

Le Monde Covid

Analyse de documents textuels grâce à l'intelligence artificielle

Topic Modeling

François LABASTIE

Décembre 2020

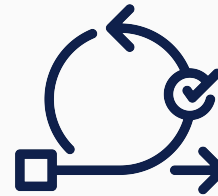
Titre professionnel "Développeur en intelligence artificielle" de niveau 6 enregistré au RNCP sous le n°34757

Passage par la voie de la formation - parcours de 19 mois achevé le 20 octobre 2020

Concept du projet & méthodologie

Topic modeling · Modélisation de thèmes
à partir de documents du média lemonde.fr

- Scraping de la page “recherche” du Monde
- Méthodologie → SCRUM & UML



DEMO

Du point de vue **utilisateur**

...

LeMondeCovid

Modeling

Topic Modeling

Recherche de topics dans un corpus

Section

planete - 1298 docs

Model

NMF (Negative Matrix Factorisation)

Nb. Top words

5

Nb. Topics

5

Nb. Documents

1298

Envoyer

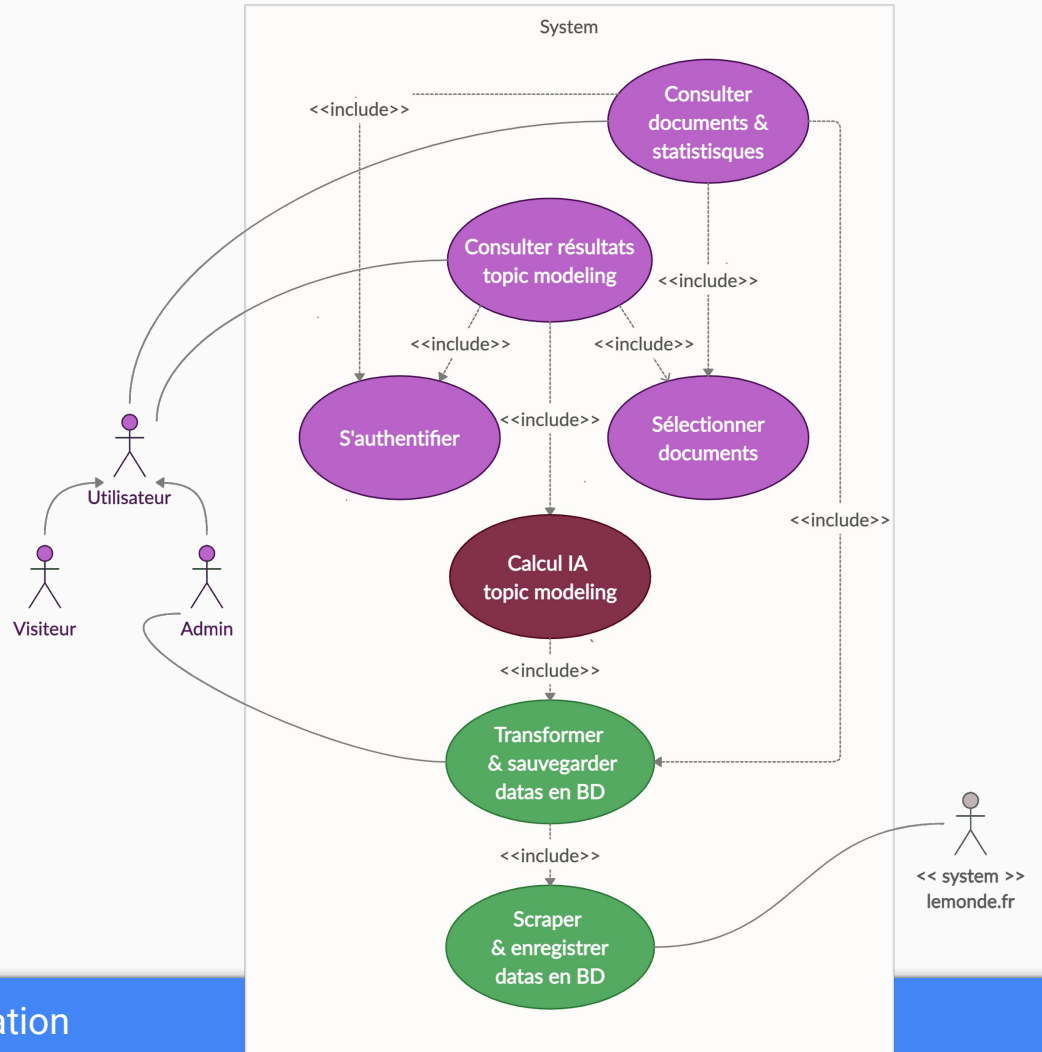
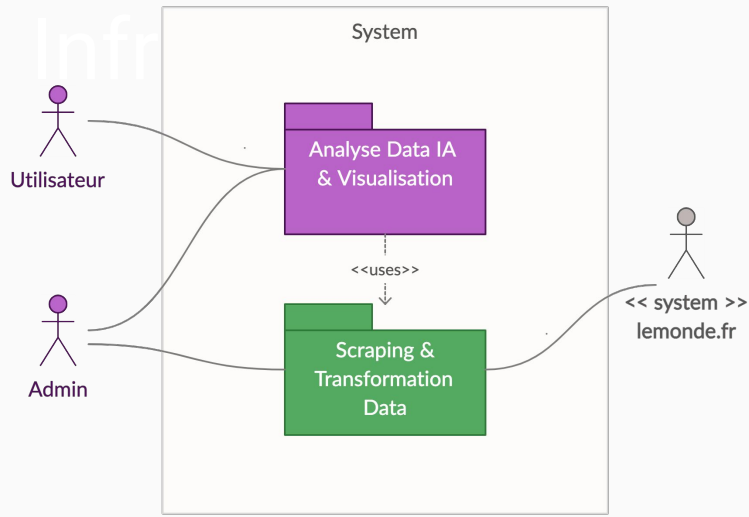
Reset

