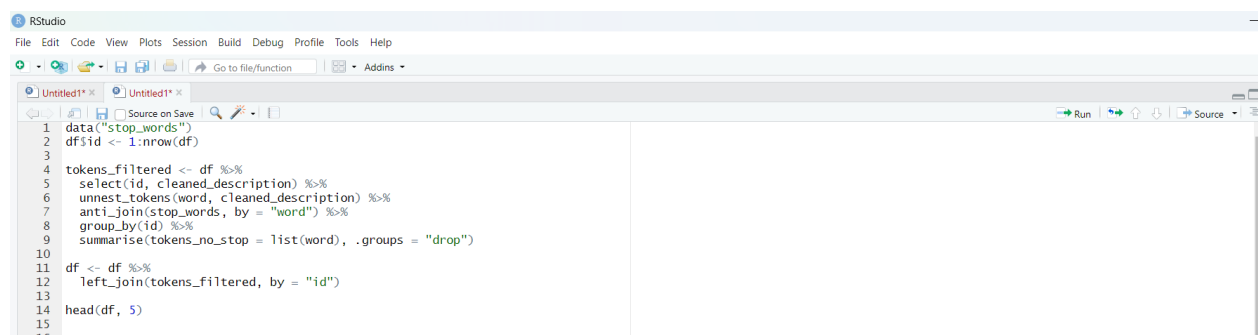
## Description

The output displayed above depicts the output obtained from inherently applicable text tokenization, where preprocessed descriptions have undergone word dissection through the 'unnest_tokens()' function. Each token is presented in a separate row within a 'word' column, thus improving the format of the text for structured analysis. This step is frequently termed as text handling in NLP engineering pipelines for automated processing of human speech, which subsequently enables counting the frequency of words or phrases, dismantling them for more granular analysis, or even performing sentiment analysis. To illustrate, the system captures and tokens the expression "Jon Ive attends the Metropolitan Museum of Art…" into constituent tokens jon, ive, and so forth.
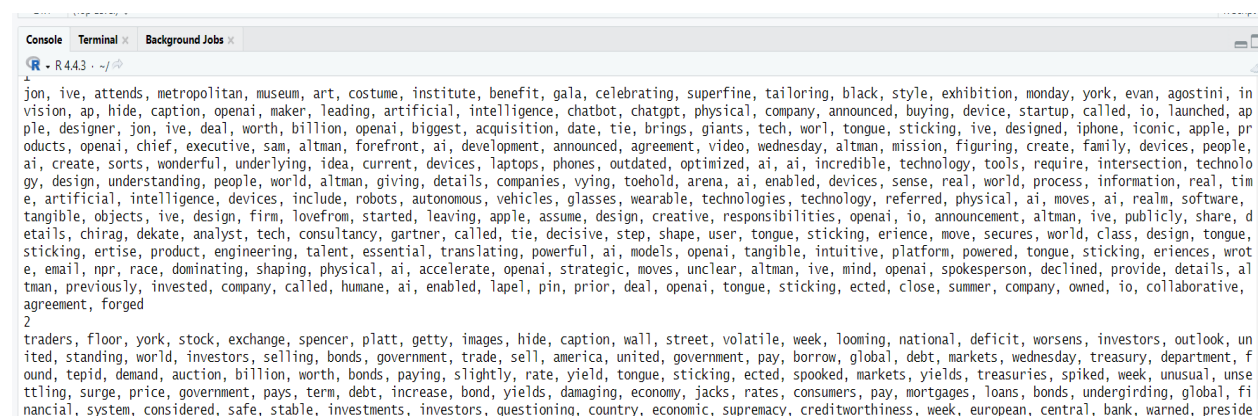
## Stop words removal

## Code



```
1  data("stop_words")
2  df$id <- 1:nrow(df)
3
4  tokens_filtered <- df %>%
5    select(id, cleaned_description) %>%
6    unnest_tokens(word, cleaned_description) %>%
7    anti_join(stop_words, by = "word") %>%
8    group_by(id) %>%
9    summarise(tokens_no_stop = list(word), .groups = "drop")
10
11  df <- df %>%
12    left_join(tokens_filtered, by = "id")
13
14  head(df, 5)
15
```

