



Sharif University of Technology

Masoud Tahmasbi Fard



Student ID: 402200275

EE577: Massive Data Analytics

Assignment #1- Introduction to [PySpark](#)

December 19, 2023

Table of Contents

1	Data Preparation	1
2	Preprocessing	1
3	Exploration	2
3.1	What are the titles and urls of 5 longest news?	2
3.2	What are the 20 most frequent words?	3
3.3	Find the <i>word clouds</i> of the news	4
3.4	Plot a timeline which shows total count of news per each day	5
3.5	Plot a timeline which shows total count of news per each day that have at least one of the word clouds	6
4	A-Priori Algorithm	7
4.1	TF-IDF	9

1 Data Preparation

To begin the analysis, the PySpark libraries were installed on Google Colab to allow for distributed processing capabilities. The Persian news dataset was then downloaded from Google Drive for use in this analysis.

A Spark session was initialized to enable access to PySpark functionality. The news data files were stored in JSON format across several files. Using the spark session, these JSON files were read into a Resilient Distributed Dataset (RDD) called *news-rdd* to facilitate parallel data processing.

An example of one element of the *news-rdd* is shown below:

[illegible]

Figure 1: An element of the *news-rdd*

2 Preprocessing

To prepare the raw news data for analysis, multiple preprocessing steps were applied using PySpark transformations and custom functions.

First, the *news-rdd* which contained the raw JSON data was parsed into a new RDD called *parsed-rdd*.

Two preprocessing functions were then defined and applied to clean and normalize the data. The *persian-only* function removed any non-Persian characters and unicode strings, while the *txt-cleansing* function standardized punctuation and removed special characters.

By applying these functions to the *parsed-rdd*, a *clean-news-rdd* RDD was produced with cleaned and normalized Persian news text. An example element is shown below:

```
print(clean_news_rdd.first())  
{'body': ['نساجی' و 'مستطیل' و 'پای' و 'دقیقه' و 'گزارش' و 'خود' و 'سبز' و 'دگه' و 'این' و 'دریازه' و 'فوتبال' و 'تیک' و 'نما' و 'مردمان' و 'گرفت']
```

Figure 2: An element of the *clean-news-rdd*

Next, to improve search and analysis, a *stop-remove-lemmatizer* function leveraged the [hazm](#) library to remove stop words and lemmatize the terms. A stopwords text file gathered for Persian was attached and utilized to remove common words. Applying this function resulted in the *processed-news-rdd* which contained cleaned, normalized, and simplified news content. An example row after these preprocessing steps is:

```
print(processed_news_rdd.first())  
{'body': ['و' و 'ادامه' و 'وصفانندی' و 'میخندیم' و 'مهدی' و 'فوتبال' و 'برایمان' و 'سکو' و 'تویزیون' و 'تیک' و 'تاریخ' و 'برسر' و 'داده' و 'نمیآید' و 'مرد']
```

Figure 3: An element of the *processed-news-rdd*

The full Persian text was now ready for next-step analysis and modeling with unnecessary characters, stop words and unlemmatized terms removed through PySpark and hazm data manipulation.

3 Exploration

3.1 What are the titles and urls of 5 longest news?

To identify the longest news stories, a custom *word-count* function was implemented to return a tuple containing the news length as number of words, along with the title and url. This

leveraged the built-in Python `len` function to count words in each news excerpt.

The *word-count* function was applied to the *clean-news-rdd* using a *map* transformation, creating keyword-length tuples for each row. The results were reduced by the keyword to consolidate news with same titles.

Finally, the *reduceByKey* method was used to sort the RDD by length and then the *takeOrdered* method was utilized to retrieve the top 5 longest news. (For simplicity in further analyses, the *top* method was utilized instead of *takeOrdered* to find top results based on the given ordering).

