



ISW creates informational merchandise for diplomatic and intelligence specialists to benefit a deeper know-how of conflicts happening across the world.

## II. METHOD

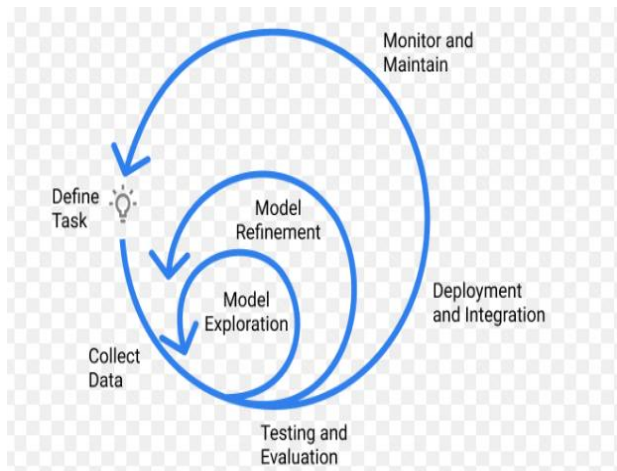


Fig 2.1

### Extracting

Hrefs We might be extracting our files from ISW's manufacturing library. First, we can scrape the 'browse' web page to get person href hyperlinks for every product. Then we keep the ones hyperlinks in a listing for our extraction features to go to later.

### Web Extraction

The first few capabilities we'll write are truthful textual content extraction. This educational isn't meant to be an instructional on using BeautifulSoup

### Get the Date

For our first feature, we can be extracting the guide date. It scans via the html record extracted from the product's webpage and reveals a discipline with the magnificence of 'submitted'. This carries our manufacturing date.

### Get the Title

Next, we need the product title. Again, this discipline is with ease classified with a category of 'title'.

### Get All the Text

Finally, we can extract the whole textual content of the record. When extract textual content, frequently comply with an 'extract-first, filter-later' fashion of net extraction. That manner that, in preliminary textual content extraction, carry out minimum filtering and processing of the textual content. favor to behavior that processing afterward in evaluation because it turns into necessary. However, in case you are extra advanced, you could need to behavior extra pre-processing of the **extracted** textual content than the beneath feature demonstrates. For get\_contents feature, caught to the naked bones — indexed some html dad and mom in a blacklist, for textual content that don't need to be extracted. Then extract all the textual content from the web page and append it right into a brief string, which in flip is appended into the listing content\_text.

## III. NATURAL LANGUAGE PROCESSING

Next, we can discern out what nations are referenced withinside the product. There are many APIs that would be utilized in checking textual content material for nations, however right here we can use a easy method: a listing of all of the nation's withinside the world. This listing is derived from Wikipedia. After the feature diagnosed all\_mentioned\_countries withinside the file, it makes use of simple statistical evaluation to pick out which nations are featured maximum prominently — those nations are maximum likely to be the factor of consciousness for the file's narrative. To do this, the feature counts the quantity of instances a rustic is noted during the file after which reveals nations noted extra instances than the average. These nations are then appended to a key country listing.