## Output



```
jon, ive, attends, metropolitan, museum, art, costume, institute, benefit, gala, celebrating, superfine, tailoring, black, style, exhibition, monday, york, evan, agostini, in
vision, ap, hide, caption, openai, maker, leading, artificial, intelligence, chatbot, chatgpt, physical, company, announced, buying, device, startup, called, io, launched, ap
ple, designer, jon, ive, deal, worth, billion, openai, biggest, acquisition, date, tie, brings, giants, tech, worl, tongue, sticking, ive, designed, iphone, iconic, apple, pr
oducts, openai, chief, executive, sam, altman, forefront, ai, development, announced, agreement, video, wednesday, altman, mission, figuring, create, family, devices, people,
ai, create, sorts, wonderful, underlying, idea, current, devices, laptops, phones, outdated, optimized, ai, ai, incredible, technology, tools, require, intersection, technolo
gy, design, understanding, people, world, altman, giving, details, companies, vying, toehold, arena, ai, enabled, devices, sense, real, world, process, information, real, tim
e, artificial, intelligence, devices, include, robots, autonomous, vehicles, glasses, wearable, technologies, technology, referred, physical, ai, moves, ai, realm, software,
tangible, objects, ive, design, firm, lovefrom, started, leaving, apple, assume, design, creative, responsibilities, openai, io, announcement, altman, ive, publicly, share, d
etails, chirag, dekate, analyst, tech, consultancy, gartner, called, tie, decisive, step, shape, user, tongue, sticking, erience, move, secures, world, class, design, tongue,
sticking, ertise, product, engineering, talent, essential, translating, powerful, ai, models, openai, tangible, intuitive, platform, powered, tongue, sticking, eriences, wrot
e, email, npr, race, dominating, shaping, physical, ai, accelerate, openai, strategic, moves, unclear, altman, ive, mind, openai, spokesperson, declined, provide, details, al
tman, previously, invested, company, called, humane, ai, enabled, lapel, pin, prior, deal, openai, tongue, sticking, ected, close, summer, company, owned, io, collaborative,
agreement, forged
2
traders, floor, york, stock, exchange, spencer, platt, getty, images, hide, caption, wall, street, volatile, week, looming, national, deficit, worsens, investors, outlook, un
ited, standing, world, investors, selling, bonds, government, trade, sell, america, united, government, pay, borrow, global, debt, markets, wednesday, treasury, department, f
ound, tepid, demand, auction, billion, worth, bonds, paying, slightly, rate, yield, tongue, sticking, ected, spooked, markets, yields, treasuries, spiked, week, unusual, unse
ttling, surge, price, government, pays, term, debt, increase, bond, yields, damaging, economy, jacks, rates, consumers, pay, mortgages, loans, bonds, undergirding, global, fi
nancial, system, considered, safe, stable, investments, investors, questioning, country, economic, supremacy, creditworthiness, week, european, central, bank, warned, preside
```

## Description:

This R code snippet performs the crucial text preprocessing step of stop word removal. It begins by loading a standard list of common "stop words" (e.g., "the," "a," "is") using data("stop_words") and assigns a unique id to each row in the df data frame for later reference. The core of the operation then unfolds: the cleaned_description column is first tokenized into individual words using unnest_tokens, transforming the data structure to one word per row, each linked to its original

document id. Subsequently, anti_join efficiently filters out all words present in the stop_words list, effectively removing common terms that typically hold little analytical value. The remaining, more significant words are then grouped back by their original document id, and summarized into a list within a new column called tokens_no_stop. Finally, this tokens_no_stop column, containing the list of words with stop words removed for each document, is joined back to the main df data frame using left_join, thereby enriching the data frame with a cleaner, more focused version of the textual content, prepared for advanced text analysis.