planete - NMF model - 5 top-words - 5 topics - 1298 documents

#	Top words				
1	france	santé	personnes	tests	ministre
2	pays	morts	etats	pandémie	décès
3	vaccin	vaccins	essais	pfizer	doses
4	patients	réanimation	lits	hôpital	hôpitaux
5	chine	février	chinois	wuhan	hubei

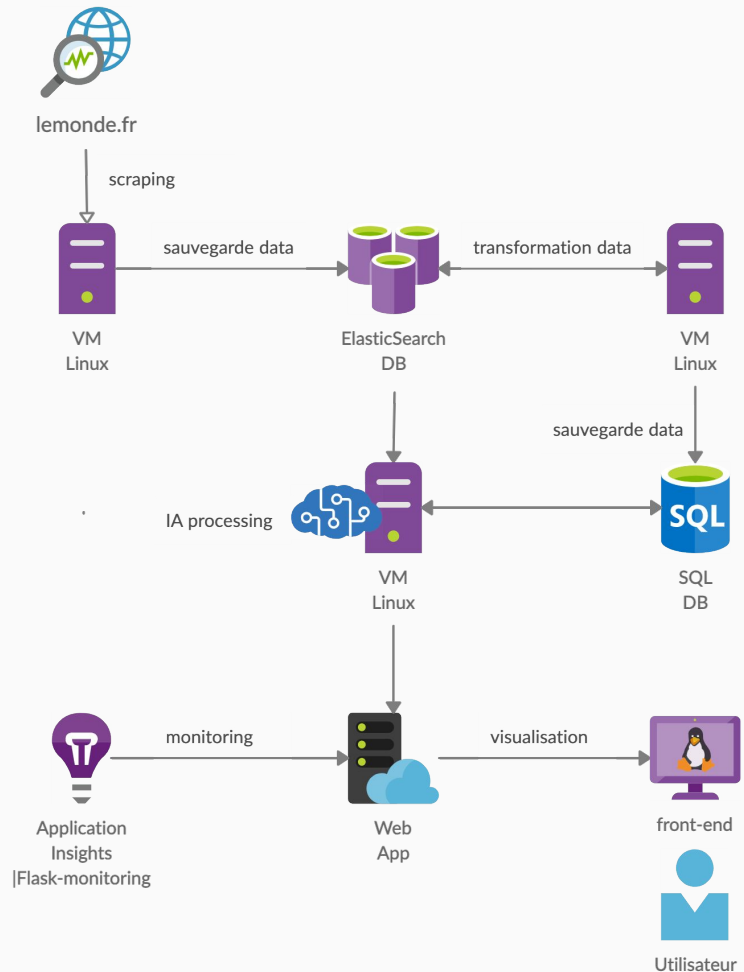
Words Selection

Recherche de documents par mots-clés



Infrastructure

- Scraping → scrapy
- Elasticsearch
- Azure SQL database
- VM sur Azure ML Studio
- Azure Web Apps
- Flask + bootstrap
- Flask-monitoring



