

A Big Data Framework for Sentiment Analysis in Healthcare: the case of Covid-19

Par Ulrich TCHUENKAM
Sous la direction de Wajdi DHIFLI

OUTLINE

- MOTIVATION AND CONTEXT
- METHODS AND MATERIAL
 - Layers of the framework
 - Data management in the framework
- CASE STUDY ON COVID-19
- RESULTS
- DISCUSSION
- CONCLUSION

Huge amount of data
on healthcare
subjects

Feelings and thoughts
classification models

Scalable, distributed
and parallelizable
technologies

Topics generation
models

Huge amount of data
on healthcare
subjects

Big data

Feelings and thoughts
classification models

Sentiment analyzers

Scalable, distributed
and parallelizable
technologies

Big data technologies

Topics generation
models

Topic modelers

OBJECTIVE

Big data
on healthcare subjects

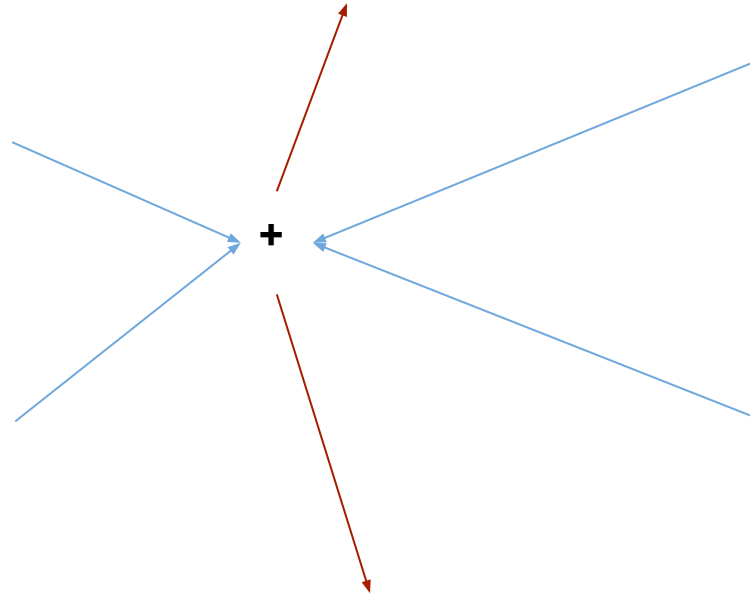
Big data technologies

Sentiment analysis

Sentiment analyzers

Topics modelers

Topic modeling



Layers architecture