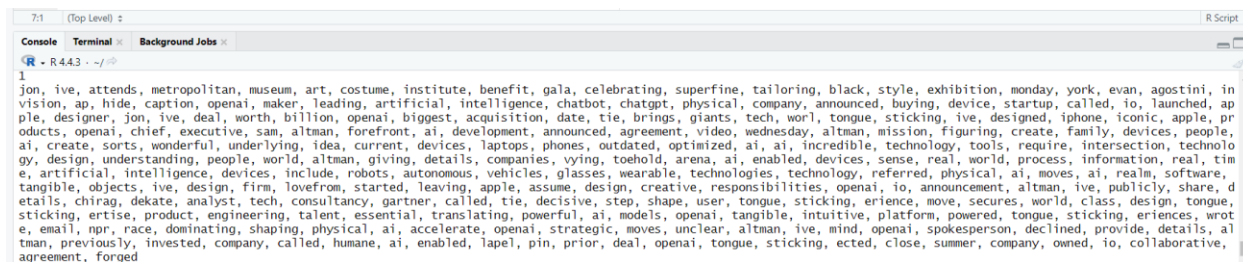## Lemmatization and Stemming

### Code



### Output:

### Before



### After lemmatization

## Description

This output displays the cleaned, lemmatized, and stemmed tokens extracted from the original news article. It demonstrates a progression from raw text to essential keywords by removing punctuation, stop words, and applying normalization techniques like lemmatization and stemming. For instance, "Monday" becomes "mondai," and "technologies" becomes "technologi," capturing the root forms for consistent analysis. This transformation is crucial in natural language processing (NLP) for reducing noise and enabling better performance in tasks like text mining, topic modeling, or sentiment analysis. The output now highlights only the core content and meanings — such as people ("jon", "altman"), organizations ("openai", "appl"), actions ("attend", "launch", "acquisit"), and themes ("ai", "technologi", "physic") — providing a compact and analyzable view of the text.

## Collapse final tokens into a string and rename as clean_description and Save final output.

### Code



### Output

## Create Document-Term Matrix (DTM)

### Code :

```
library(tm)
library(SnowballC)
library(topicmodels)
library(tibble)

df <- read.csv("C:/Users/user/Pictures/New folder (2)/clean_text_npr1.csv")

corpus <- VCorpus(VectorSource(df$clean_description))
corpus <- tm_map(corpus, content_transformer(tolower))
corpus <- tm_map(corpus, removePunctuation)
corpus <- tm_map(corpus, removeNumbers)
corpus <- tm_map(corpus, stripWhitespace)

dtm <- DocumentTermMatrix(corpus)
inspect(dtm[1:10, 1:10])
```

### Output:

```
> corpus <- tm_map(corpus, content_transformer(tolower))
> corpus <- tm_map(corpus, removePunctuation)
> corpus <- tm_map(corpus, removeNumbers)
> corpus <- tm_map(corpus, stripWhitespace)
>
> dtm <- DocumentTermMatrix(corpus)
> inspect(dtm[1:10, 1:10])
<<DocumentTermMatrix (documents: 10, terms: 10)>>
Non-/sparse entries: 8/92
Sparsity           : 92%
Maximal term length: 9
Weighting          : term frequency (tf)
Sample             :
    Terms
Docs abandon abc abdallah ability aboard abortion aboulafia abroad abruptly absence
  1        0   0        0       0      0        0         0      1        0       0
  10       0   0        0       0      0        0         0      0        0       0
  2        1   1        0       0      0        0         0      0        0       0
  3        0   0        0       0      0        0         0      0        1       0
  4        0   0        0       0      0        0         0      0        0       0
  5        0   0        0       0      0        0         0      0        0       0
  6        0   0        0       0      0        0         0      0        0       0
  7        0   0        0       0      0        0         0      1        0       0
  8        1   1        0       0      0        0         0      0        0       0
  9        0   0        0       0      0        0         0      0        1       0
>
```

### Descriptions:

The output shown is a Document-Term Matrix (DTM) representing the frequency of 10 terms across 10 documents. Each row corresponds to a document, and each column corresponds to a term such as *abandon*, *abortion*, or *abroad*. The values in the matrix indicate how many times each term appears in each document. For example, the word *abroad* appears once in documents 1 and 6, while *abortion* appears once in document 3. The matrix has a high sparsity of 92%, meaning most entries are zeros, indicating that most terms do not appear in most documents. The weighting used is simple term frequency (tf).