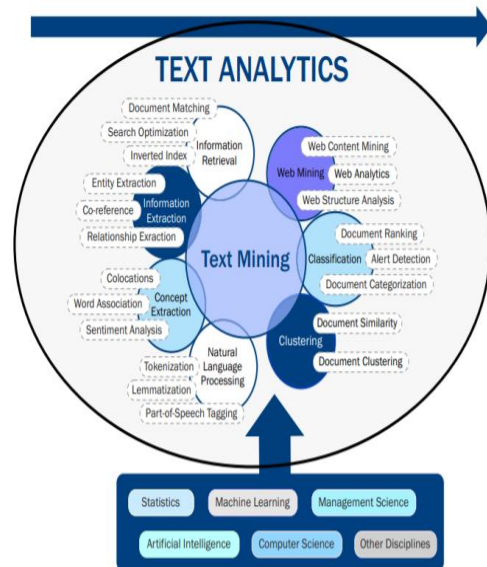


Fig 3.1

## IV. KEYWORD EXTRACTION:

**Term Frequency-Inverse Document Frequency** Our subsequent assignment is to extract key phrases from the text. The maximum not unusual place approach of doing that is via way of means of the usage of a way known as Term Frequency-Inverse Document Frequency (TF-IDF). Basically, TF-IDF fashions degree how frequently a period or phrase turned into used inside an unmarried report, then compares that to its common utilization for the duration of the complete corpus of documents. If a period is used often in an unmarried report, and seldom throughout the complete corpus of documents, then it's miles probable that time represents a keyword particular to that precise report. This article isn't always supposed to be a complete review of TF-IDF fashions. First, our characteristic will create what's typically called a 'bag-of-phrases'. This will music each phrase utilized in each report. Then, it's going to rely on each utilization of each phrase in every report — the time frequency. Then, it takes the not unusual place logarithm of each sentence in each report containing the time — the inverse report frequency.



Those values are then written to coordinates in a matrix, that's then looked after to assist us to locate the phrase's maximum probable to symbolize key phrases for our report.

$$tfidf(t, d, D) = tf(t, d) \times idf(t, D)$$

	Doc 1	Doc 2	...	Doc n
Term(s) 1	12	2	...	1
Term(s) 2	0	1	...	0
...	...	...	...	...
Term(s) n	0	6	...	3

Fig 4.1

$$IDF_i = \log \left( 1 + \frac{N_D}{f_i} \right)$$

Inverse Document Frequency for the search term  $i$  within the corpus of documents  
 The number of documents in the corpus of documents that contain the term  $D$   
 The number of documents that contain the search term

Fig 4.2

## V. TOPIC MODELING

One of the maxima not unusual place obligations in NLP is called subject matter modeling. This is a shape of clustering that tries to mechanically type files into classes primarily based totally on their textual content. In this instance, would really like to recognize at-a-look what subjects ISW is covering. By sorting files into classes primarily based totally on textual content, can without difficulty benefit an at-a-look knowledge of the document's foremost ideas.

### Topic Modeling Pipeline

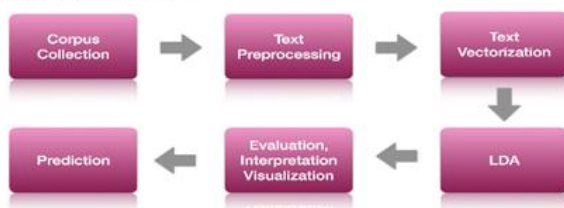


Fig 5.1

### Vectorization

For this situation, may be the use of a ok-way clustering set of rules to behavior subject matter modeling. First, will use a TF-IDF set of rules once more to vectorize every document. Vectorization is a device-mastering time that refers to the transformation of non-numeric facts into numeric spatial facts that the pc can use to behavior device mastering obligations.

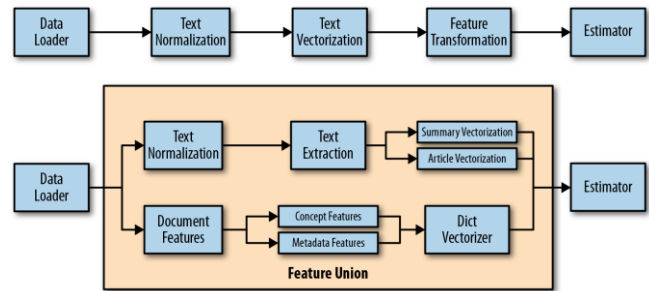


Fig 5.2

### Optimization

Once files are vectorized, helper capabilities test to look what the most excellent variety of clusters are. (The ok in ok way). In this case, the most excellent variety changed into 50. Once determined the most excellent variety, in this situation commented out that line of code and manually adjusted the parameters to identical 50. That is due to the fact the dataset reading does now no longer extrude often, so can assume the variety of most excellent clusters to live the equal over time. For facts that adjustments extra frequently, you must go back the most excellent variety of clusters as a variable — this could assist your clustering set of rules to mechanically set its most excellent parameters. an instance of this in time-collection evaluation article.