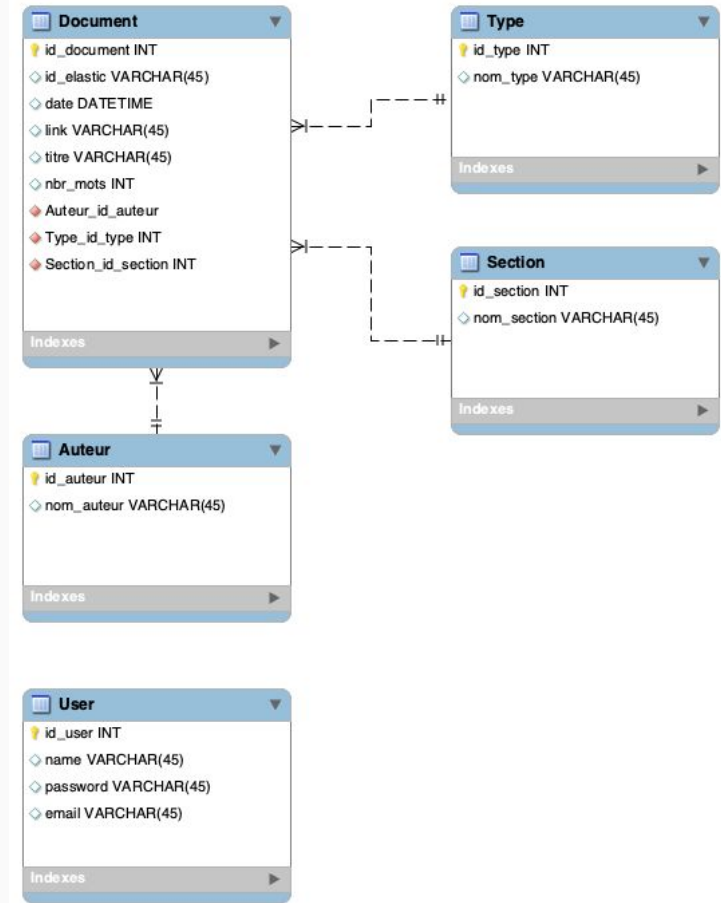
Recueil et traitement des données

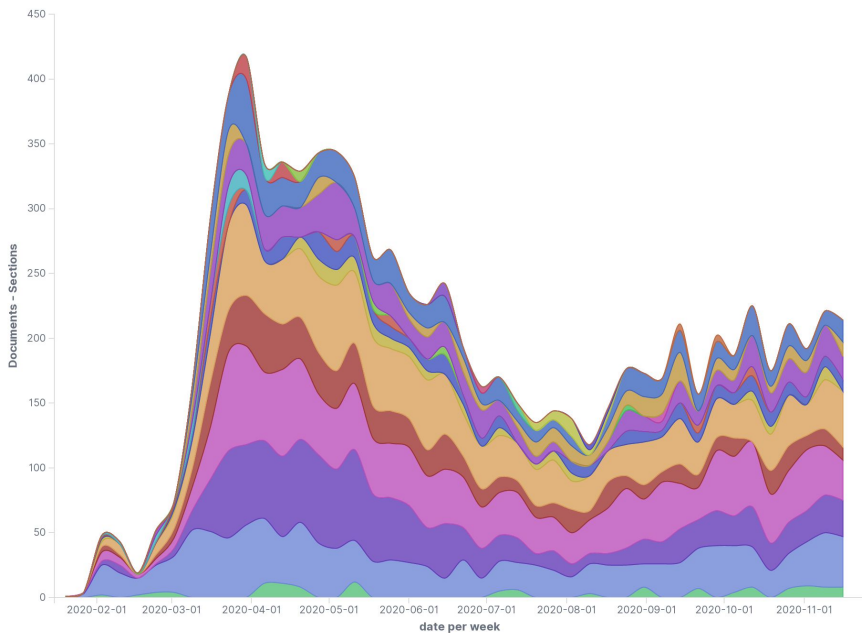
- Scraping
- Insertion data → bd [Elasticsearch](#) (tokenization, stopwords, etc.)
- Insertion data → bd Azure SQL database