This returned the following 5 longest news based on word count:

```
[
  (2359,
    {'title': 'الاول بدائيہ منافع دقيق و معتبر u200c\ہرچہ بايد دربارہ اہم رتبہ',
      'url': 'https://www.borna.news/بخش-%D9%82%D8%B1%D8%A2%D9%86-%D9%85%D8%B9%D8%A7%D8%B1%D9%81-53/920622-%D9%87%D8%B1%DA%86%D9%87-%D8%A8%D8%A7%D8%8C%D8%AF-%D8%AF%D8%B1%D8%A8%D8%A7%D8%B1%D9%87-%D9%86%D9%87%D9%85-%D8%B1%D8%A8%DB%8C%D8%B9-%D8%A7%D9%84%D8%A7%D9%88%D9%84-%D8%A8%D8%AF%D8%A7%D9%86%DB%8C%D8%AF-%D9%85%D9%86%D8%A7%D8%A8%D8%B9-%D8%AF%D9%82%DB%8C%D9%82-%D9%85%D8%B9%D8%AA%D8%A8%D8%B1'}),
  (2185,
    {'title': 'از امنيت جانوران درندہ تا آباداني کرات ديگر در عصر امام زمان(عج)',
      'url': 'https://www.khabaronline.ir/news/1815236/عج-امام-زمان-در-عصر-امام-کرات-ديگر-در-امنيت-جانوران-درندہ-تا-آباداني-کرات-ديگر'}),
  (2160,
    {'title': '\n',
      'url': 'http://www.mazandnume.com/fullcontent/81555/%D8%A7%D8%B2-%D9%85%DB%8C%D8%A7%D9%86%D9%90-%D8%A8%D8%A7%D9%86%DA%AF-%DA%AF%D8%B1%D8%AF%D8%B4%E2%80%8C%D9%87%D8%A7%DB%8C-%DA%86%D8%B1%D8%AE/'}),
  (2160,
    {'title': '\n',
      'url': 'http://mazandnume.com/fullcontent/81555/%D8%A7%D8%B2-%D9%85%DB%8C%D8%A7%D9%86%D9%90-%D8%A8%D8%A7%D9%86%DA%AF-%DA%AF%D8%B1%D8%AF%D8%B4%E2%80%8C%D9%87%D8%A7%DB%8C-%DA%86%D8%B1%D8%AE/'}),
  (2145,
    {'title': 'خبرگزاری فارس - آيا دولت سیزدهم گنتمان و راهبرد فرهنگی دارد؟',
      'url': 'https://www.farsnews.ir/news/14020625000557/دارد-فرهنگي-و-راهبرد-گنتمان-آيا-دولت-سیزدهم-گنتمان-دارد'})]
```

Figure 4: titles and urls of 5 longest news

3.2 What are the 20 most frequent words?

A custom *words-Key Value* function was created to tokenize each news article and output tuples of (*word*, 1) to count individual word frequencies.

This function was applied using a *flatMap* transformation to split all text into discrete words and counts. The results were reduced by the word to sum all 1s and consolidate unique terms. Finally, the *sortByKey* method sorted the aggregated word counts RDD and returned the 20 most common tokens overall, visualized in the following bar chart:

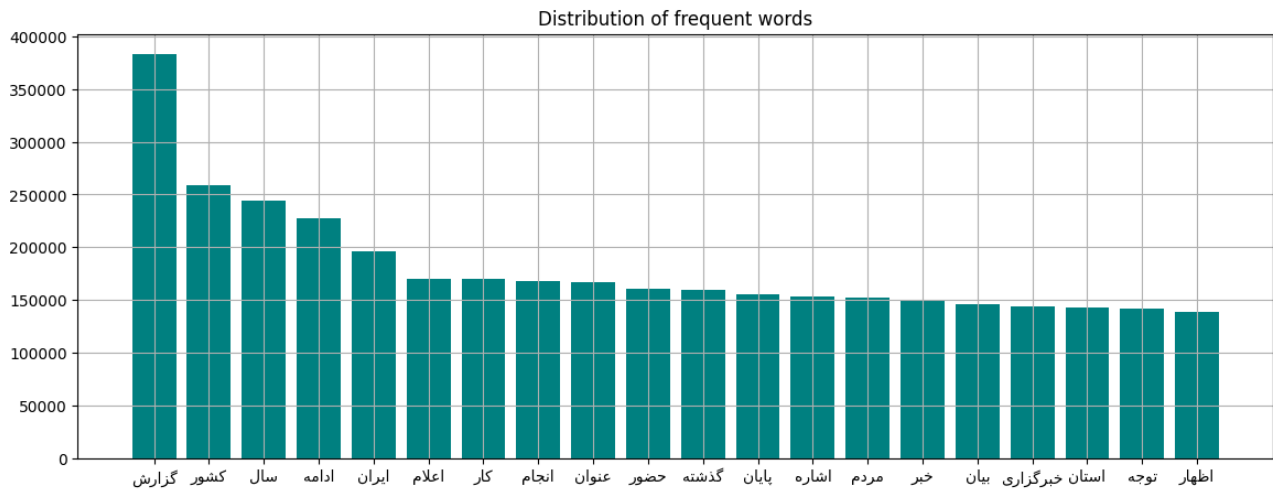


Figure 5: 20 most frequent words

To properly display the Persian text, the *bidi.algorithm* and *arabic-reshaper* libraries were employed to render the words.