### Clustering

Once every cluster is complete, keep the variety of every cluster (1–50) to a listing of cluster numbers and the key phrases making up every cluster to a listing of cluster keywords. These cluster key phrases may be used later to feature a name to every subject matter cluster.

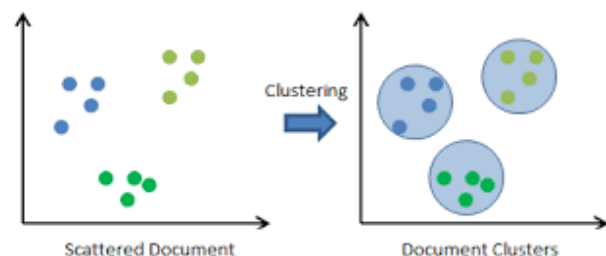


Fig 5.3

### Putting it Together

Finally, we can extract our data. Using the listing of hrefs we were given earlier, it's time to use all our extraction capabilities to the net content.

```

def get_topics(content_list):
    processed_text=[pre_process(text) for text in content_list]
    stop_words=set(stopwords.words('english'))
    cv=CountVecorizer(max_df=0.85,stop_words=stop_words)
    word_count_vector=cv.fit_transform(processed_text)

    tfidf_transformer=TfidfTransformer(smooth_idf=True,use_idf=True)
    tfidf_transformer.fit(word_count_vector)

    feature_names=cv.get_feature_names()
    vector=tfidf_transformer.transform(cv.transform(processed_text))

    #find_optimal_clusters(vector,50)

    clusters = MiniBatchKMeans(n_clusters=50, init_size=1024, batch_size=2048, random_state=20)
    for cluster in clusters:
        cluster_number.append(int(cluster))
  
```



## VI. TOPIC MODELING ENRICHMENT

Our clusters gave us a listing of phrases which are related to every cluster; however, the clusters are titled surely with numbers. This offers us the possibility to plan a phrase cloud or different thrilling visualization that could assist us apprehend every cluster, however it now no longer as beneficial for at-a-look information in a based dataset.

Additionally, consider that a few files may also fall inside a couple of subject matter classes. Multiple clustering isn't supported with the aid of using k-means, so will pick out those files manually. First, I'll print the primary few rows of key phrases to get a concept of the facts I'm dealing with. After enormous experimentation with plenty of techniques, determined on a completely easy method. scanned every listing of key phrases touching on every cluster and cited enormous key phrases in every that associated with a selected subject matter. At this stage, area information became key. Know, for example, that Aleppo in an ISW report is nearly sincerely cited in connection with the Syrian Civil War. For this example, however, the easy method

works well. After making notice of numerous enormous key phrases gift withinside the cluster lists, made some lists of very own that contained key phrases related to the very last subject matter classes desired withinside the based facts. The characteristic surely compares every cluster's listing of key phrases with the lists created, then assigned a subject call primarily based totally on suits withinside the lists. It then appends the ones very last subjects to a listing of topic categories.

$$IDF_i = \log \left( 1 + \frac{N_D}{f_i} \right)$$

Inverse Document Frequency for the search term  $i$  within the corpus of documents  
 The number of documents in the corpus of documents that contain the term  $D$   
 The number of documents that contain the search term

Fig 6.1

Topic Modelling Enrichment

```
[ 'ret,islamic,political,province,military,groups,security,us,state,states,syrian,diyala,american,isis,n
'ar,raqqa,eastern,northern,al,campaign,city,province,update,forces,iraq,may,syrian,ez,zour,regime,deir
'haq,ahl,ib,gains,asa,january,activity,deir,control,ez,zour,isis,iraq,november,update,report,situation
'eu,political,president,bugayova,intsum,sanctions,nataliya,nato,west,elections,review,africa,informati
'damascus,isis,jabhat,military,homs,russian,province,pro,russia,groups,forces,rebel,al,city,assad,oppo
```

