

Fake news detection .

by YAC agency

Notre équipe



Youcef Boucheta
data analyst



Amani Ellafi
data scientist



Clément Charp
data engineer

Sommaire

Rappel de la mission

Interface utilisateur

Notre méthode

Déploiement 2022/2023



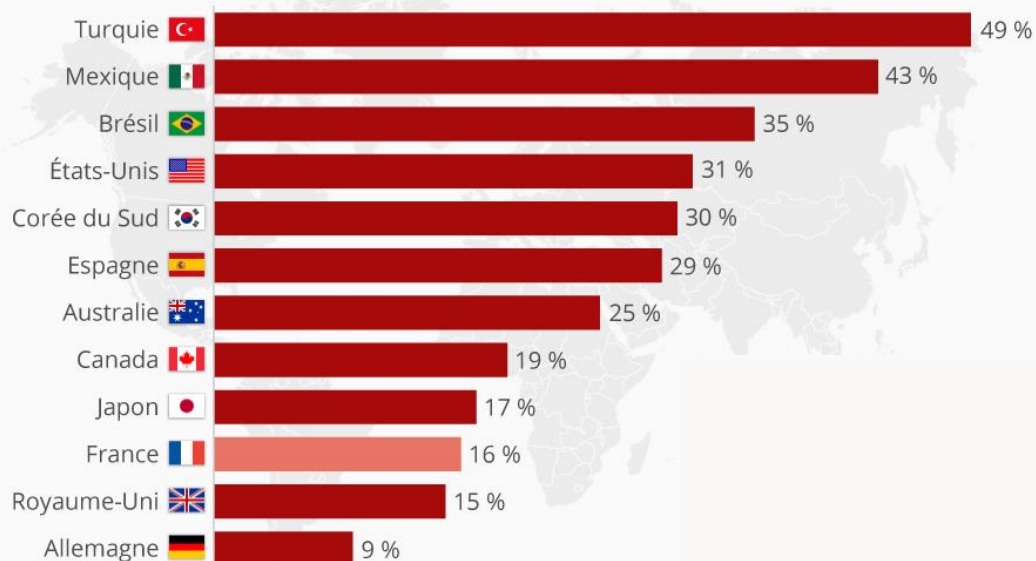


Créer un modèle
capable de
détecter les
“fake news”

Les “fake news”, un problème mondial

L'exposition aux "fake news" dans le monde

Répondants déclarant avoir été exposés à de fausses informations *



Source : Reuters Institute Digital News Report 2018

L'importance de détecter les “fake news” pour REUTERS ?



ADN

Reuters fournit des informations fiables qui permettent aux personnes comme aux machines de prendre des décisions éclairées.



Journalistes

Donner le moyen au journaliste de garantir une information fiable, rapide et éviter de se faire striker

!FAKERS

Track down the real news by Reuters

Document
technique



Notre méthode

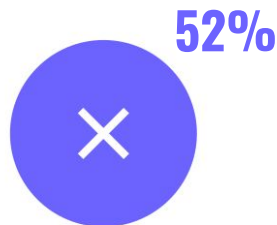
Choix/analyse/préparation du dataset

Entraînement et test du modèle

Sélection du meilleur modèle

Déploiement de l'interface utilisateur

Le dataset



Fake news

23 481 articles

Articles collectés sur Politifact (une organisation de *fact checking* aux USA)



