

---

---

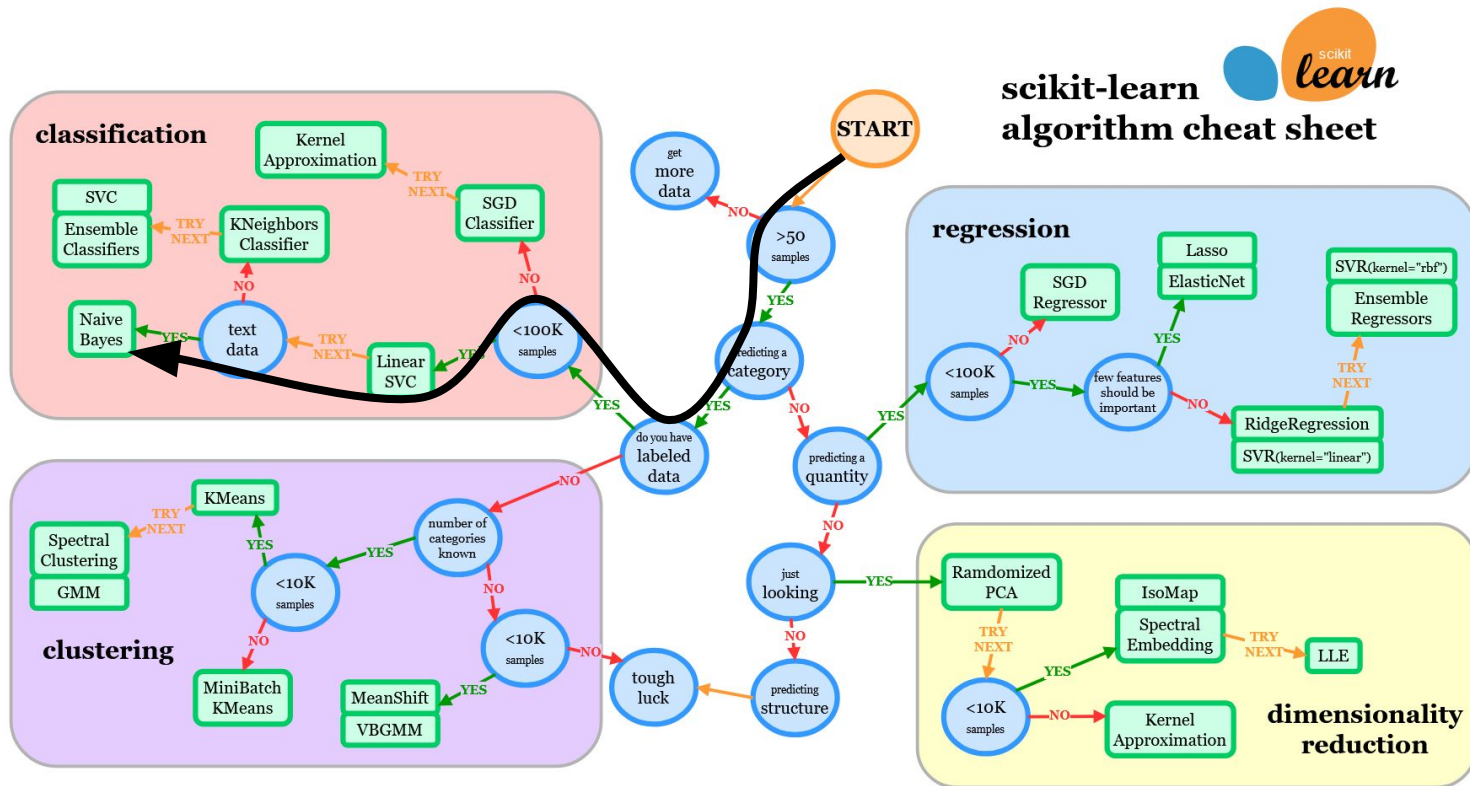
# SpamClassifier

Fabien Herry  
Melody Duplaix

---

---

# Modèles testés



# Premières comparaisons

## LinearSVC Scikit-Learn après vectorisation

	precision	recall	f1-score	support
0	0.98	1.00	0.99	966
1	0.99	0.86	0.92	149
accuracy			0.98	1115
macro avg	0.99	0.93	0.95	1115
weighted avg	0.98	0.98	0.98	1115

	text	résult
0	Join us today ! Flexible work without constrain...	ham
1	Chronopost : your package 7d6595466533 is wait...	ham
2	I'm a nigerian prince i have 250 millions doll...	spam
3	There was a hack on your amazon account clic o...	ham
4	Your energy provider want to talk to you about...	ham
5	URGENT: Your account has been compromised. Cal...	ham

# Premières comparaisons

## Naive Bayes simple

```
Exactitude : 0.92  
Rappel : 0.99  
Rappel négative: 0.65  
Precision : 0.92  
Précision négative: 0.93  
F1-Score : 0.95  
F1-Score Négatif: 0.76
```

	text	result
0	Join us today ! Flexible work without constrain...	spam
1	Chronopost : your package 7d6595466533 is wait...	spam
2	I'm a nigerian prince i have 250 millions doll...	ham
3	There was a hack on your amazon account clic o...	ham
4	Your energy provider want to talk to you about...	ham
5	URGENT: Your account has been compromised. Cal...	spam

Pourtant, Modèle simple plus performant sur des spams écrit à la main, non dans le dataset

↪ Overfitting

# Problématique exploratoire

## Question de models:

- Quel modèle choisir
- Quel features
- Sur quel jeu de données entraîner

## Question de Data :

- Chercher plus de données
- Définir un test de satisfaction

# Automatisations

Trouver les combinaisons de features / modèles les plus efficaces

I - Module pour créer et entraîner un modèle différent rapidement avec le même preprocessing  
→ pipeline