SOME OF THE KEYWORDS ASSOCIATED WITH EACH TOPIC CLUSTER. WE'LL USE THESE KEYWORDS TO SORT CLUSTERS INTO PRE-DEFINED CATEGORIES

```
oir=['OIR Iraq','yezidis','mosul','peshmerga','isis','iraq','sinjar','baghdad','maliki',
'daquq','anbar','isf','abadi','malaki','ramadi','iraqi','fallujah','dabiq']
```

```
terrorism=['Terrorism','jihadi','islamic','salafi','qaeda',
'caliphate','isis','terrorist','terrorism']
```

```
syrian_conflict=['Syrian Conflict','sana','syria','assad',
'idlib','aftrin','aleppo']
```

```
russia=['Russia','russia','belarus','slavic','kremlin','russian',
'minsk','ukraine','putin']
```

```
iran=['Iran','iran','iranian','proxy','militias','militia','marjah']
```

```
turkey=['Turkey','erdogan','turkish','turkey']
```

## VII. DATABASE CREATION

The final step is to carry collectively all our extracted statistics. For these statistics, select the JSON layout. This is due to the fact desired to shape sure sorts of statistics differently — for example, the places subject will consist of a listing of dictionaries of vicinity names, latitudes, and longitudes. In JSON layout is the only manner to save such formatted statistics to a neighborhood disk



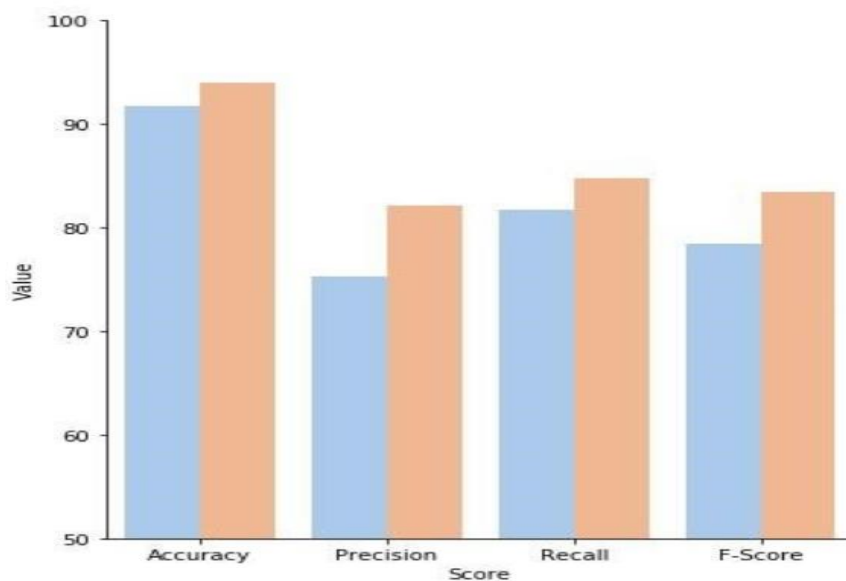
```
#Add all of the defined lists from our functions
#To the new storage list

doc={
    '_id': len(hrefs) - i,
    'title': titles[i],
    'date': dates[i],
    'places': coord_list[i],
    'people': people[i],
    'keywords': keywords[i],
    'countries': countries,
    'full text': content_text[i],
    'url': links[i],
    'topic cluster': cluster_number[i],
    'categories': topic_categories[i]
}
```