Visualization Layer

Dashboard, data statistics

Big data processing and Analytics Layer

Batch processing, sentiment analysis, topic modeling

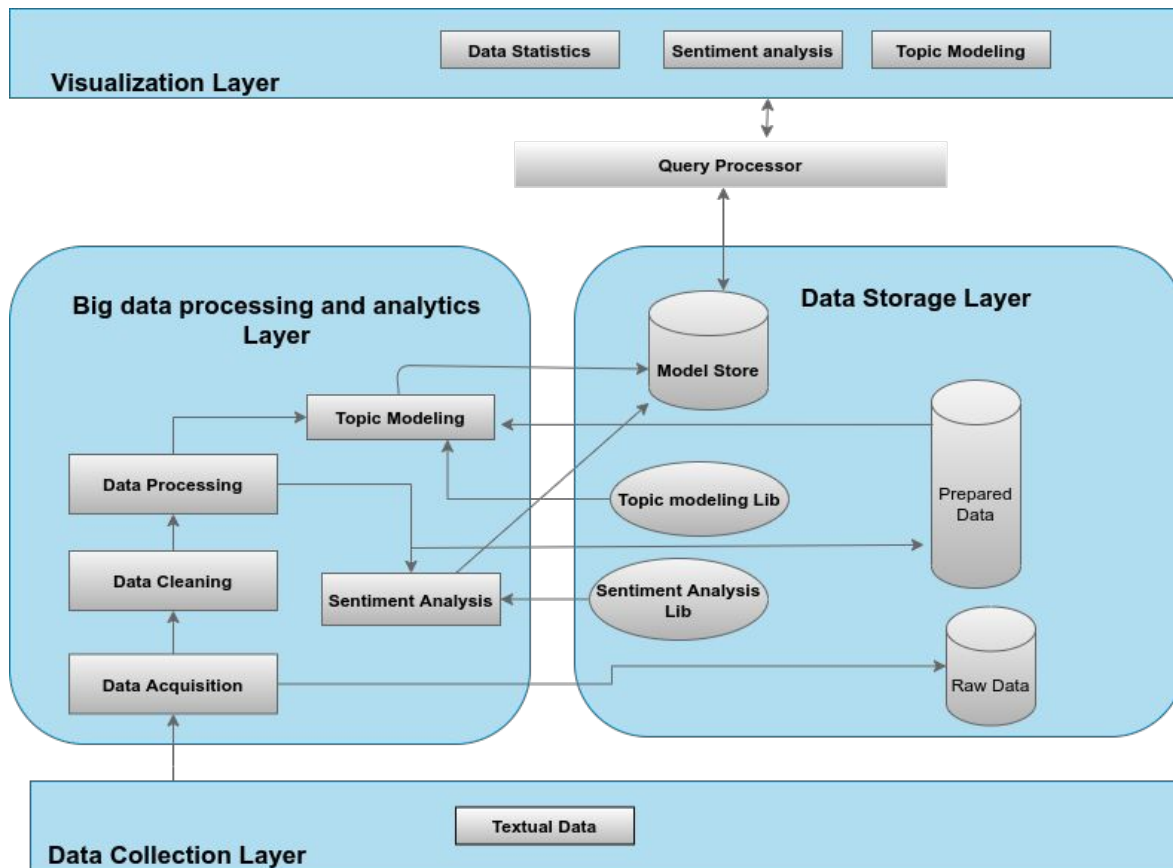
Data Storage Layer

Relational databases, JSON documents

Data Collection Layer

Textual data (Social media, news articles, RSS)

METHODS AND MATERIAL



*Data management in
the framework*

Data acquisition

- Using Twitter's API
extract tweets about
covid-19.
+
- A public dataset
containing tweets about
corona.

Data storage

- PostgreSQL
- Elasticsearch

Data description

- *Size* : 633549 tweets
- *Language* : English
- *Period* : April 1st to April 14th
- *Topic* : CORONA
- *Variables* :

<i>Text</i>	<i>Location</i>	<i>Date</i>
-------------	-----------------	-------------

Data wrangling

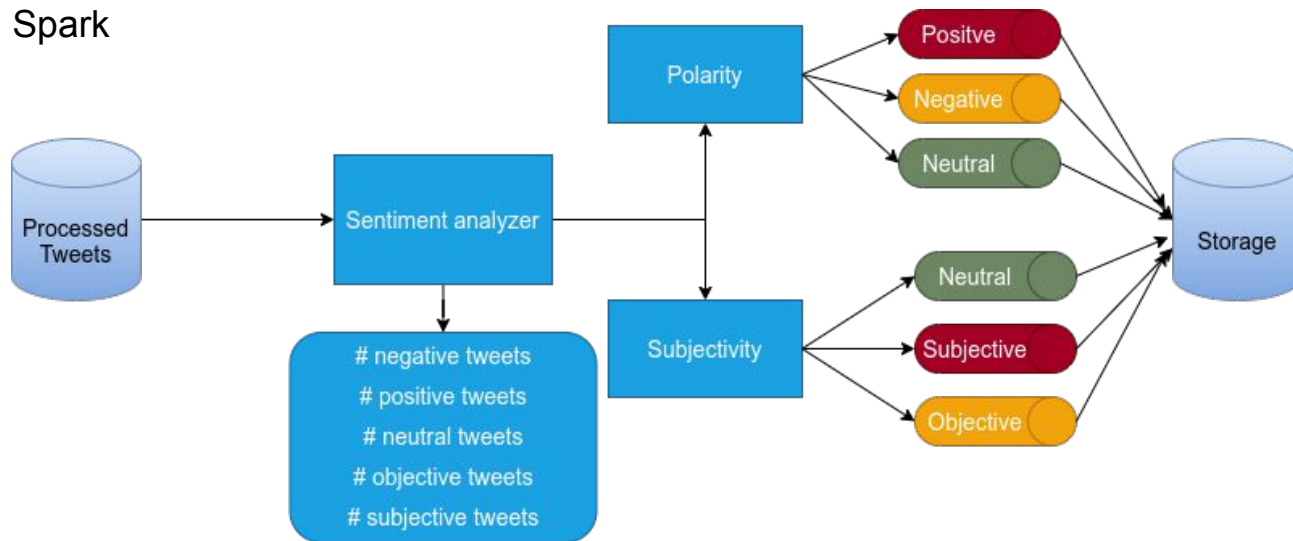
- Removal of non-alphabetical characters
- Removal of duplicates
- Removal of emoticons
- Removal of punctuation marks