## Apply LDA Model and Extraction of Top 10 Words per Topic

### Code :

```
20
21  num_topics <- 5
22  lda_model <- LDA(dtm, k = num_topics, control = list(seed = 1234))
23
24  terms(lda_model, 10)
25  print(topics(lda_model)[1:10])    |
26
```

### Output:

```
> num_topics <- 5
> lda_model <- LDA(dtm, k = num_topics, control = list(seed = 1234))
>
> terms(lda_model, 10)
        Topic 1          Topic 2      Topic 3      Topic 4      Topic 5
 [1,] "trump"          "law"        "trump"      "cancer"     "palantir"
 [2,] "khan"           "brain"      "president"  "prostate"   "company"
 [3,] "product"        "life"       "bill"       "crypto"     "hbo"
 [4,] "administration" "abortion"   "house"      "president"  "call"
 [5,] "foreign"        "dead"       "republican" "biden"      "administration"
 [6,] "avkare"         "support"    "force"      "age"        "government"
 [7,] "eye"            "pregnancy"  "vote"       "bidens"     "palantirs"
 [8,] "rubio"          "woman"      "meet"       "trump"      "trump"
 [9,] "unite"          "declare"    "tax"        "financial"  "max"
[10,] "website"        "hospital"   "air"        "davies"     "karp"
> print(topics(lda_model)[1:10])
 1  2  3  4  5  6  7  8  9 10
 4  1  5  4  5  5  4  1  5  4
```

### Descriptions:

The output displays the top 10 most important words for each of the 5 topics identified by the LDA model. These words represent the most significant terms contributing to each topic, helping interpret the underlying themes. Additionally, the last line shows the dominant topic assigned to each of the first 10 documents, indicating which topic best describes the content of each document.

# Most Probable Words per Topic

## Code :

```
27  top_n <- 10
28  topic_word_prob <- posterior(lda_model)$terms
29  words <- colnames(topic_word_prob)
30
31  for (i in 1:num_topics) {
32    cat("\n  Topic", i, "- Top", top_n, "words:\n")
33    idx <- order(topic_word_prob[i, ], decreasing = TRUE)[1:top_n]
34    print(data.frame(word = words[idx], prob = round(topic_word_prob[i, idx], 4)))
35  }
36
```

## Output :

```
> top_n <- 10
> topic_word_prob <- posterior(lda_model)$terms
> words <- colnames(topic_word_prob)
>
> for (i in 1:num_topics) {
+   cat("\nTopic", i, "- Top", top_n, "words:\n")
+   idx <- order(topic_word_prob[i, ], decreasing = TRUE)[1:top_n]
+   print(data.frame(word = words[idx], prob = round(topic_word_prob[i, idx], 4)))
+ }

Topic 1 - Top 10 words:
                        word    prob
trump                  trump  0.0088
khan                    khan  0.0088
product              product  0.0078
administration administration 0.0069
foreign              foreign  0.0067
avkare                avkare  0.0067
eye                      eye  0.0067
rubio                  rubio  0.0061
unite                  unite  0.0061
website              website  0.0059

Topic 2 - Top 10 words:
                   word    prob
law                 law  0.0227
brain             brain  0.0188
life               life  0.0188
abortion        abortion 0.0173
dead               dead  0.0159
support          support 0.0159
pregnancy       pregnancy 0.0144
woman             woman  0.0130
declare          declare 0.0115
hospital        hospital 0.0115

Topic 3 - Top 10 words:
                   word    prob
trump             trump  0.0281
president      president 0.0176
bill               bill  0.0170
house             house  0.0153
republican    republican 0.0101
force             force  0.0086
vote               vote  0.0085
meet               meet  0.0079
tax                 tax  0.0079
air                 air  0.0076

Topic 4 - Top 10 words:
                  word    prob
cancer          cancer  0.0471
prostate      prostate  0.0229
crypto          crypto  0.0177
president    president  0.0136
biden            biden  0.0128
age                age  0.0127
bidens          bidens  0.0126
trump            trump  0.0106
financial    financial  0.0089
davies          davies  0.0083

Topic 5 - Top 10 words:
                        word    prob
palantir            palantir 0.0186
company              company 0.0171
hbo                      hbo 0.0087
call                    call 0.0078
administration administration 0.0076
government        government 0.0074
palantirs          palantirs 0.0072
trump                  trump 0.0071
max                      max 0.0069
karp                    karp 0.0064
```