II - Automatisation d'une analyse de base de donnée

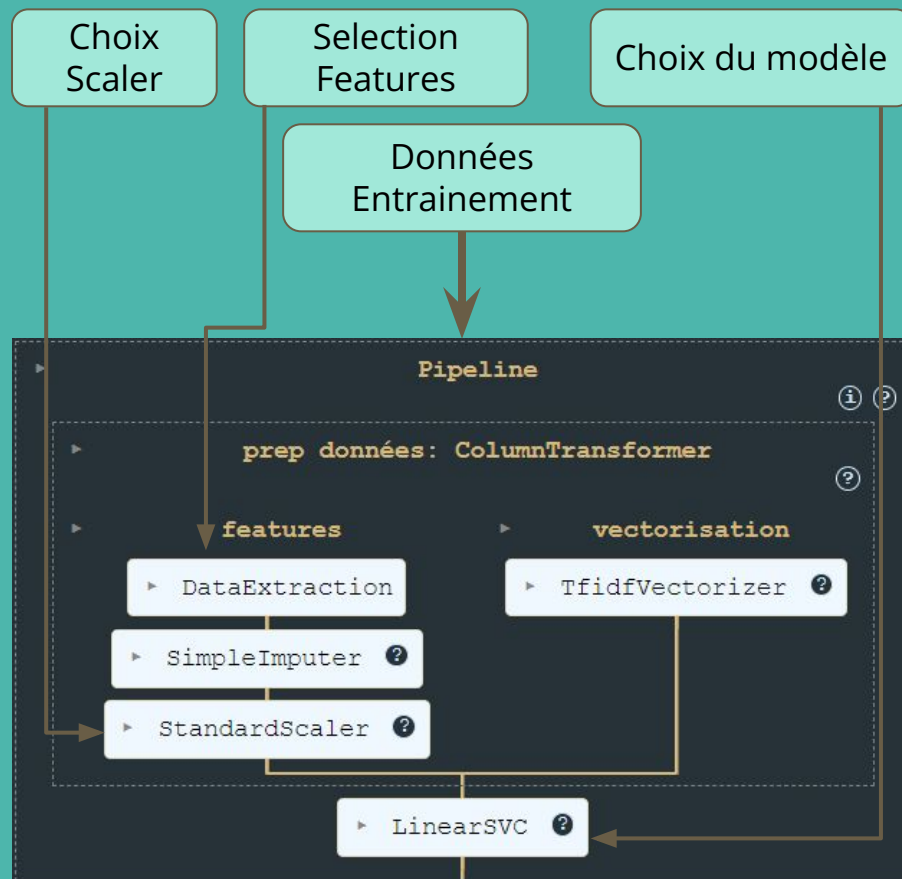
III - Modules pour tester l'efficacité de toutes les combinaisons de features possibles pour chaque modèle, pour chaque combinaison de dataset → validation croisée

---

# Pipeline de base

Modulaire pour  
automatisation

## Paramètres d'entrée



# Données disponible

- Base de données d'origine
- Base de données étudiant nigerien
- Telegram
- Email Enron . Probleme de formatage

---

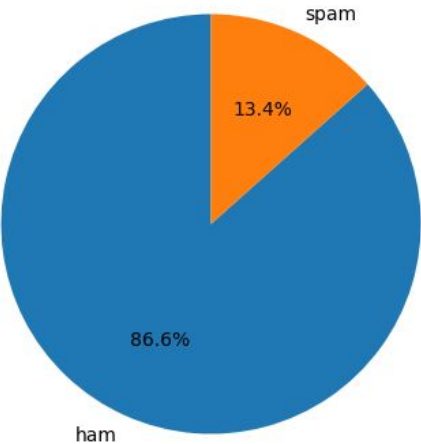


# Exploration des données et features

- Répartition spam / ham
- Nombre caractères par sms
- Nombre mots par sms
- Nombre de symboles monétaires
- Nombre de chiffres
- Nombre de liens