Data processing

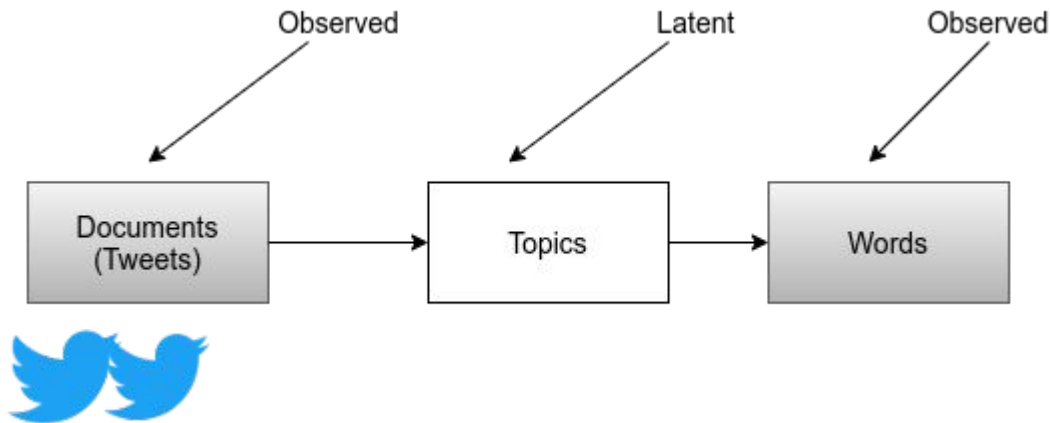
- Tokenization
- Removal of stopwords
- Converting tokens to lowercase characters
- Creation of bigrams
- Stemming
- Lemmatization
- Creating corpus and dictionary.

Sentiment analysis

Packages : TextBlob, Spark



Topic Modeling



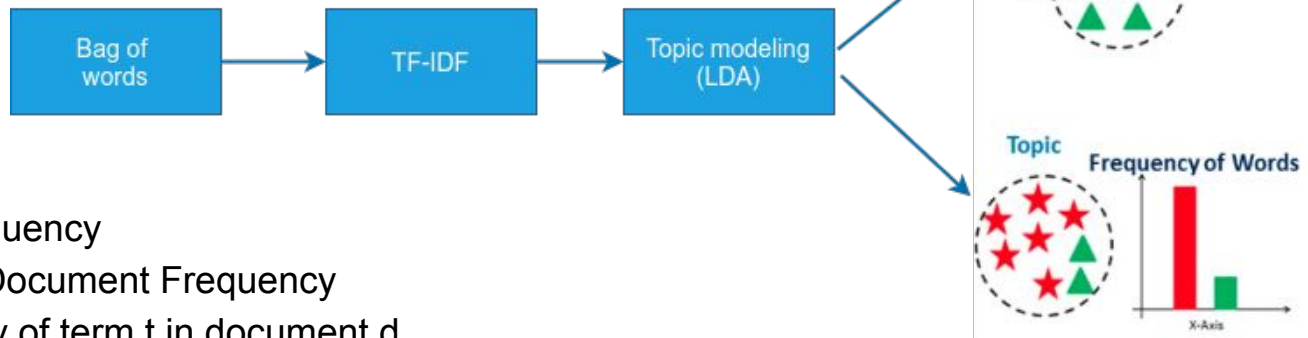
- Documents are about several topics at the same time.
- Topics are associated with different words.
- Topics in the documents are expressed through the words used.

Topic Modeling

$$TF(t) = \frac{f_{t,d}}{T_d} \quad , \quad IDF(t) = \log\left(\frac{N}{D_t}\right)$$

Package : Gensim

$$TF-IDF(t) = TF(t) * IDF(t)$$



Notations:

- TF = Term Frequency
- IDF = Inverse Document Frequency
- $f_{t,d}$ = Frequency of term t in document d ,
- T_d = Number of terms in document d ,
- N = Total number of documents,
- D_t = Number of documents with term t

Evaluation metric : Topic coherence

$$\text{Coherence} = \sum_{i < j} \text{score}(w_i, w_j)$$

$$\text{score}_{\text{UMass}}(w_i, w_j) = \log \frac{D(w_i, w_j) + 1}{D(w_i)} \quad \text{for } w_1 \dots w_k \text{ k top words}$$

Notations

- $D(w_i)$: count of documents containing the word w_i
- $D(w_i, w_j)$: count of documents containing both words w_i and w_j
- D : The total number of documents in the corpus.