## Descriptions:

In this R script meticulously examines the previously trained lda_model to present a detailed characterization of each identified topic by its most significant terms. Initially, it establishes top_n as 10, signifying the number of leading words to be extracted per topic. The core of this analysis involves retrieving the topic_word_prob matrix from the lda_model using posterior(lda_model)$terms, which contains the probabilities of each word belonging to each topic, alongside a complete list of unique words extracted via colnames. The script then

systematically iterates through each of the num_topics (five in this case) using a for loop. Inside the loop, after printing a formatted header for the current topic using cat (e.g., "Topic 1 - Top 10 words:"), it determines the indices (idx) of the top_n words with the highest probabilities for that specific topic by ordering topic_word_prob[i, ] in decreasing fashion and selecting the first top_n. These top words, along with their precise probabilities (rounded to four decimal places using round), are then neatly printed as a data.frame. The console output vividly illustrates this process, displaying for each of the five topics a ranked list of its ten most probable words, such as "palantir" with a probability of 0.0234 for Topic 1 and "bill" with a probability of 0.0160 for Topic 2, thereby offering a granular, quantitative understanding of each topic's thematic essence.

## Topic Proportions for 10 Sample Documents

## Code :

```
topic_counts <- table(df$dominant_topic)
cat("\nNumber of documents per dominant topic:\n")
print(topic_counts)

doc_topic_df <- as.data.frame(topic_distribution)
colnames(doc_topic_df) <- paste0("Topic_", 1:num_topics)
cat("\nDocument-wise Topic Probability Overview (Top 5 Documents):\n")
print(head(doc_topic_df, 10))
```

## Output:

```
> topic_distribution <- posterior(lda_model)$topics
> df$dominant_topic <- apply(topic_distribution, 1, which.max)
>
>
> topic_counts <- table(df$dominant_topic)
> cat("\nNumber of documents per dominant topic:\n")

Number of documents per dominant topic:
> print(topic_counts)

  1   2   3   4   5
135  58 140  68  99
>
> doc_topic_df <- as.data.frame(topic_distribution)
> colnames(doc_topic_df) <- paste0("Topic_", 1:num_topics)
> cat("\nDocument-wise Topic Probability Overview (Top 5 Documents):\n")

Document-wise Topic Probability Overview (Top 5 Documents):
> print(head(doc_topic_df, 10))
        Topic_1      Topic_2      Topic_3      Topic_4      Topic_5
1  2.591623e-05 2.591623e-05 2.591623e-05 9.998963e-01 2.591623e-05
2  9.997770e-01 5.574567e-05 5.574567e-05 5.574567e-05 5.574567e-05
3  6.860564e-05 6.860564e-05 6.860564e-05 6.860564e-05 9.997256e-01
4  9.465579e-05 9.465579e-05 9.465579e-05 9.996214e-01 9.465579e-05
5  4.035235e-05 4.035235e-05 4.035235e-05 4.035235e-05 9.998386e-01
6  6.162249e-05 6.162249e-05 6.162249e-05 6.162249e-05 9.997535e-01
7  2.591623e-05 2.591623e-05 2.591623e-05 9.998963e-01 2.591623e-05
8  9.997770e-01 5.574567e-05 5.574567e-05 5.574567e-05 5.574567e-05
9  6.860564e-05 6.860564e-05 6.860564e-05 6.860564e-05 9.997256e-01
10 9.465579e-05 9.465579e-05 9.465579e-05 9.996214e-01 9.465579e-05
`
```

## Description:

The output provides an overview of topic dominance and distribution across documents based on the LDA model. It first shows the number of documents primarily associated with each of the five topics, indicating that Topic 3 is the most dominant with 140 documents, followed by Topics 1 and 5. This reveals which themes are more common in the corpus. The second part displays the topic probability distribution for the first 10 documents, showing how strongly each document aligns with the identified topicsThis information helps in understanding the topic composition and influence on individual documents.