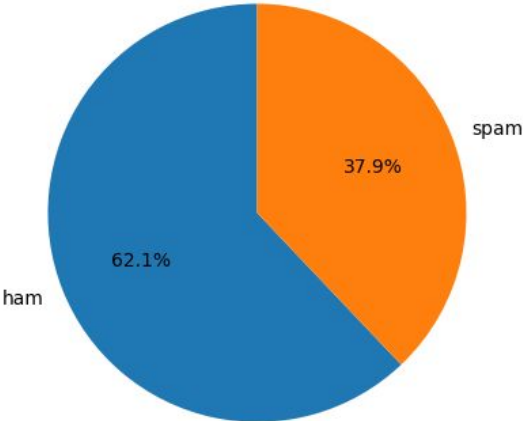
# Visualisation des données

Repartition des données



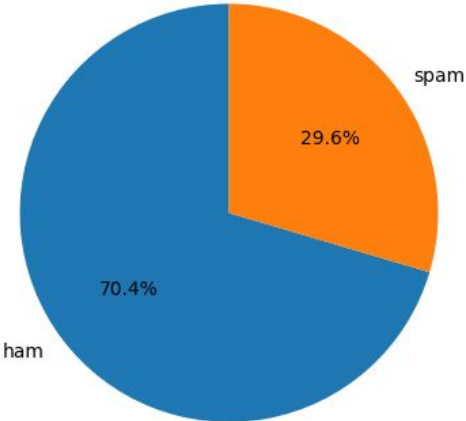
Dataset original

Repartition des données



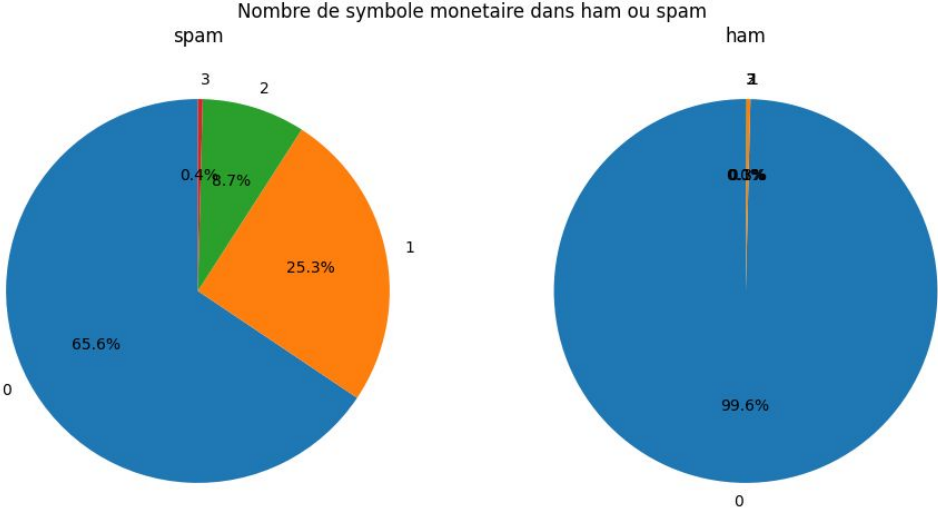
Dataset Nigérien

Repartition des données

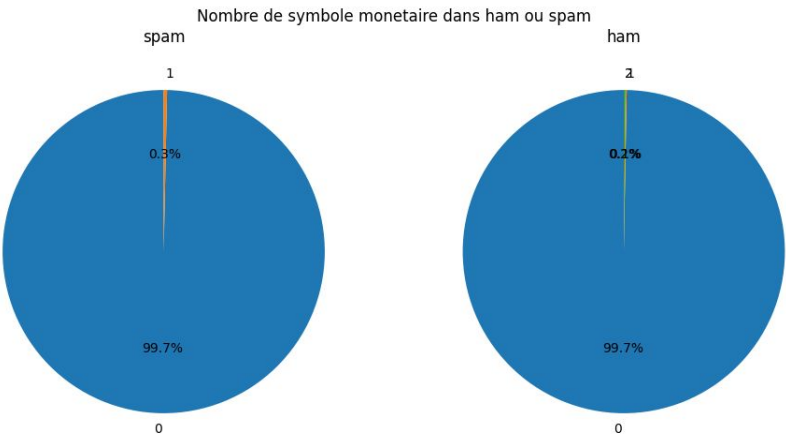


Dataset Télégram

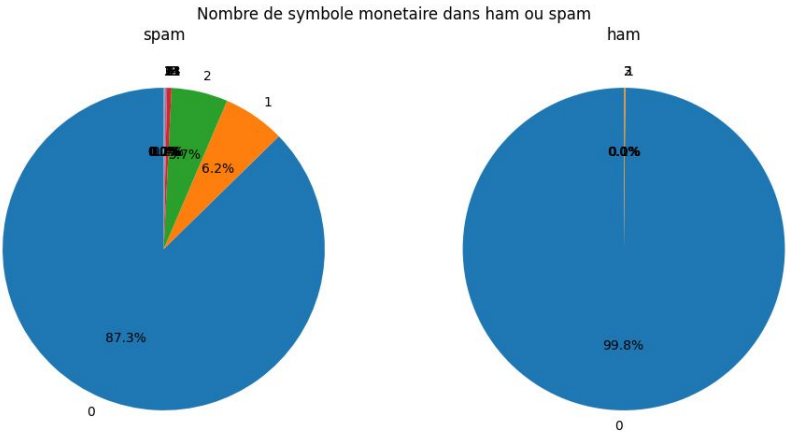
# Visualisation des données



Dataset original

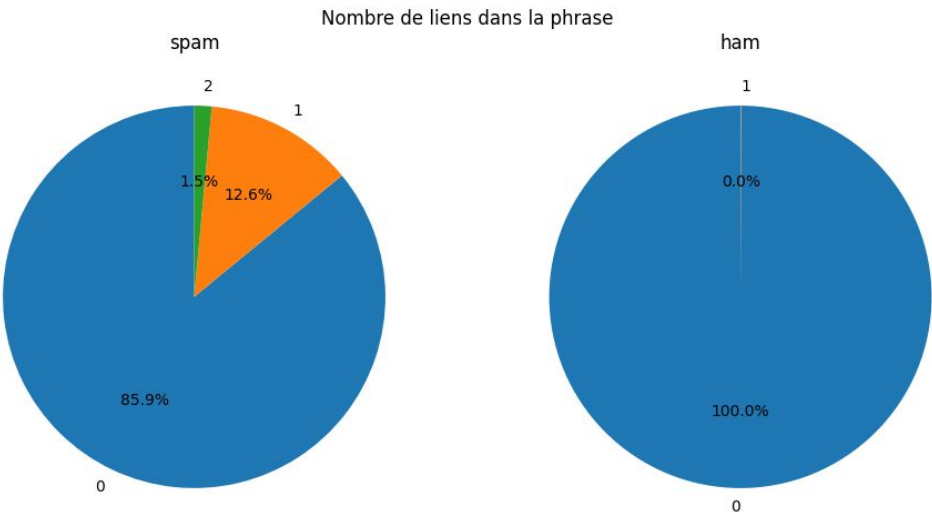