48%

Real news

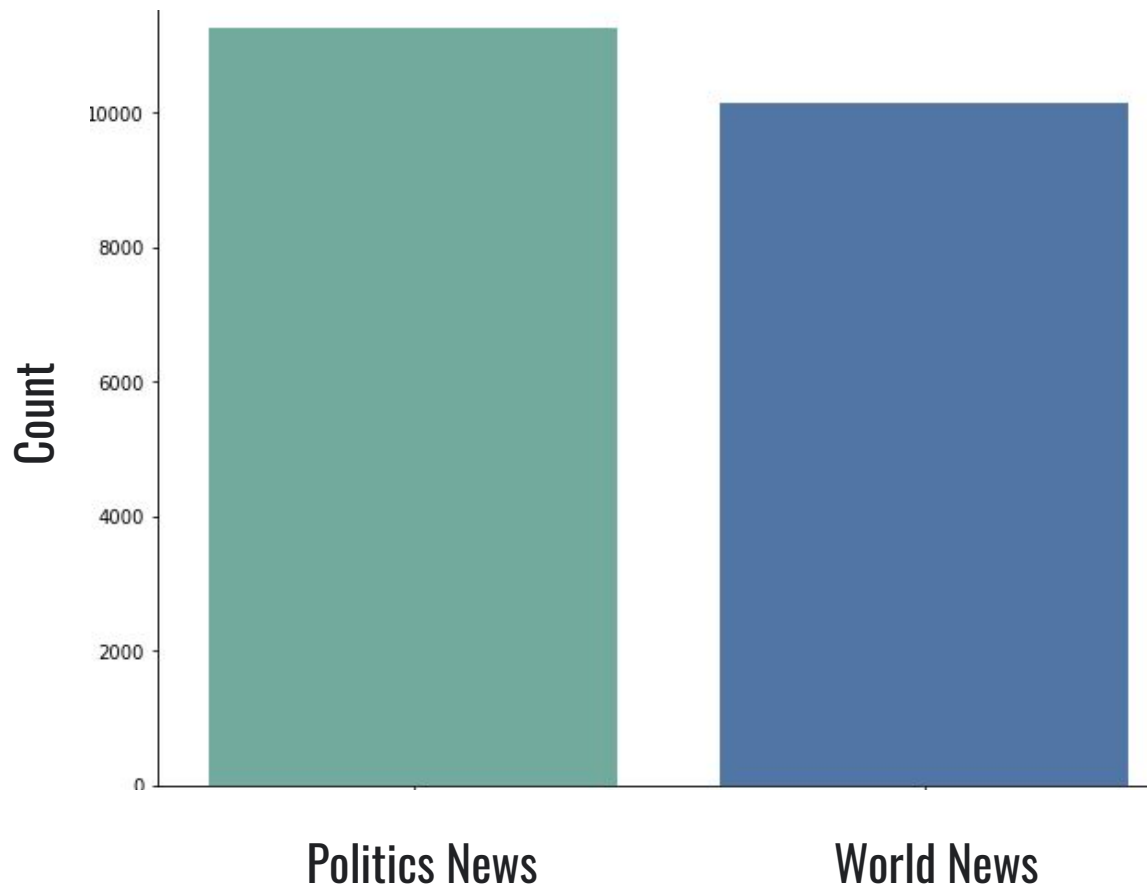
21 417 articles

Articles récoltés sur le site Reuters.com (agence de presse international)