## VIII. RESULT

```
[{'_id': 1745,
  'categories': ['Russia'],
  'countries': {'all mentioned countries': ['Russia',
    'Belarus',
    'Azerbaijan',
    'Kyrgyzstan',
    'Serbia'],
    'focus area': ['Russia', 'Belarus']}},
  'date': '2020-10-20T15:07:13-04:00',
  'full text': 'Russia's Unprecedentedly Expansive Military Exercises in Fall 2020 Seek to Recreate Soviet-Style Multinational Army | Institute for the Study of War',
  'keywords': ['exercises',
    'kremlin',
    'russian',
    'fsu',
    'brotherhood',
    'deployments',
    'belarus',
    'dot',
    'unbreakable',
    'wmd'],
  'people': ['Vladimir Putin'],
  'places': [{'Lat': 64.6063136, 'Location': 'Russia', 'Lon': 97.7453061},
    {'Lat': 55.7514952, 'Location': 'Kremlin', 'Lon': 37.6181531},
    {'Lat': 55.4805995, 'Location': 'Soviet Union', 'Lon': 28.773013},
    {'Lat': 53.4250605, 'Location': 'Belarus', 'Lon': 27.6971358},
    {'Lat': 55.7504461, 'Location': 'Moscow', 'Lon': 37.6174943},
    {'Lat': 40.3936294, 'Location': 'Azerbaijan', 'Lon': 47.7872508},
    {'Lat': 41.5089324, 'Location': 'Kyrgyzstan', 'Lon': 74.724091},
    {'Lat': 44.0243229, 'Location': 'Serbia', 'Lon': 21.0765743},
    {'Lat': 52.2809939, 'Location': 'bratstvo-2020', 'Lon': 43.6981349}],
  'title': 'Russia's Unprecedentedly Expansive Military Exercises in Fall 2020 Seek to Recreate Soviet-Style Multinational Army',
  'topic cluster': 47,
  'url': 'http://www.understandingwar.org/backgrounder/russia%E2%80%99s-unprecedentedly-expansive-military-exercises-fall-2020-seek-recreate-soviet'}
```

Result Fig 8.1



Result Fig 8.2

## IX. CONCLUSION

Through this paper, we've got tested an efficient and correct structured-records extraction from websites. Such correct records extraction is viable through the usage of NLP strategies with the utility of Machine Learning (ML) and Deep Learning (DL) fashions. With the reference of textual files like activity emptiness records and CV, that are frequently used in applicant monitoring structures, we've got evolved this gadget and examined it with files of comparable type. The assessment metrics received are better and the execution time for such records extraction is likewise better. Such highly optimized and correct records extraction structures may be utilized in various other fields like studies publications, activity portals etc. However, the records used for education ML and DL fashions must be tweaked in keeping with corresponding applications of the gadget. Although the tactics that we've got located now have solved our issue, we might should ex trade it withinside the close to future in case the requirement modifications or new sorts of records are received. There are nevertheless a few strategies that are but to be examined inclusive of embedding the usage of the BERT model, attempting record similarity instead of token similarity and so on. Also, withinside the subsequent stage, we can be focusing on extracting records from unstructured files of different domain names.

## REFERENCES

1. <https://morioh.com/p/4fdf57959638>
2. <https://towardsdatascience.com/something-from-nothing-use-nlp-and-ml-to-extract-and-structure-web-data-3f49b2f72b13>
3. <https://www.dataquest.io/blog/nlp-project-scraping-the-web-to-gather-data/>
4. <http://www.webextractor.com/>
5. <https://medium.com/geekculture/4-web-scraping-projects-that-will-help-automate-your-life-6c6d43aefeb5>
6. <https://www.ibm.com/cloud/learn/text-mining#:~:text=Text%20mining%2C%20also%20known%20as,meaningful%20patterns%20and%20new%20insights>
7. <https://www.geeksforgeeks.org/data-mining-techniques/>
8. <https://seleritysas.com/blog/2019/06/14/five-key-methods-for-big-data-optimisation/>

## AUTHORS PROFILE



**Gaurav Sharma**, Research & Development Software Engineering at OnePlus India. Graduated from Birla Institute of technologies. My Field of Research Includes NLP, Image Processing, Deep learning, Machine Learning, NFC, AI Algorithms. I am also a content creator at Udemy. By implementing my ideas, I have also won National, Regional, and international level Hackathons in the same Domain. my experience in this field includes an internship as a Machine Learning Internship. Being an AI enthusiast, I always aim to pursue masters in the coming year. I am enthusiastic software developer with love for design patterns, algorithms and, I also have an avid interest in math as it keeps my head running. Passionate about Coding. Senior Software Engineer at OnePlus with excellent problem-solving skills and ability to perform well in a team with a Bachelor of Engineering focused in Computer Science. LinkedIn: <https://www.linkedin.com/in/gaurav-sharma21/>