- package Analyse → <https://github.com/flabastie/news-analysis>
- package Scraping → <https://github.com/flabastie/news-analysis-data>

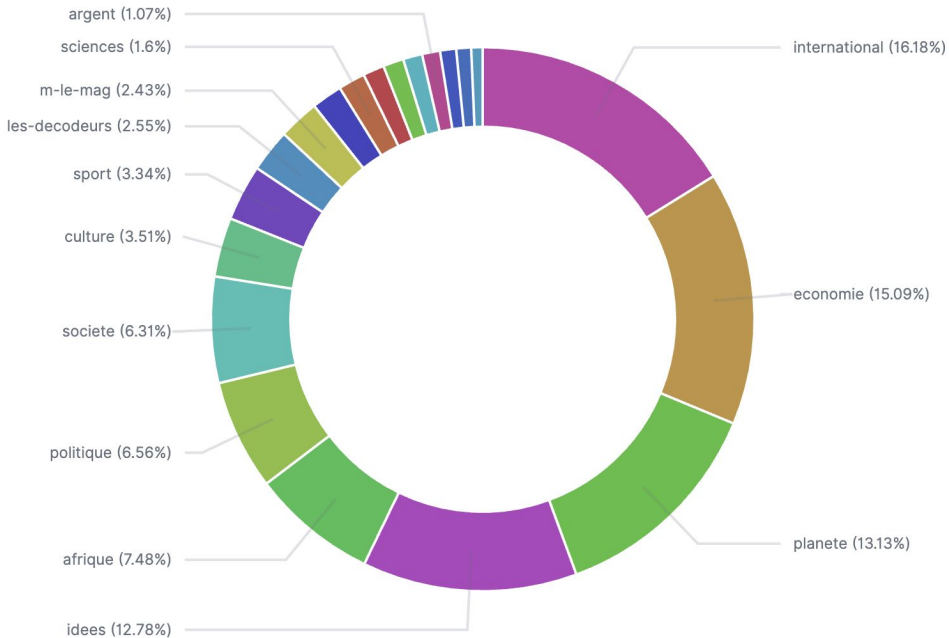
Dev tools Elastic → GET /news_analysis/_mapping

```
{
  "news_analysis" : {
    "mappings" : {
      "properties" : {
        "author" : {
          "type" : "text",
          "fields" : {
            "keyword" : {
              "type" : "keyword",
              "ignore_above" : 256
            }
          }
        },
        "content_all" : {
          "type" : "text",
          "fields" : {
            "keyword" : {
              "type" : "keyword",
              "ignore_above" : 256
            }
          }
        },
        "content_html" : {
          "type" : "text",
          "fields" : {
            "keyword" : {
              "type" : "keyword",
              "ignore_above" : 256
            }
          }
        }
      }
    }
  }
}
```





- les-decode
- planete
- idees
- international
- afrique
- economie
- m-le-mag
- culture
- pixels
- sante
- campus
- politique
- sciences
- education
- sport
- societe
- m-perso
- disparitions
- argent
- smart-cities
- livres
- m-styles
- series-d-et
- football
- emploi



Code Scraping

- scrapy_projects
- transform_projects

- [lemonde_covid/spiders/lm_covid.py](#)
- [lemonde_covid/pipelines.py](#)
- [transform_projects/LeMonde_stop_words.ipynb](#)
- [transform_projects/elastic-to-sql.py](#)

Code IA

- TopicModelingLDA
- TopicModelingNMF

- [project/main.py](#)
- [project/queries/selection.py](#)
- [project/processing/modeling.py](#)

Model Evaluation

- *Subjectivité ...*
- Topic coherence
- Coherence score

Topic coherence

Attribuent un score à chaque topic en mesurant le degré de similitude sémantique entre les mots les mieux notés du topic.