Dataset Nigérien

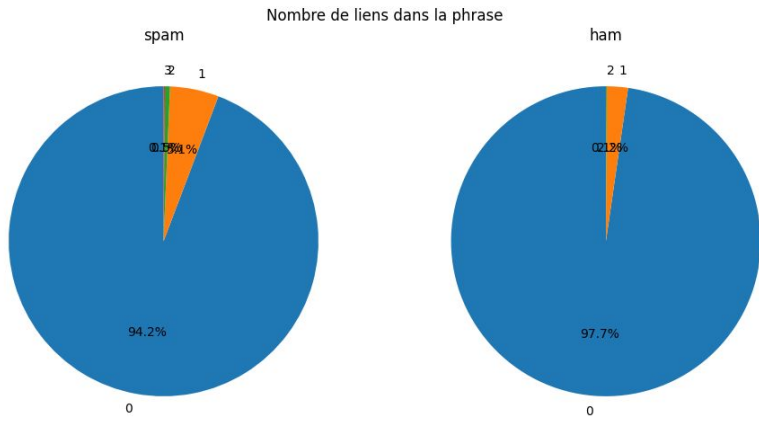


Dataset Télégram

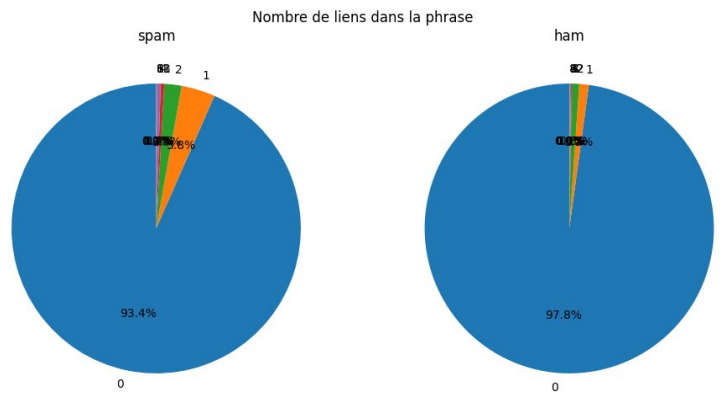
# Visualisation des données



Dataset original



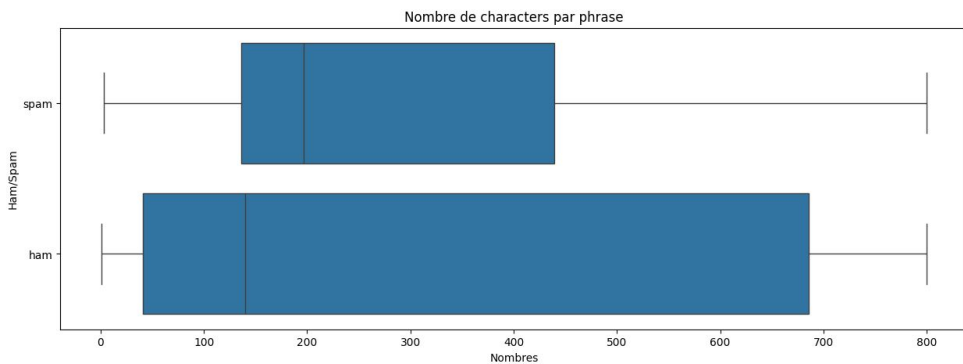
Dataset Nigérien



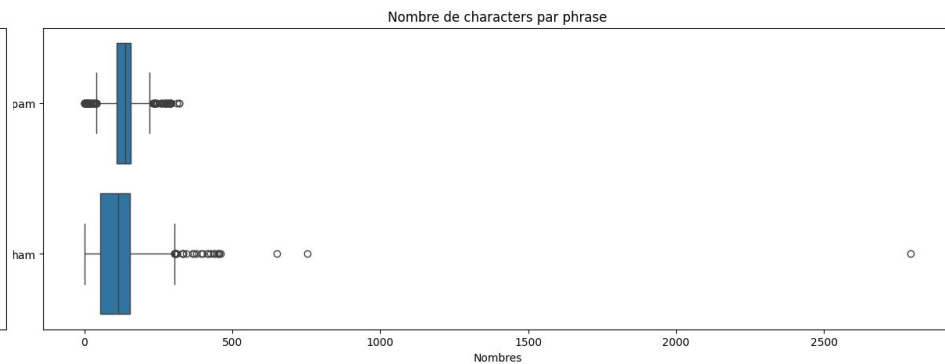
Dataset Télégram

# II Visualisation des données

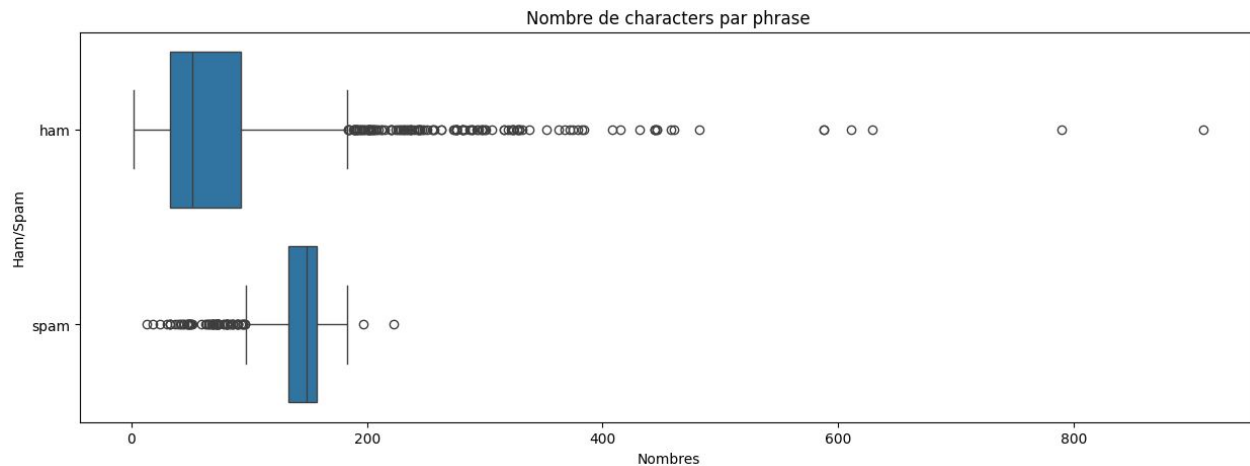
Dataset Niégerien



Dataset Télégram

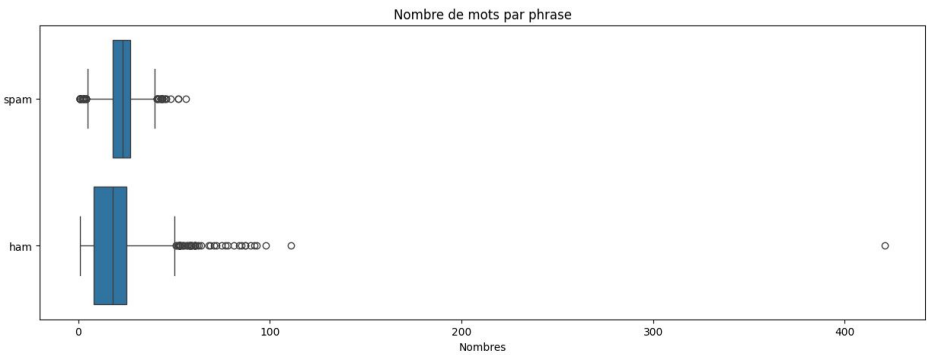