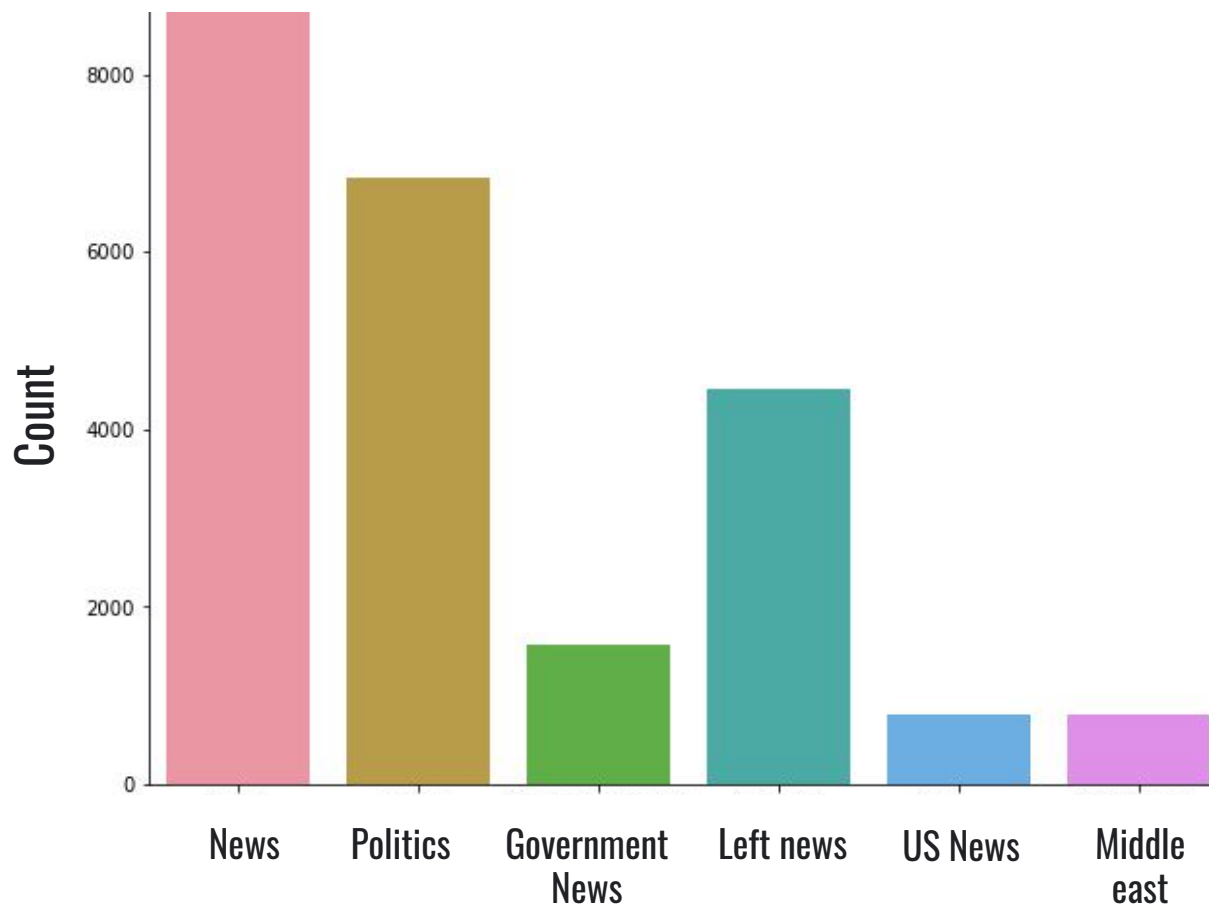
Le dataset

- 44 898 articles (lignes)
- 4 colonnes

title titre de l'article	text texte de l'article	subject sujet de l'article	date date de publication de l'article	Fake or not ? Fiabilité de l'article
Donald Trump Sends Out Embarrassing New Year's Eve Message; This is Disturbing	Donald Trump just couldn't wish all Americans a Happy New Year and leave it at that. Instead, he had...	News	December 31, 2017	1
As U.S. budget fight looms, Republicans flip their fiscal script	WASHINGTON (Reuters) - The head of a conservative Republican faction in the U.S. Congress, who voted...	politicsNews	December 29, 2017	0