Mesures qui aident à distinguer entre les topics qui sont des sujets sémantiquement interprétables et les topics qui sont des artefacts d'inférence statistique.

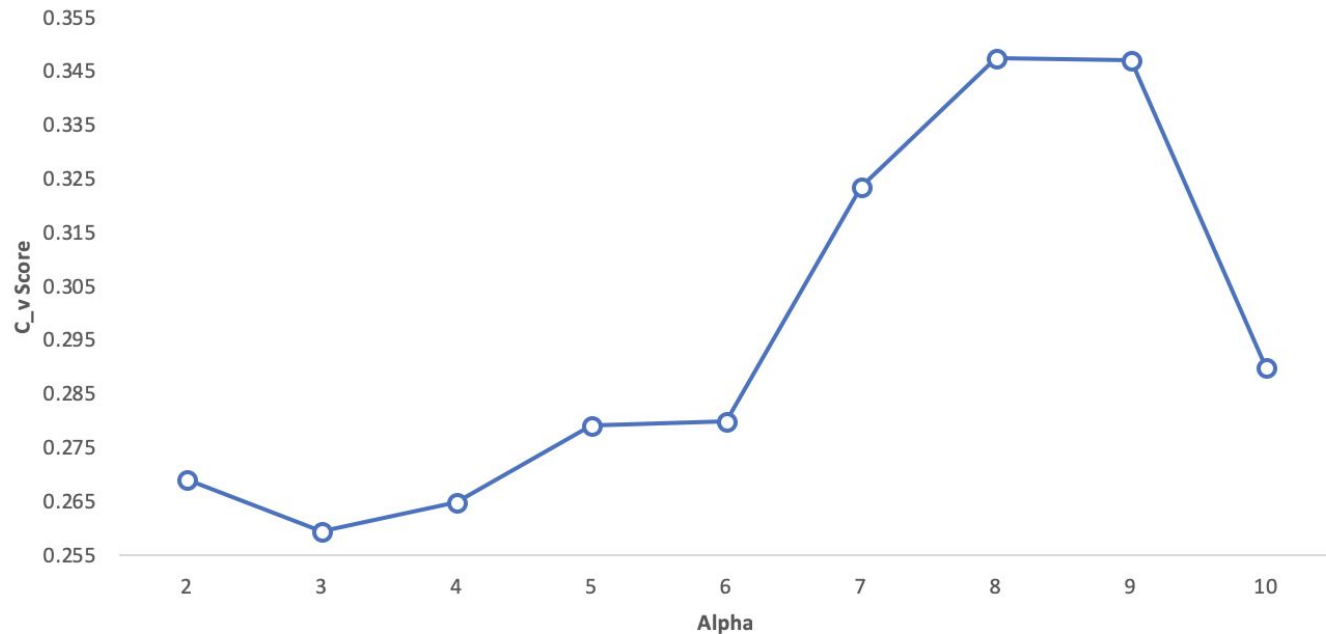
Coherence score C_v

La mesure C_v est basée sur une fenêtre glissante, une segmentation en un seul ensemble des premiers mots et une mesure de confirmation indirecte qui utilise des informations mutuelles ponctuelles normalisées (NPMI) et la similitude cosinus.

Intertopic Distance Map (via multidimensional scaling)



Topic Coherence: Determining optimal number of topics



DEMO

Du point de vue **développeur**

...

+ monitoring

LeMondeCovid

Modeling

Topic Modeling

Recherche de topics dans un corpus

Section

planete - 1298 docs

Model

NMF (Negative Matrix Factorisation)

Nb. Top words

5

Nb. Topics

5

Nb. Documents

1298

Envoyer

Reset

planete - NMF model - 5 top-words - 5 topics - 1298 documents

#	Top words				
1	france	santé	personnes	tests	ministre
2	pays	morts	etats	pandémie	décès
3	vaccin	vaccins	essais	pfizer	doses
4	patients	réanimation	lits	hôpital	hôpitaux
5	chine	février	chinois	wuhan	hubei

Words Selection

Recherche de documents par mots-clés

Merci !