Frequent terms match expectation for a Persian news dataset, indicating meaningful text analysis.

3.3 Find the *word clouds* of the news

To visualize prominent topics, *word clouds* were generated from the keyword tags associated with each news.

A *keywords-Key Value* function was implemented that checked if a given keyword existed in an news' keywords or not.

This function was applied using map and reduce transformations on the full dataset to count the occurrences of each keyword across all news stories.

The results were then scored and sorted by importance according to:

$$\text{score} = 2^{\log_{10}(\text{count})}$$

Finally, the scores were normalized by the maximum value to derive weights. Generating word clouds using these weights produced the following visualizations, sized by significance:

```
[('رژیم صهیونیستی', 1.0),
 ('اسرائیل', 0.9153835361235024),
 ('جمهوری اسلامی ایران', 0.8212246921831099),
 ('غزه', 0.7451676515885112),
 ('لیگ قهرمانان آسیا', 0.7394615313215953),
 ('تیم ملی', 0.7363438588848646),
 ('آموزش و پرورش', 0.7239635116902534),
 ('جمهوری اسلامی', 0.6955555918790012),
 ('فلسطین', 0.6675887614508967),
 ('ایالات متحده', 0.6580118424770381),
 ('حماس', 0.655299914935207),
 ('اسلام', 0.6533581120125306),
 ('طوفان الاقصی', 0.6400008319358191),
 ('استقلال', 0.6399173808248128),
 ('رئیس جمهور', 0.6394579473932916),
 ('سیستان و بلوچستان', 0.6250389114260585),
 ('بانک مرکزی', 0.6229149527751284),
 ('سازمان ملل', 0.6116172477928972),
 ('لایحه برنامه هفتم توسعه', 0.5958543790760519),
 ('مجلس شورای اسلامی', 0.5955093942849485)]
```

Figure 6: word clouds



Figure 7: word clouds

3.4 Plot a timeline which shows total count of news per each day

To examine trends over time, a timeline was generated showing the total number of articles published each day.

The *date-published* field containing the date for each news story was extracted as the key. The

value was set to 1 to represent one occurrence per article.

These $(date, 1)$ tuples were mapped across the dataset and then reduced by the date key to sum all instances published on each day.

The following figure showing the resulting time series plot is provided to visualize the trends over time in news volume.