**True
data**

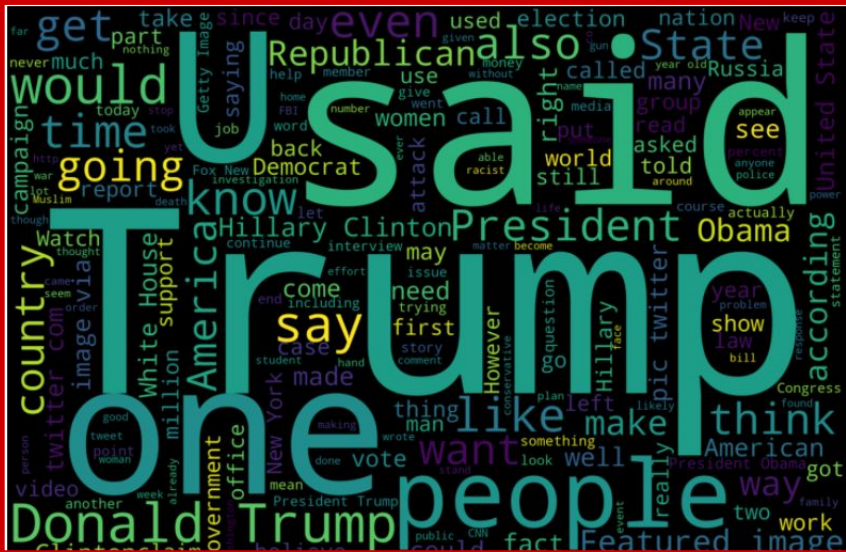


Fake data

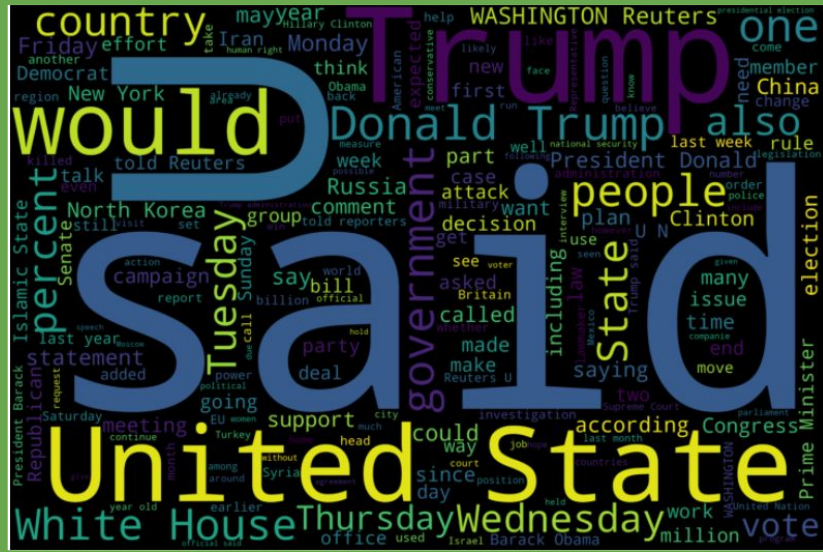
Exploration du dataset

YAC agency

Fake data



True data



Preprocessing

True data

Real news avaient une source de publication de (“Washington (Reuters), Twitter tweet”).

- Enlever les détails de la publication “-”
- Supprimer le Text vide de la ligne 8970

WASHINGTON (Reuters) - Transgender President Donald Trump's administration decision to put on hold orders

```
8771 In a speech weighted with America's complicate...
8970
9008 The following timeline charts the origin and s...
9009 Global health officials are racing to better u...
```

Preprocessing

Fake data

Les fausses nouvelles n'ont pas une source de publication.

- Les titres en capital cases
- Fake data contient 630 textes vides

NS: MARYLAND GOVERNOR BRINGS IN NATIONAL G
FULL VIDEO: THE BLOC
(VIDEO) HILLARY CLINTON: RELIGIOUS BELIEFS MU
E PROTECTING OBAMA: WON'T RELEASE NAMES OF
RICAL SNL TAKE ON HILLARY'S ANNOUNCEMENT: 'B

title	text	subject
AGENCY		left-news
N CASH		left-news
ORTION		left-news
MERICA		left-news
E BACK!'		left-news