Dataset original

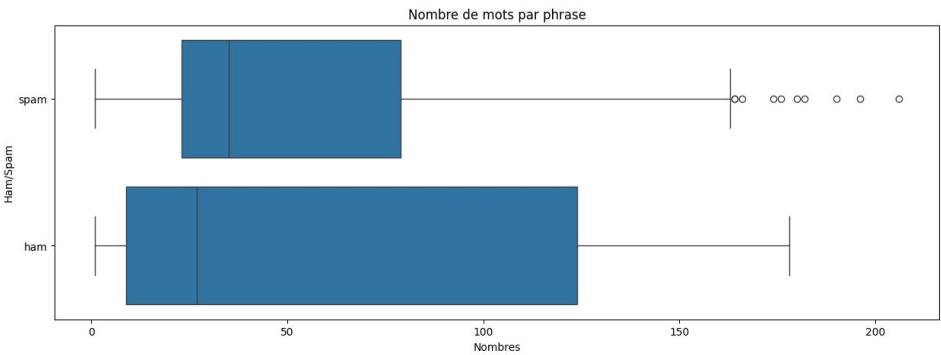


# Visualisation des données

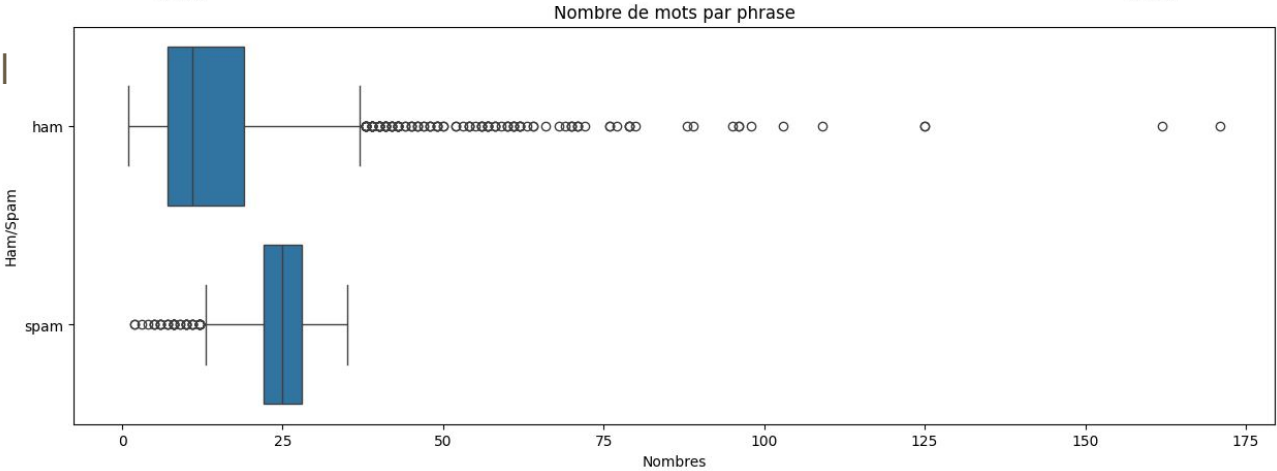
Dataset Niégerien



Dataset Télégram

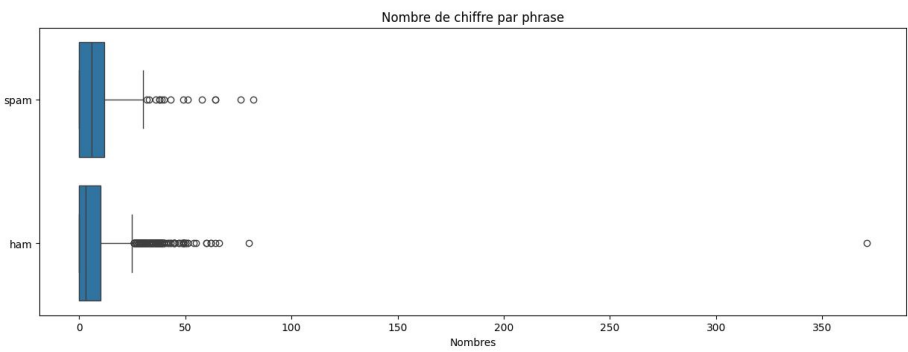


Dataset original

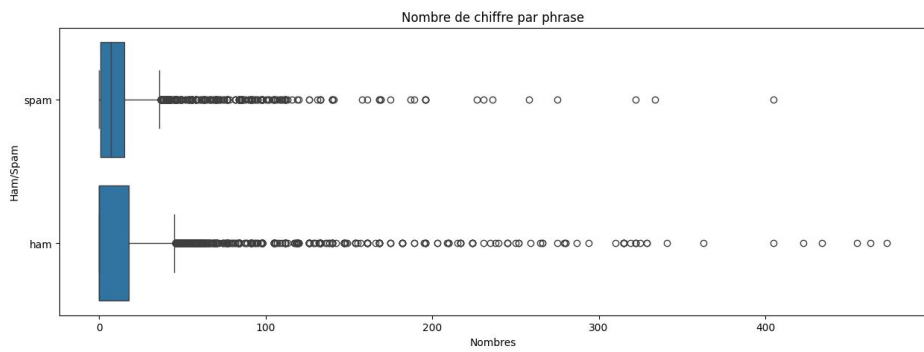


# Visualisation des données

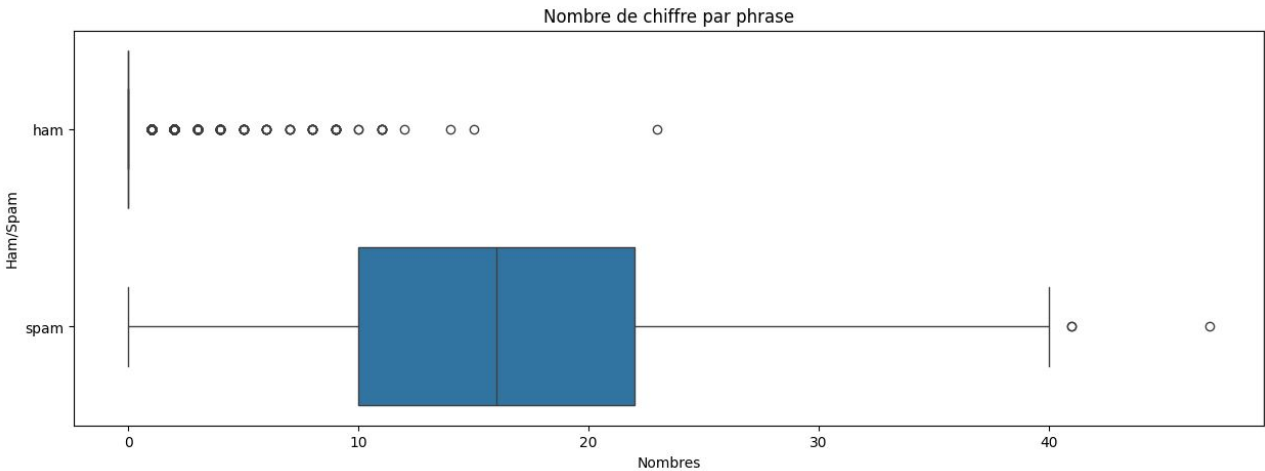
Dataset Niégerien



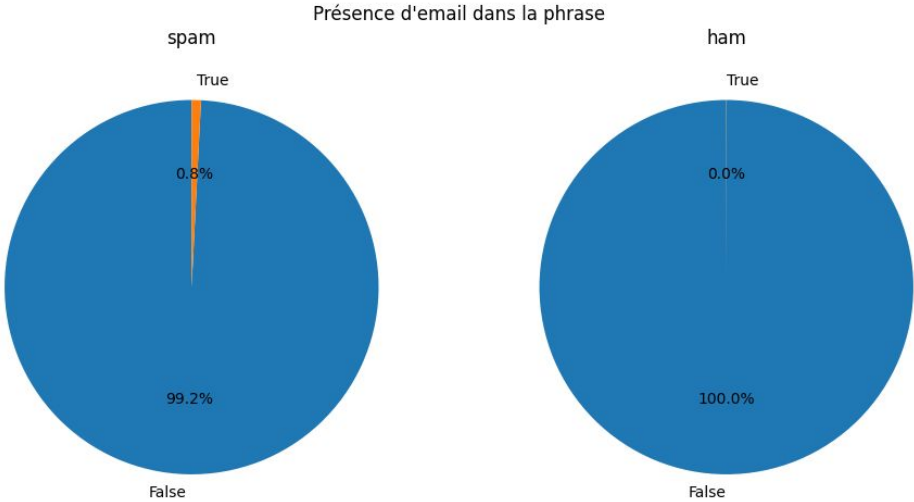
Dataset Télégram