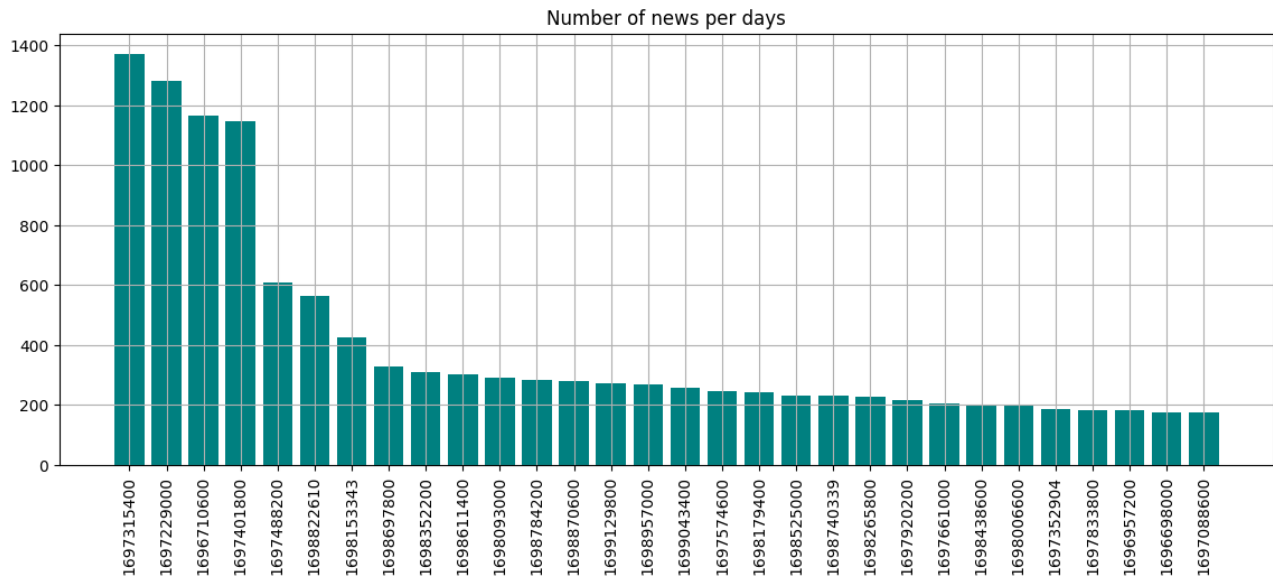


Figure 8: *word clouds*

3.5 Plot a timeline which shows total count of news per each day that have at least one of the word clouds

To narrow down on highly relevant and trending news cycles, the time series analysis was extended to only count articles containing at least one keyword from the significant terms identified in the word clouds.

The same date extraction and counting process was applied, but now combined with a check if any word cloud words appeared in the keywords for that article. So the date tally only incremented if overlapping semantic keywords existed.

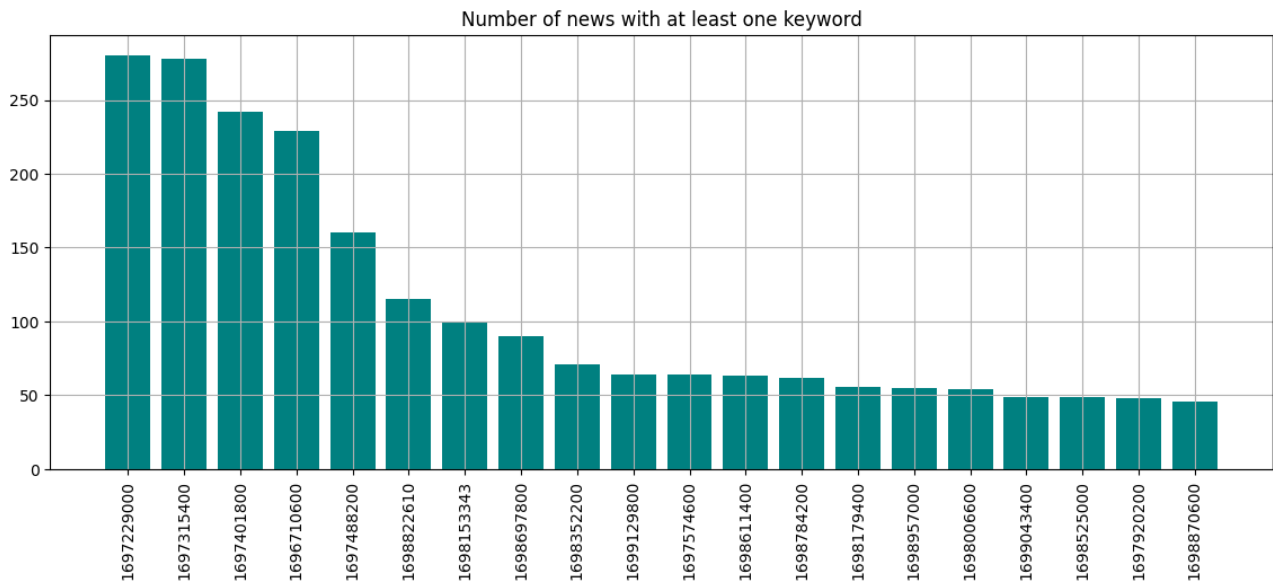


Figure 9: Number of news per each day with at least one keyword

As evident by comparing the unfiltered and filtered date volumes, the peaks in overall news volume correlate strongly to increases in semantically relevant articles containing trending terms. The strong correlation between the unfiltered and filtered date volumes implies the word clouds effectively identified hot topics and keywords that saw greater prevalence during major Persian language news cycles. This determines that the high news volume on certain dates was largely driven by significant coverage of trending subjects captured in the word cloud visualization.