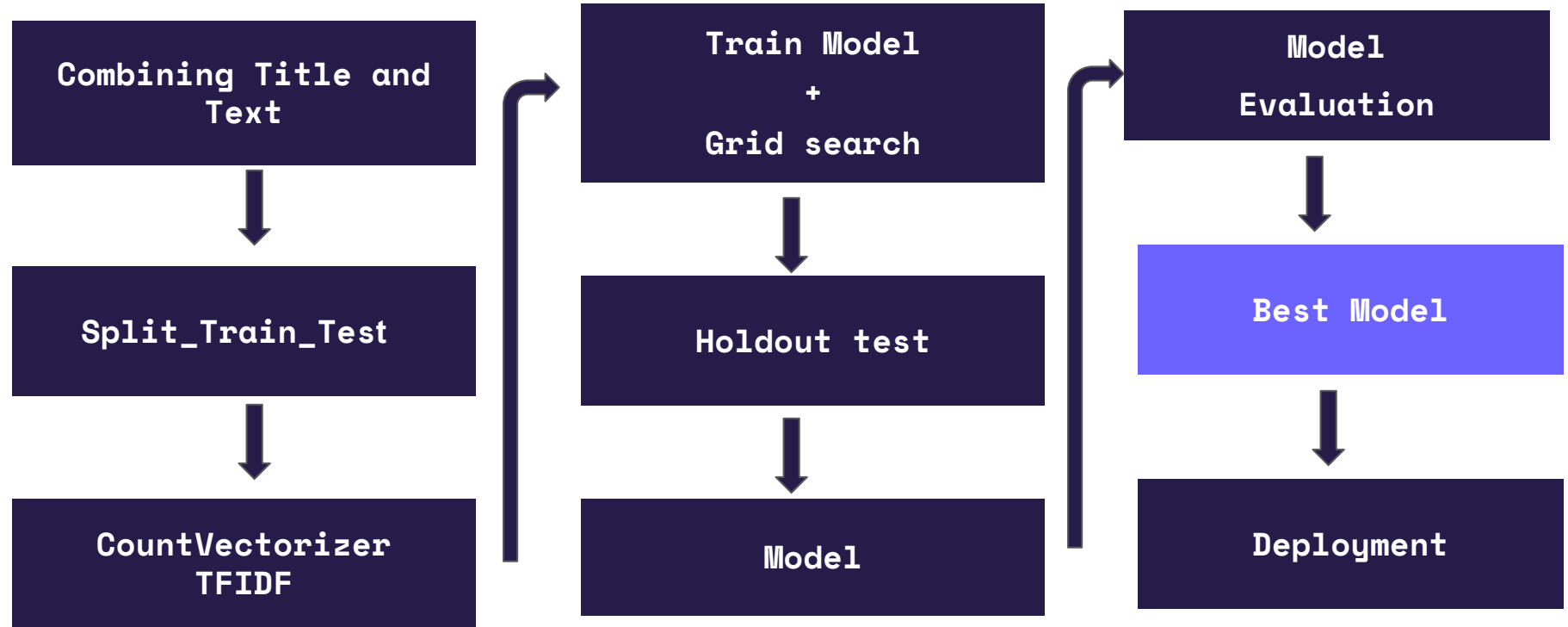
Preprocessing

Feature selection

CountVectorizer pour compter le nombre de mots (fréquence des termes), limiter la taille de votre vocabulaire.

Tfidftransformer et Tfidfvectorizer permettent de convertir une collection de documents bruts en une matrice de fonctionnalités TF-IDF.

Méthode



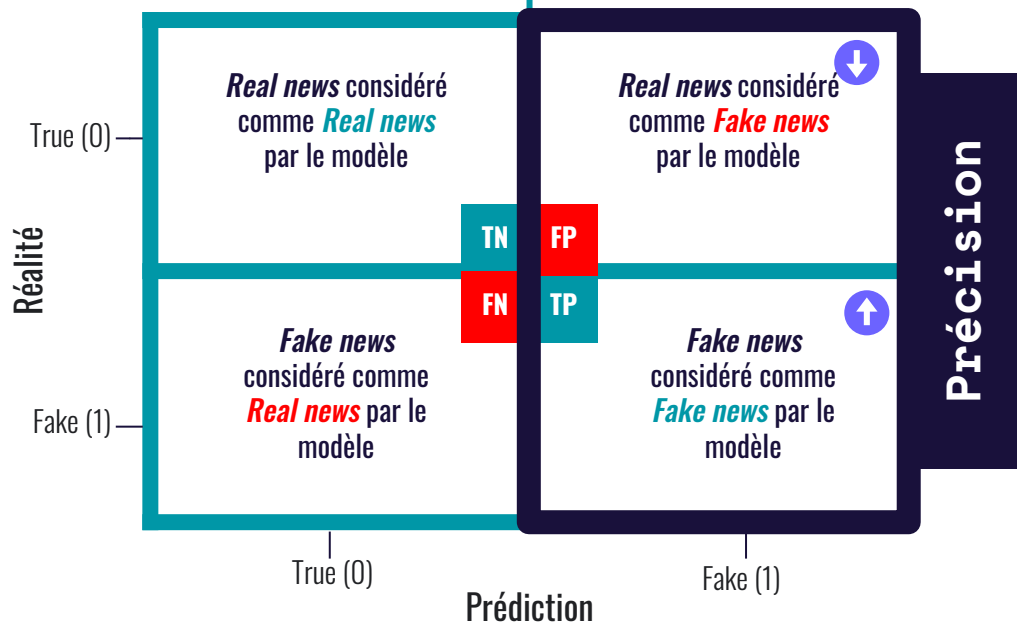
La classification binaire, quèsaco ?

44 898 articles

75% entraînement

25% test

Calcul de la performance avec les 25% test



Nos différents modèles

MultinomialNB

RandomForest

LSTM

Best model ?



1er modèle : MultinomialNB

Train score : 95% - Test Score : 95% - **Score de précision : 99%**

Réalité	True (0)	5125	229
		TN	FP
Fake (1)	311	5125	TP
		FN	
		True (0)	Fake (1)
		Prédiction	

- Tfidfvectorizer
- MultinomialNB
- GridSearch

2ème modèle : RandomForest

Accuracy : 98% - F1 Score : 99% - Recall : 98% - **Precision : 99%**

Réalité	True (0)	5302	52
		TN	FP
Fake (1)	63	5808	
		FN	TP
	Prédiction	True (0)	Fake (1)