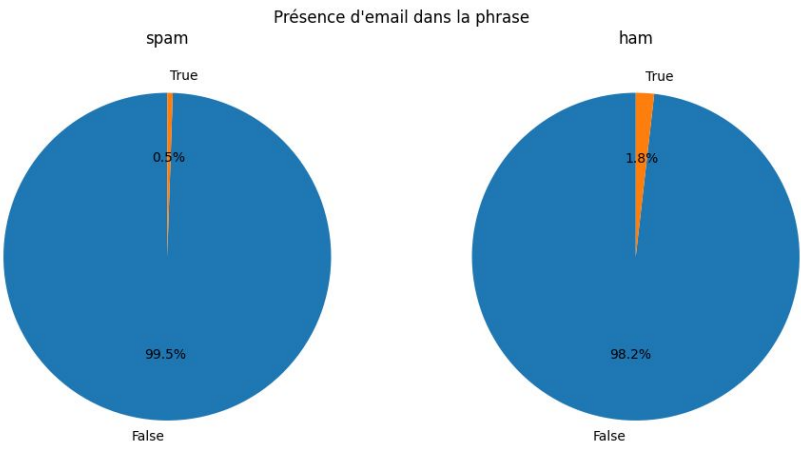
Dataset original



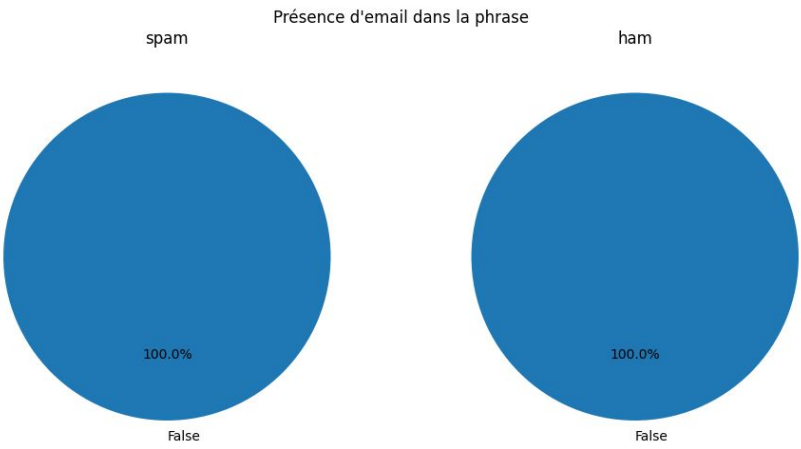
# Visualisation des données



Dataset original



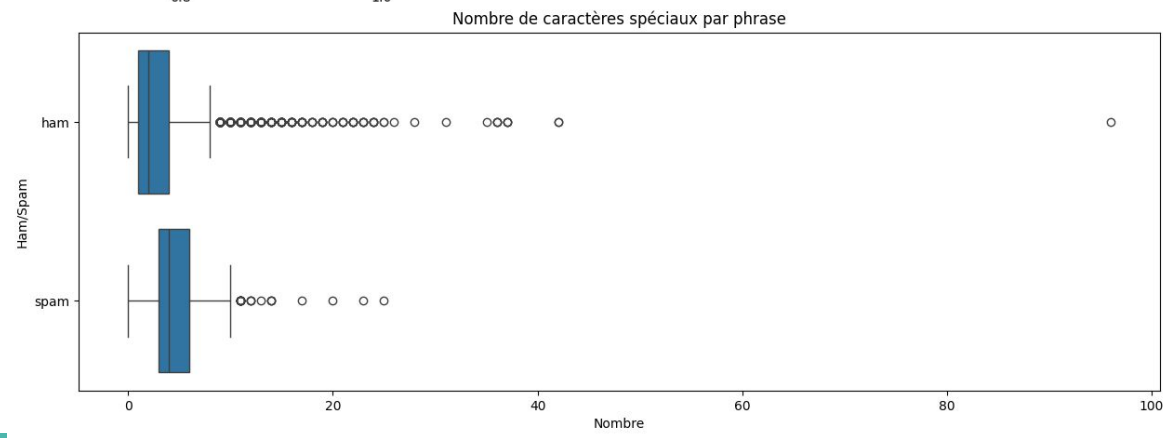
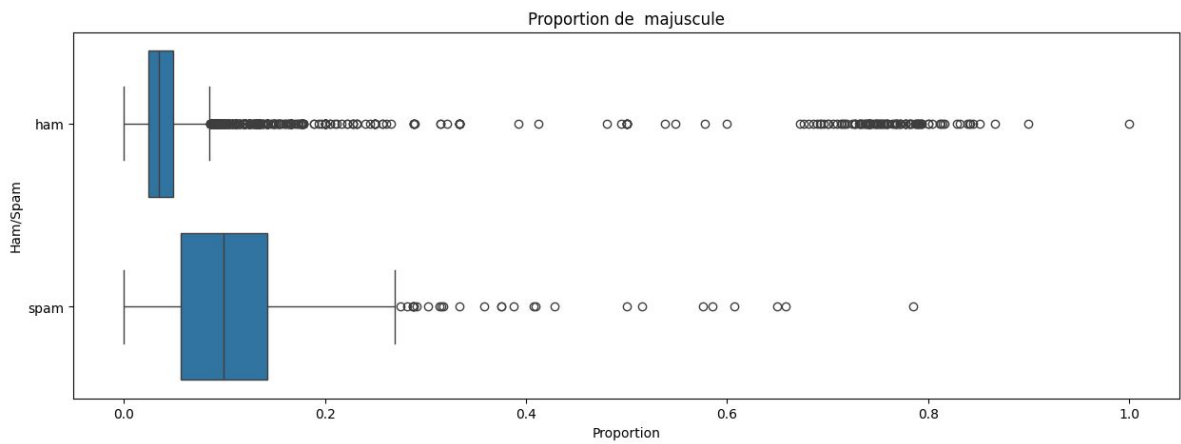
Dataset Nigérien



Dataset Télégram



# Visualisation des données





# Tests des modèles

Guidelines de test :

- Entraînement sur le même découpage des jeux de données
- Test sur les mêmes jeux de données non disponible en entraînement
- Test final sur des spams écrits mains

---

# Gestion base de données

```
import Automatisation.DataSetManagement as DataSetManagement
import Automatisation.EvaluateModelsFeatures as EvaluateModelsFeatures

dataSet = {"first":"DataSetBrut/BD1.txt","nigerian":"DataSetBrut/DataSmsSpamNH.csv","telegram":"DataSetBrut/telegram_spam_dataset.csv"}

dataManag = DataSetManagement.DataSetManage(dataSet,42)
```