Data Summary tab

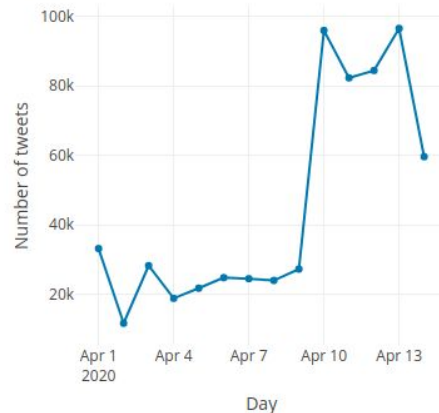
Package : Dash

Drag and Drop or [Select dataset file](#)

Preprocessing april_1_14_2020.csv completed and now using it

Total Number of Tweets: 633549

Size of dataframe: 324.22MB



Topic modeling tab (Model 1)

How many topics to model?



How many top terms per topic?



Slide to select the document-topic density (ALPHA)

ALPHA = 0.001



Slide to select the topic-word density (ETA)

ETA = 0.001



Slide to select minimum Probability of topic?

Topics with Probability < 0.01 shall be rejected?



Run Model(Basic)

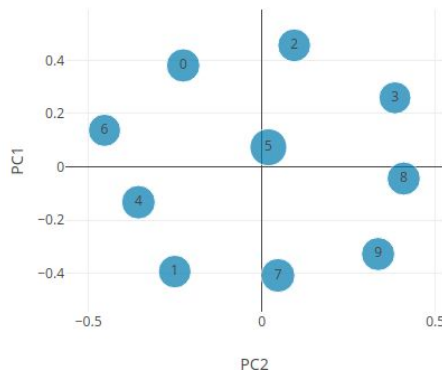
Run Model(Advance)

Metrics

Perplexity = -626.21

Coherence score = 0.71

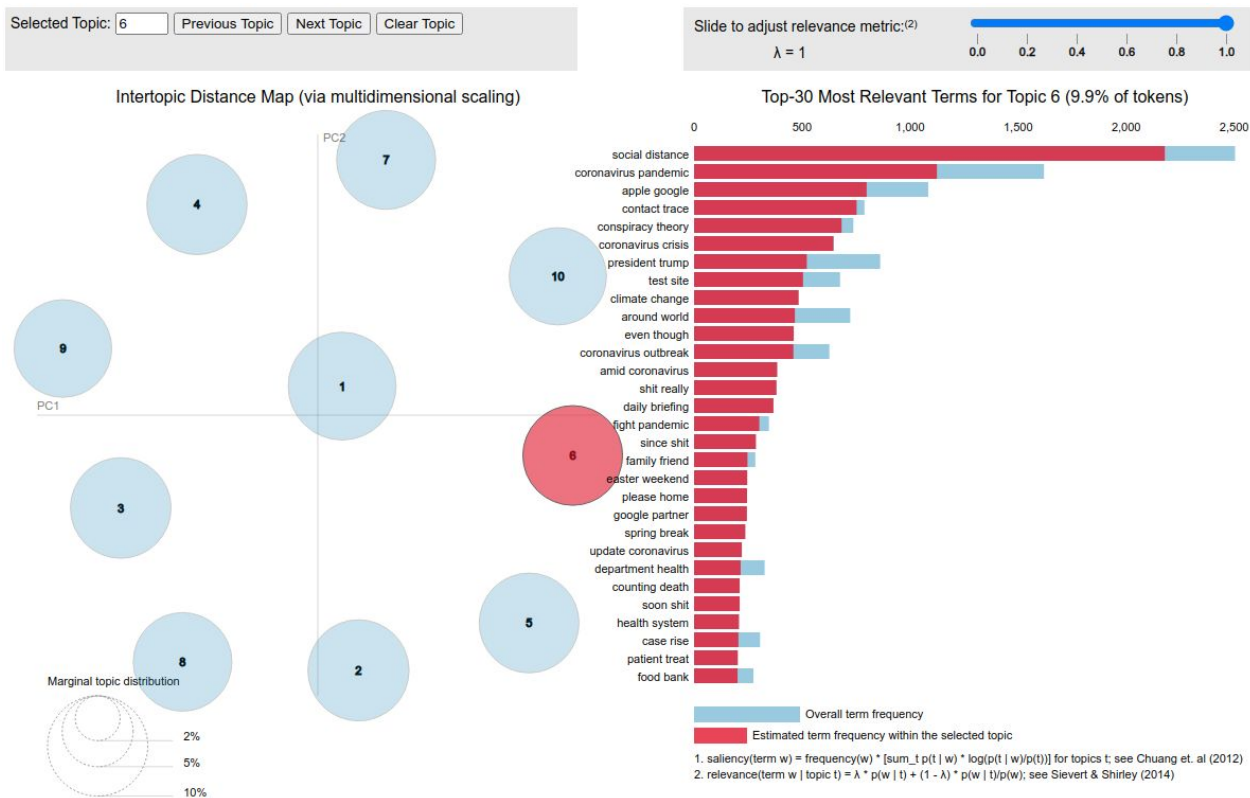
Intertopic Distance Map (via multidimensional scaling)



novel coronavirus
real estate
related death
wuhan china
care worker
raise money
care patient
track fight
health care
response pandemic