4 A-Priori Algorithm

The next analysis sought to uncover sets of words that commonly co-occur together across the news stories using association rule learning. Specifically, the goal was finding pertinent 3-word combinations that appear together frequently.

To accomplish this, the A-priori algorithm was implemented in PySpark. The A-priori technique works by first identifying individual items that appear frequently, then combining them iteratively into itemsets, pruning combinations that fall below minimum thresholds.

First, the code calculated individual word frequencies to determine popular terms. The 20 most common words were selected as candidates to build word pairs from.

Then to generate 3-element tuples, the frequent word pairs were combined with the top single words in a constrained manner. The *A-priori-3tuples* function ensured new triplets met the Apriori monotonic condition - all 2-word subsets must also be frequent. Applying ranking and filtering resulted in pertinent 3-word associations like:

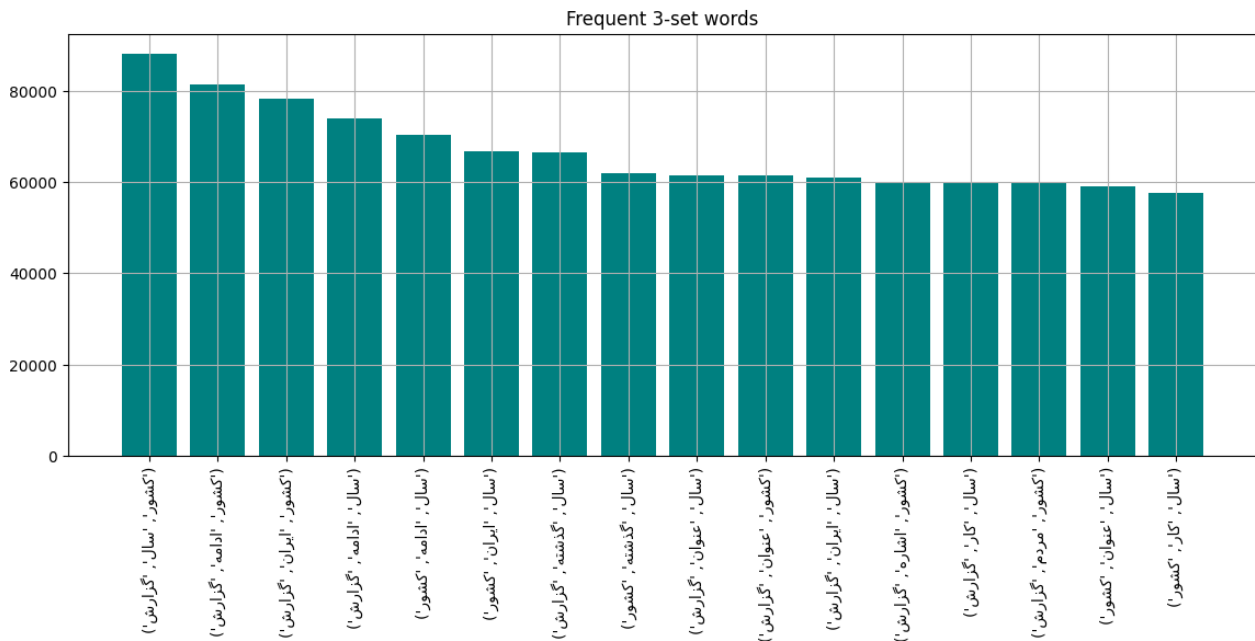


Figure 11: Frequent 3-set words

The output demonstrates meaningful multi-word terms around linked concepts.

4.1 TF-IDF

To evaluate the relevance of the identified 3-word phrase associations, a TF-IDF analysis was conducted.

TF-IDF stands for Term Frequency-Inverse Document Frequency, and is often used in text mining to score the importance of words/phrases while balancing their frequency and specificity. The TF component measures how often a term occurs in an article, while the IDF adjusts for common words across the corpus.