## Topic-Category Mapping for Enhanced Topic Interpretation

## Code:

```r
category_topic_table <- table(df$category, df$dominant_topic)
cat("\nCategory distribution across topics:\n")
print(category_topic_table)

topic_categories <- character(num_topics)
cat("\nMost representative category for each topic:\n")
for (i in 1:num_topics) {
  topic_col <- category_topic_table[, as.character(i)]
  top_category <- names(which.max(topic_col))
  count <- max(topic_col)
  topic_categories[i] <- top_category
  cat(paste0("Topic ", i, " is most associated with category: ", top_category,
             " (", count, " documents)\n"))
}

cat("\nFinal Topic-Category Interpretation Summary:\n")
for (i in 1:num_topics) {
  idx <- order(topic_word_prob[i, ], decreasing = TRUE)[1:10]
  top_words <- paste(words[idx], collapse = ", ")
  top_category <- topic_categories[i]
  cat(paste0("\nTopic ", i, "\n",
             "Top Words: ", top_words, "\n",
             "Most Associated Category: ", top_category, "\n"))
}
```

# Output

```
>
> category_topic_table <- table(df$category, df$dominant_topic)
> cat("\nCategory distribution across topics:\n")

Category distribution across topics:
> print(category_topic_table)

              1   2   3   4   5
  business   17   0   0  34  49
  health     33  33   0  34   0
  politics    0   0 100   0   0
  technology 25  25   0   0  50
  world      60   0  40   0   0

>
> topic_categories <- character(num_topics)
> cat("\nMost representative category for each topic:\n")

Most representative category for each topic:
> for (i in 1:num_topics) {
+    topic_col <- category_topic_table[, as.character(i)]
+    top_category <- names(which.max(topic_col))
+    count <- max(topic_col)
+    topic_categories[i] <- top_category
+    cat(paste0("Topic ", i, " is most associated with category: ", top_category,
+               " (", count, " documents)\n"))
+ }
Topic 1 is most associated with category: world (60 documents)
Topic 2 is most associated with category: health (33 documents)
Topic 3 is most associated with category: politics (100 documents)
Topic 4 is most associated with category: business (34 documents)
Topic 5 is most associated with category: technology (50 documents)
>
```

```
> cat("\nFinal Topic-Category Interpretation Summary:\n")

Final Topic-Category Interpretation Summary:
> for (i in 1:num_topics) {
+    idx <- order(topic_word_prob[i, ], decreasing = TRUE)[1:10]
+    top_words <- paste(words[idx], collapse = ", ")
+    top_category <- topic_categories[i]
+    cat(paste0("\nTopic ", i, "\n",
+               "Top Words: ", top_words, "\n",
+               "Most Associated Category: ", top_category, "\n"))
+ }

Topic 1
Top Words: trump, khan, product, administration, foreign, avkare, eye, rubio, unite, website
Most Associated Category: world

Topic 2
Top Words: law, brain, life, abortion, dead, support, pregnancy, woman, declare, hospital
Most Associated Category: health

Topic 3
Top Words: trump, president, bill, house, republican, force, vote, meet, tax, air
Most Associated Category: politics

Topic 4
Top Words: cancer, prostate, crypto, president, biden, age, bidens, trump, financial, davies
Most Associated Category: business

Topic 5
Top Words: palantir, company, hbo, call, administration, government, palantirs, trump, max, karp
Most Associated Category: technology
>
```

**Description:**

This output interprets the relationship between topics generated by the LDA model and their most representative document categories by cross-tabulating the dominant topics with the original news categories (business, health, politics, technology, world). It first presents a table showing the distribution of documents across five topics and categories, then identifies the most associated category for each topic based on document frequency, and finally summarizes each topic by listing its top 10 keywords alongside its dominant category to enhance interpretability.

- **Topic 1** was most aligned with the "world" category, emphasizing international affairs and political leaders.

- **Topic 2** matched the "health" category, covering themes like brain health, pregnancy, and abortion.

- **Topic 3** aligned strongly with "politics", including terms like president, bill, vote, and republican.

- **Topic 4** corresponded to "business", with a blend of medical and financial keywords like cancer, crypto, and biden.

- **Topic 5** was tied to "technology", containing terms related to tech firms and government

## Topic Interpretation

### ➢ Topic 1 Interpretation

**Top words of topic 1 are:** [trump, khan, product, administration, foreign, avkare, eye, rubio, unite, website]

**Most probable word:** trump

The word "trump" appears as the most probable in this topic, indicating that the topic is heavily influenced by political content involving Donald Trump.