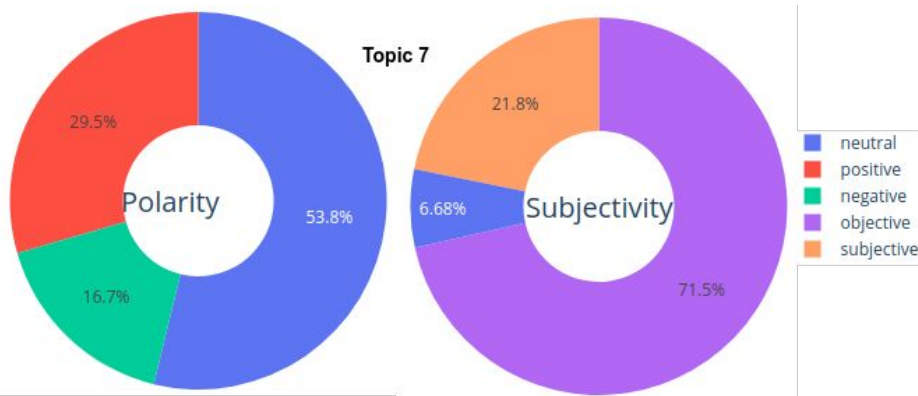
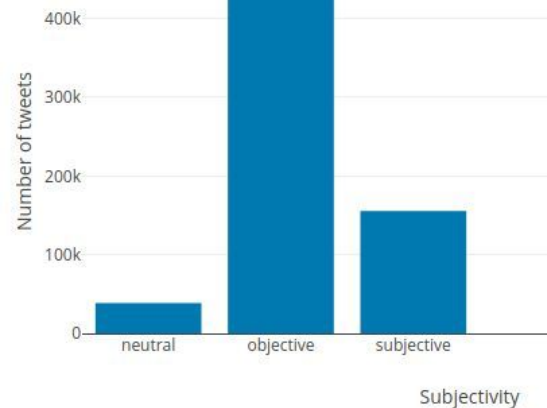
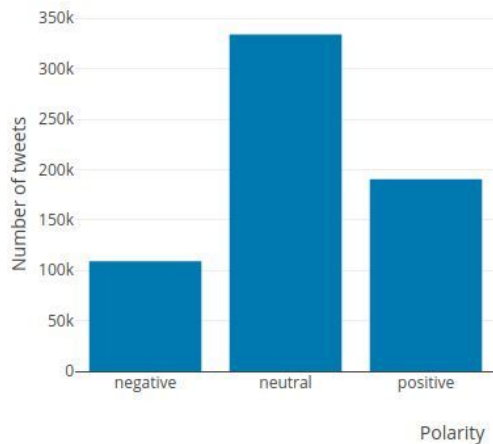
Rename Topic 0	Rename Topic 1	Rename Topic 2	Rename Topic 3
Rename Topic 4	Rename Topic 5	Rename Topic 6	Rename Topic 7
Rename Topic 8	Rename Topic 9		

Topic modeling tab (Model 2)



RESULTS

Sentiment analysis tab



DISCUSSION

- The API of Twitter provides a great way to retrieve tweets but has some limitations.
- Twitterscraper and GetOldTweets3 give access to older tweets but they only get a random number of them.

POTENTIAL FUTURE WORKS

- Leveraging Spark Streaming API.
- Add a user mapping feature to the framework.
- Add support for other corpus and languages.
- Use ontologies and knowledge graphs to enhance sentiment analysis and topic modeling performances.

CONCLUSION

- The work leverages big data tools and provides an abstracted data collection layer which can be adapted to several data sources (twitter, Facebook posts, news articles, RSS fluxes, scraping modules).
- The proposed framework is scalable, **big data** ready, and has been applied on real-world covid-19 data.
- And I have been able to use the technologies taught in class (Dash, Spark, Elasticsearch, PostgreSQL, machine learning, *etc.*) to develop a prototype of the framework.

Thank you for your attention !

REFERENCES

Twitterscraper : <https://pypi.org/project/twitterscraper/0.2.7/>

GetOldTweets3 : <https://pypi.org/project/GetOldTweets3/>

Spacy : <https://pypi.org/project/spacy/>

Spark : <https://spark.apache.org/>