The scoring is defined by:

$$\text{score} = f_{t,d} / \text{len}(d) \times \log\left(1 + \frac{N}{n_t}\right)$$

This equation was applied to score documents containing the top 3-term phrases from Apriori. A threshold was set to filter for news stories strongly featuring the target phrase words according to TF-IDF. As an example, the following relevant articles were extracted with $\text{TF-IDF} > 0.06$:

[illegible]

```
{
  'body': 'شماره ۱۴۰۲/۶/۲۸/۱۸۲۹۷۲۷۴_۵۴۴.jpg',
  'image_title_url': 'https://cdn.vjc.ir/files/afnews/1402/6/28/18297۲۷۴_۵۴۴.jpg',
  'language': 'fa',
  'source': 'خبرگزاری شبکه خبرنگاران | اولین کانال ویدیویی ایران و جهان | VJC',
  'date_published': '169513۲001',
  'id': '1b463212۲۲۲adfa8a9c49f6ea',
  'url': 'http://www.vjc.ir/afnews/855۱872',
  'crawler_timestamp': '1695132128',
  'ingestor_timestamp': '1695146666',
  'summary': '',
  'hostname': '',
  'parser_categories': [],
  'keywords': ['شماره ۱۴۰۲/۶/۲۸/۱۸۲۹۷۲۷۴', 'تلفاز'],
  'parser_keyword': [],
  'author': ''
}
```

Figure 12: News with highest TF-IDF score for the first word

[illegible][illegible]

Figure 13: News with highest TF-IDF score for the second word

```
[{'گزارش':  
  [{'body': ['الجزیره',  
              'غزه',  
              'یرداخته‌است',  
              'قتلع',  
              'گزارش',  
              'غیرنظامی',  
              'مجروح',  
              'اسرائیل'],  
    'image_title_url': 'https://cdn.donya-e-  
eqtesad.com/thumbnail/3cJeFfKrSdm3/QHn809nsSzT8qCU7RegsN6Pbb5v74eEtBKeS0h05RaYq6720rIRuQkt7ITyzEhnm/AP23280355081673.jpg',  
    'language': 'fa',  
    'source': 'روزنامه دنیای اقتصاد',  
    'title': 'علم در شهر الزهراء غزه+ فیلم 200ع گزارشی تکان دهنده از قتل',  
    'date_published': 1698190443,  
    'uid': '0283ab8596a77006eff697219',  
    'url': 'https://donya-e-eqtesad.com/بخش-%D8%B3%D8%A7%D8%8C%D8%AA-%D8%AE%D9%88%D8%A7%D9%86-62/4013982-  
%DA%AF%D8%B2%D8%A7%D8%B1%D8%B4%D8%8C-%D8%AA%DA%A9%D8%A7%D9%86-%D8%AF%D9%87%D9%86%D8%AF%D9%87-%D8%A7%D8%B2-  
%D9%82%D8%AA%D9%84-%D8%B9%D8%A7%D9%85-%D8%AF%D8%B1-%D8%B4%D9%87%D8%B1-%D8%A7%D9%84%D8%B2%D9%87%D8%B1%D8%A7%D8%8C-  
%D8%BA%D8%B2%D9%87-%D9%81%D8%8C%D9%84%D9%85',  
    'crawler_timestamp': 1698190779,  
    'ingestor_timestamp': 1698206973,  
    'summary': 'الجزیره در گزارشی به قتل علم و مجروح شدن صدها غیرنظامی در غزه توسط اسرائیل یرداخته است'.}]
```

Figure 14: News with highest TF-IDF score for the third word

As it can be seen, some of the selected news stories for different phrase words are related, while others feature the terms in different contexts. This can be attributed to the varying semantic usage and polysemy of the high-level keywords.

However, an interesting observation is that due to the IDF component penalizing frequency in more documents, the TF-IDF model tends to highlight words even in shorter news stories. So

despite lower raw term counts, major terms stand out in specific contexts.

Also as expected, the high count of a word across all documents results in low TF-IDF scores for those very frequent terms. With the same threshold, there was only one news story selected for the third word in the phrase which had highest overall frequency, while multiple news belonged to the other two less common words. This demonstrates the impact of inverse document frequency damping scores for widely used words.

This inverse document frequency effect explains why most of the top-scoring selections are relatively short articles where the key phrase words make up a significant portion of the content. To avoid overemphasizing words from smaller documents, an alternative normalization approach could involve dividing the term frequencies by the maximum word count per document before applying TF-IDF weighting. This normalized frequency would prevent article length bias and improve generalization of the selections.