

Topic Modeling

NLP

Dat'Art



Sommaire

- 01 Contexte
- 02 Organisation
- 03 Demo
- 04 Présentation du projet
- 05 Démarche
- 05 Axes d'amélioration



Contexte

Objectif :

- Créer une interface streamlit pour prédire le sujet d'un article scientifique**
- Déploiement d'une API sur Azure avec le modèle**
- Connexion avec une base de données SQL**



Organisation

Organisation

Board View Galerie Table 1 de plus...

Filtrer Trier ... Nouvelle page

To Do 2

API

To Do

- Rassembler les topics
- Pipeline
- faire l'algo pour prendre les topics
- faire un transformer en list
- API local du modèle et de la transformation pour la classif

A faire Naïs

To Do

+ Nouvelle page

Doing 4

Ahmed

Streamlit

Doing

Naïs

Prez

Doing

Connexion BDD

Doing

Camille

Done 🎉 1

Tous

EDA

Done 🎉

+ Nouvelle page

+



Demo



Présentation du projet

Base de données



API (contient le modèle)



Streamlit (interface)

En parallèle :

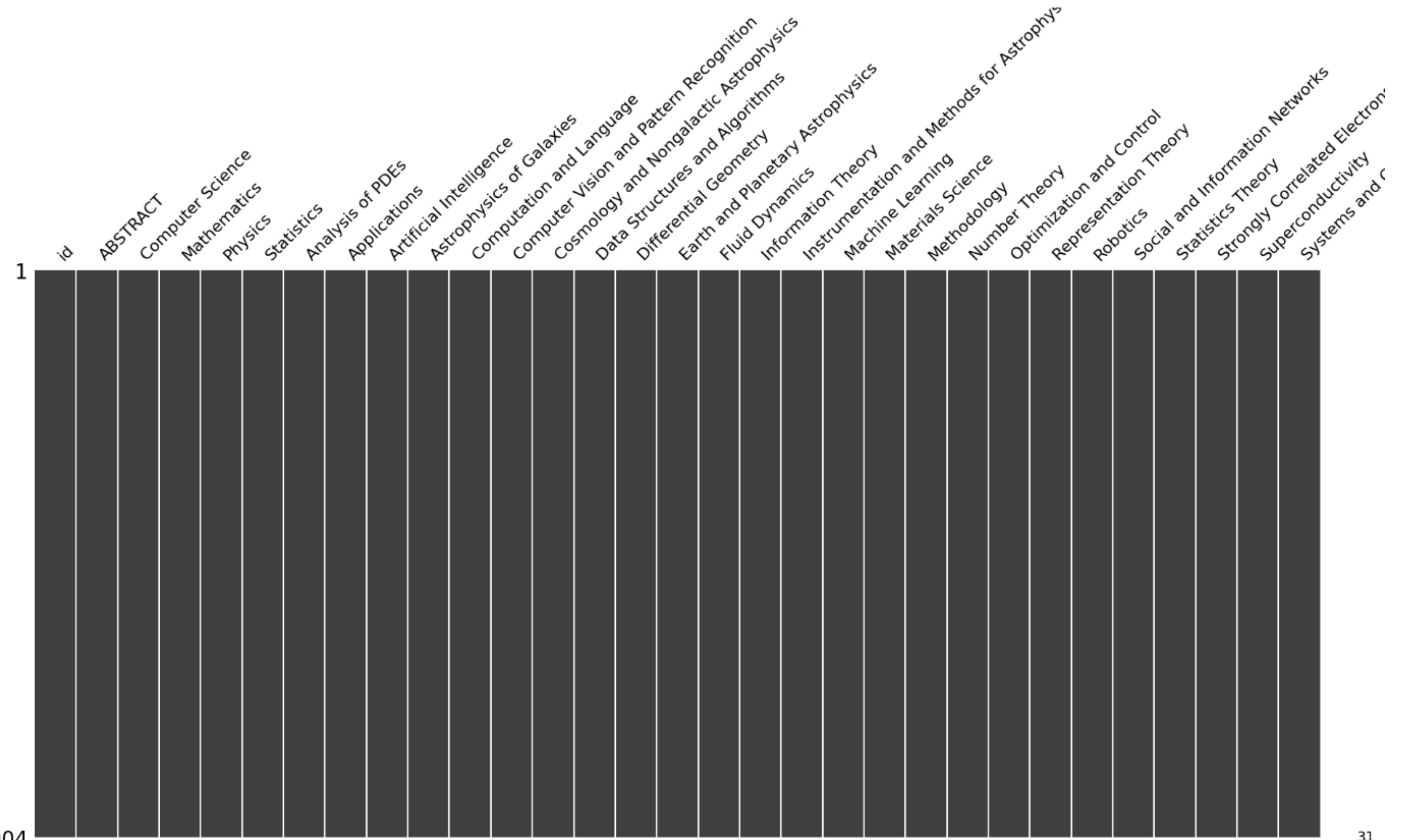
Modèle de classification sur 29 topics

API pour aller requêter le modèle



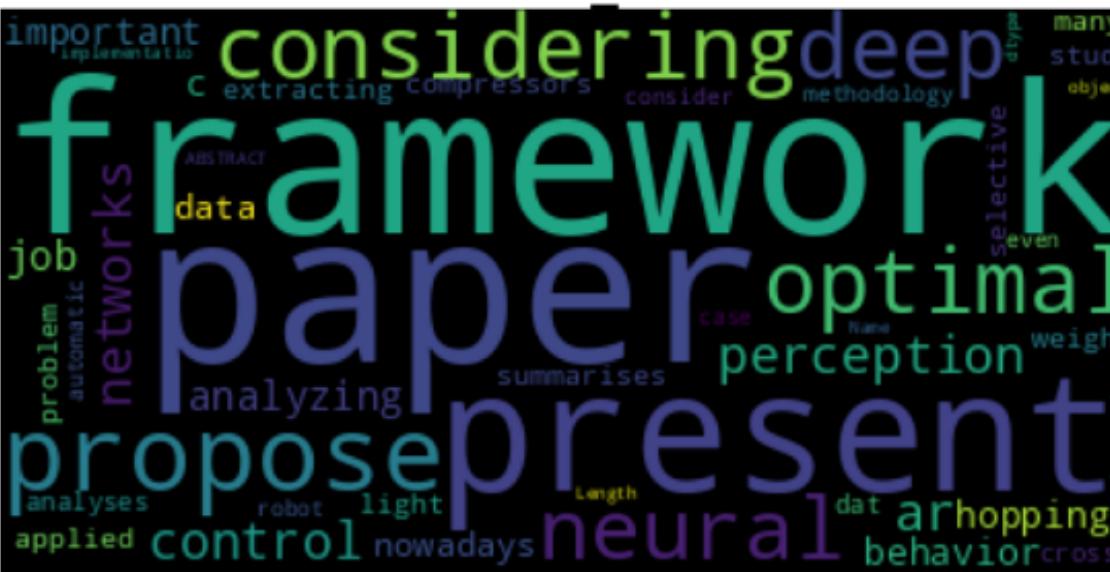
Démarche

Les données

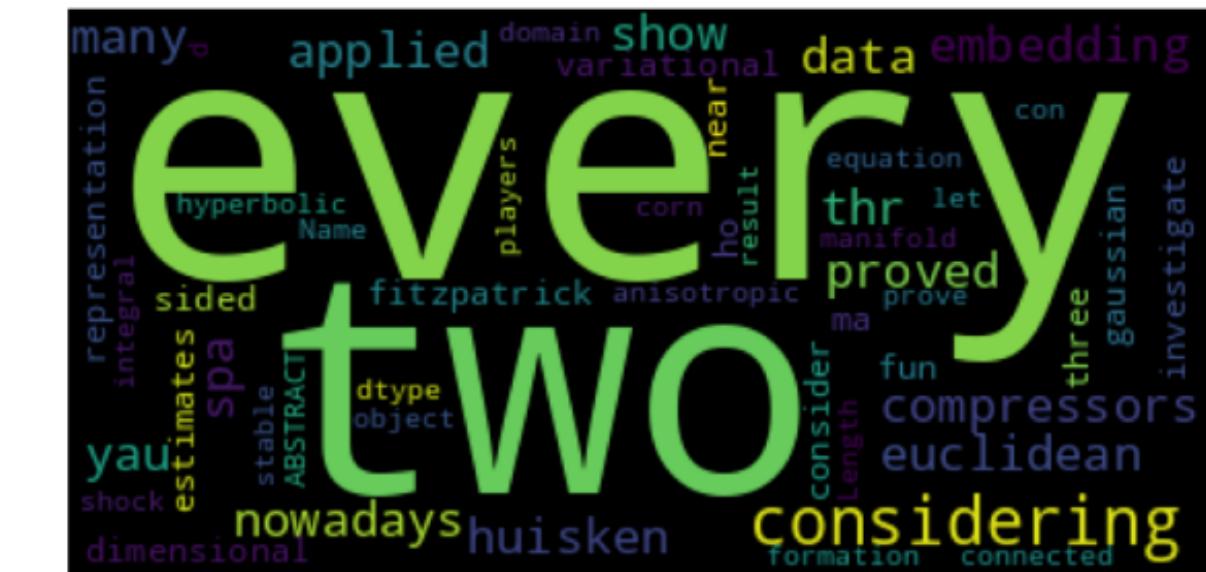


TOP WORDS FOR A GIVEN TOPIC

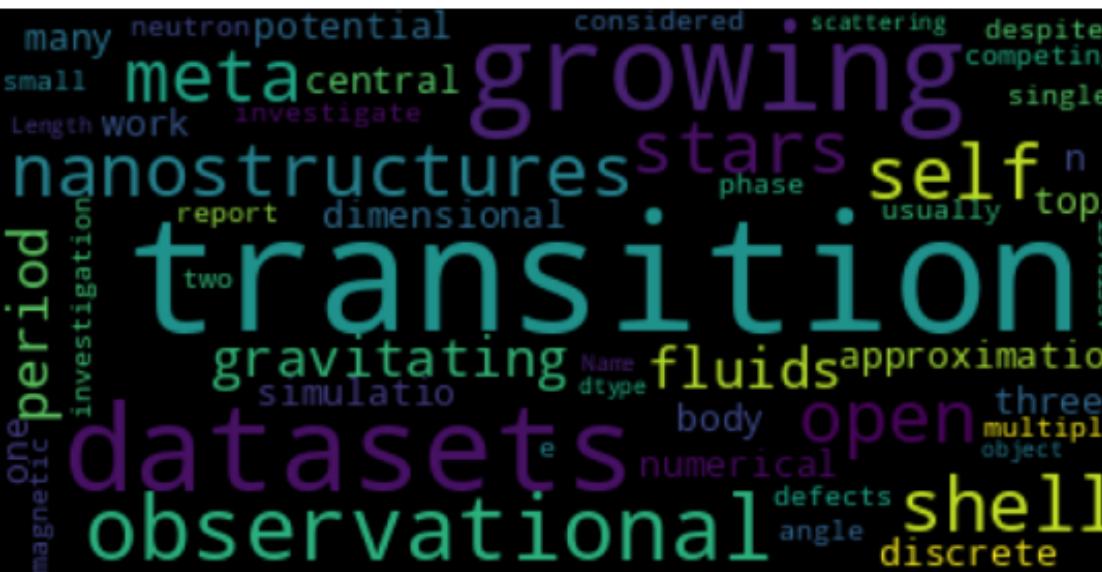
COMPUTER SCIENCE



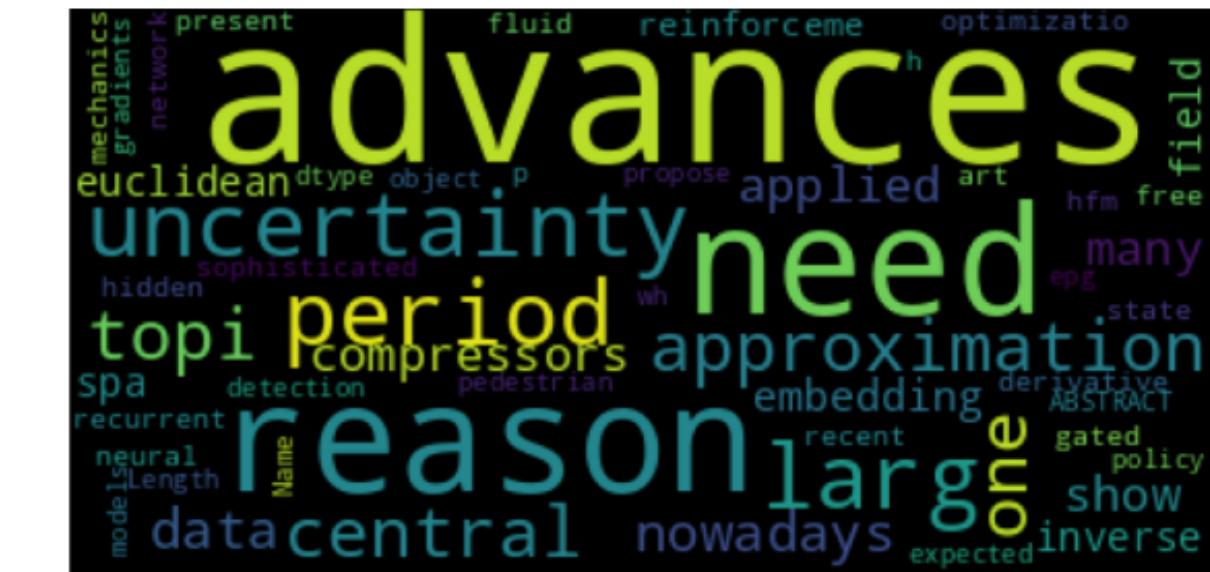
MATHEMATICS



PHYSICS



STATISTICS



Modèle LDA



Utilisation d'un modèle de probabilité basé sur la loi de Dirichlet afin de déterminer en fonction des mots contenus dans le texte initial quel sera le sujet de ce dernier, car en fonction du topic certains mots ont plus de chance d'apparaître que d'autres.