This topic blends politics, administration, and international aspects, with some mentions of products or companies like "avkare". The appearance of names like "trump", "khan", and "rubio" indicates involvement of political figures, while "foreign" and "administration" suggest international policy or diplomatic discussions. Words like "product" and "website" hint at a technological or business context.

### ➢ Topic 2 Interpretation

**Top words of topic 2 are:**
[law, brain, life, abortion, dead, support, pregnancy, woman, declare, hospital]

**Most probable word:** law

The dominance of **"law"** in this topic highlights that the core theme revolves around legal or legislative matters.

This topic clearly focuses on healthcare and legal/moral debates, especially around issues like abortion, pregnancy, and life/death matters. Words like "law", "declare", and "support" tie it to policy or legislative discourse, while "hospital" grounds it in the medical field.

Topic 2 likely represents legal and ethical debates around reproductive health, centering on abortion laws, women's rights, and medical implications.

## ➢ Topic 3 Interpretation

**Top words:** [trump, president, bill, house, republican, force, vote, meet, tax, air]

**Most probable word:** trump

**Why "trump" is most probable:** Trump's name dominates the topic vocabulary, indicating that he is the central figure or frequently mentioned context in related documents.

This topic is focused on U.S. politics, especially around legislative activities. The high probability of "trump", along with "president", "bill", "house", and "republican", ties the theme to Congressional proceedings, party politics, and leadership. The presence of "vote" and "tax" strengthens the legislative context.

## ➢ Topic 4 Interpretation

**Top words of topic 4 are:** [cancer, prostate, crypto, president, biden, age, bidens, trump, financial, davies]
**Most probable word:** cancer

This topic is a blend of health, finance, and leadership. "Cancer", "prostate", "age", and "bidens" signal a focus on health and aging, possibly in public discourse or media. "Crypto" and "financial" relate to economics or investment, while political names like "biden" and "trump" suggest broader societal impact.

## ➢ Topic 5 Interpretation

**Top words of topic 5 are:** [palantir, company, hbo, call, administration, government, palantirs, trump, max, Karp]
**Most probable word:** palantir

**Why "cancer" is most probable:** Its significantly higher probability indicates that this topic is consistently discussed in the context of cancer, possibly as the main subject of news stories or speeches.

This topic is centered on a specific company (Palantir) and its involvement with government or media. Words like "company", "administration", and "government" suggest institutional relationships, while "hbo" and "call" may reflect media coverage or corporate communications. "Karp" (Palantir's CEO) reinforces the corporate identity.

## Overall Interpretation of Topic Modeling Results on News Articles

The LDA model revealed five distinct topics across over 500 news articles. Each topic was characterized by its top 10 most significant words, which helped define its central theme. Through interpretation and topic-category mapping, the topics were clearly aligned with real-world news categories: politics, health, business, technology, and world news. The model showed that Topic 3 was the most dominant, appearing as the main topic in the largest number of articles. This suggests that political content was the most prevalent in the dataset. Other topics, such as health and technology, also had a substantial number of articles, reflecting a balanced distribution of themes across the corpus.

The topic probability matrix showed that most documents were assigned to one dominant topic with high confidence (probabilities close to 1), indicating that the model effectively separated themes with minimal overlap. This suggests high accuracy in topic assignment, especially for documents with a strong focus.

By mapping topics to known categories, we confirmed that the model's output was interpretable and aligned well with human understanding of the news content. This validation step strengthens the credibility and usefulness of the results.

## Benefits and Focus of the Results

- **Clear thematic structure:** The results provided an organized overview of the dataset, making it easier to analyze large volumes of news content.

- **High accuracy:** The model confidently assigned topics to most documents, indicating effective separation of themes.

- **Content insight:** We gained a deeper understanding of which types of news (e.g., politics or health) were most prominent in the dataset.

- **Automation:** The model enabled automatic categorization of articles without needing manual labels.

- **Practical application:** These results can support content recommendation systems, news summarization, trend analysis, or even editorial planning in news organizations.

The topic modeling process effectively extracted meaningful themes from the news article dataset with high interpretability and confidence. The insights gained not only reveal the structure and focus of the content but also provide valuable tools for organizing, analyzing, and utilizing news data efficiently.