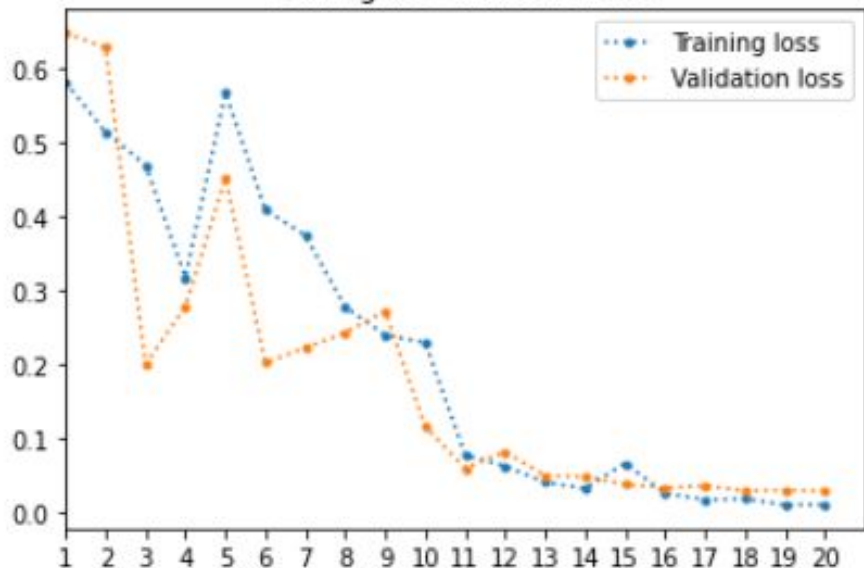
- CountVectorizer
- Tfidftransformer
- RandomForestClassifier
- GridSearch

3ème modèle : Deep learning LSTM

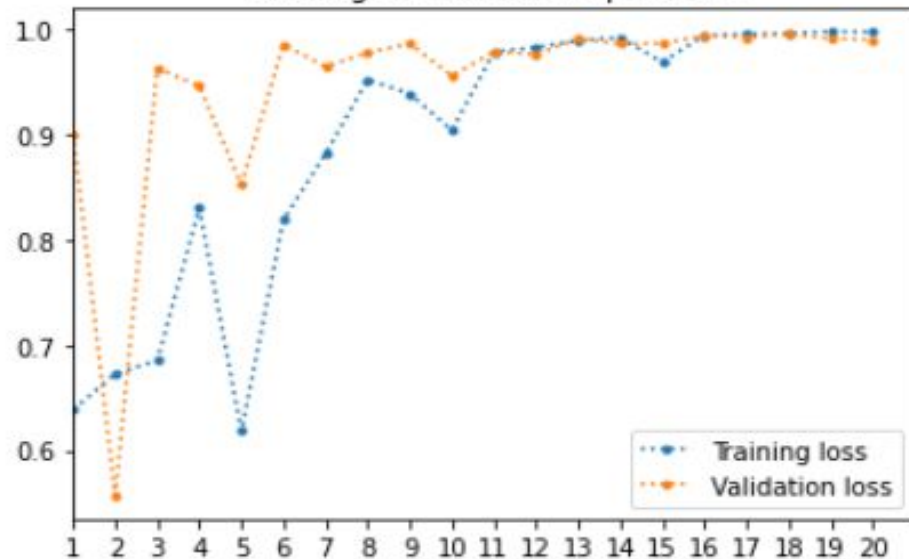
Loss function : 0,028

Precision : 99%

Training and Validation Loss



Training and Validation precision



Les différents modèles

YAC agency

3ème modèle : Deep learning LSTM

Precision : 99%

`tokenizer.texts_to_sequences(X)`

LSTM(128)

Dense Layer
Sigmoid

Optimizer:Adam

loss:binary_crossentropy

Les différents modèles

YAC agency

Modèle final : Logistic regresssion



Accuracy : 99% - F1 Score : 99% - Recall : 99% - **Precision : 99%**

Réalité	True (0)	5320	34
		TN	FP
Fake (1)	56	5815	
		FN	TP
	Prédiction	True (0)	Fake (1)

- Tfidfvectorizer
- LogisticRegression
- GridSearch

Les différents modèles

YAC agency

Et après ?



Déploiement de
Fakers en ligne
pour l'interne

**1er trimestre
2022**

Mise à jour du modèle
avec articles de 2022
+
Déploiement en widget
pour le grand public

**4ème trimestre
2022**

Développement
du modèle en
multilangage

**Courant
2023**



Merci pour votre attention

Des questions ?