Classification

- ▼ **4 Preprocessing** ¶
 - ▶ **4.1 Remove punctuations**
 - ▶ **4.2 Remove stopwords**
 - ▶ **4.3 Remove frequents words**
 - ▶ **4.4 Remove rare words**
 - ▶ **4.5 Lemmatization**

Classification

	precision	recall	f1-score	support
0	0.74	0.63	0.68	123
1	0.33	0.23	0.27	120
2	0.35	0.37	0.36	273
3	0.82	0.66	0.73	111
4	0.68	0.58	0.63	130
5	0.66	0.53	0.58	196
6	0.74	0.66	0.69	119
7	0.66	0.47	0.55	102
8	0.81	0.72	0.76	109
9	0.87	0.75	0.81	92
10	0.86	0.60	0.71	70
11	0.79	0.56	0.65	79
12	0.72	0.54	0.62	97
13	0.68	0.67	0.68	756
14	0.67	0.68	0.67	145
15	0.47	0.40	0.43	114
16	0.82	0.71	0.76	89
17	0.49	0.37	0.43	131
18	0.87	0.77	0.82	81
19	0.82	0.69	0.75	189
20	0.70	0.62	0.66	112
21	0.55	0.52	0.53	110
22	0.74	0.66	0.70	189
23	0.79	0.76	0.78	112
24	0.50	0.49	0.50	102
micro avg	0.66	0.59	0.62	3751
macro avg	0.69	0.59	0.63	3751

Deep learning

Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
embedding (Embedding)	(None, 200, 50)	3813750
batch_normalization (BatchN ormalization)	(None, 200, 50)	200
flatten (Flatten)	(None, 10000)	0
Fully_Connected (Dense)	(None, 200)	2000200
dropout (Dropout)	(None, 200)	0
la_base (Dense)	(None, 200)	40200
Output (Dense)	(None, 25)	5025
dropout_1 (Dropout)	(None, 25)	0

Epoch 6/20
351/351 [=====] - 27s 78ms/step - loss: 0.3107 - accuracy: 0.5402 - val_loss: 0.2137 - val_accuracy: 0.3920



Axes d'amélioration

Axes amélioration

- **Securiser l'API**
- **Finir le déploiement sur Azure**
- **Améliorer le modèle**
- **Développer le Deep Learning => LSTM**

Merci !

Camille, Naïs, Ahmed