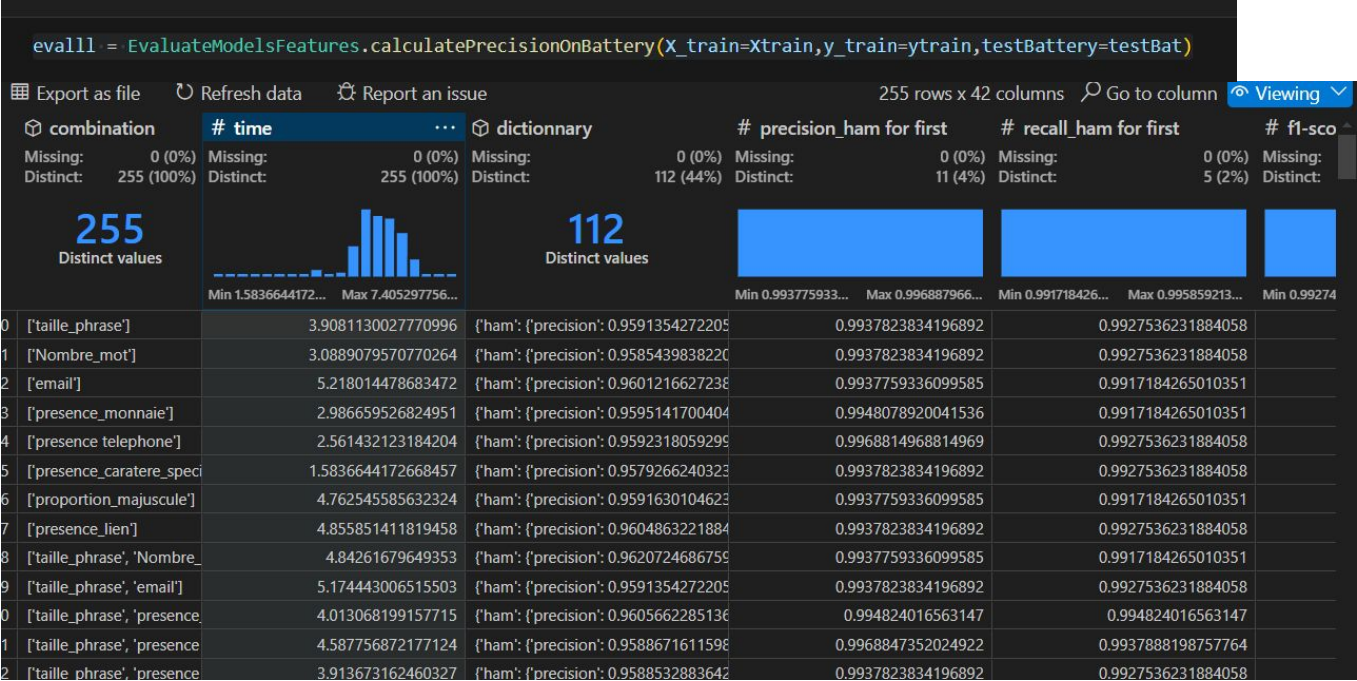
✓ 0.4s

```
testBat = dataManag.GetAllTest()
Xtrain,ytrain = dataManag.GetCombinedTrain(["first","nigerian","telegram"])
```

- Système modulaire de gestion des base de données et séparation Train/Test

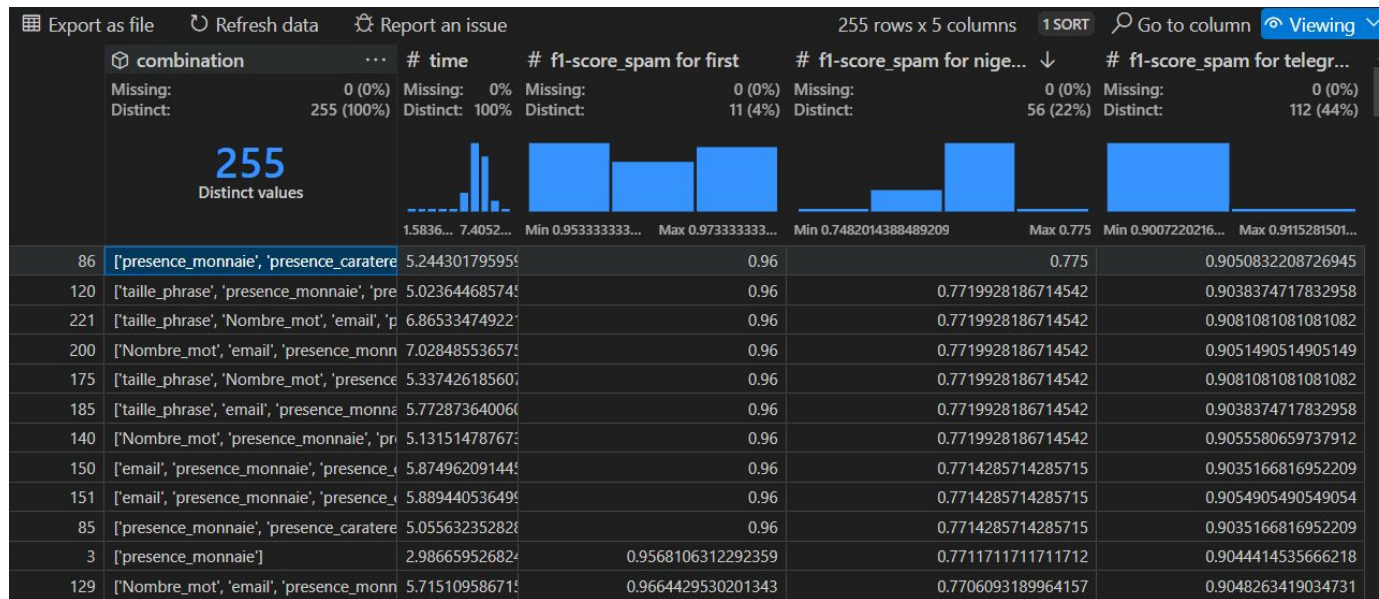


# Résultat du test automatique



- Test long à effectuer chaque modèle 5 secondes
- Beaucoup d'information a résumé

# Resumé



- Score correct sur toutes les data
- Influence faible du choix des features

# Choix du modèle

Modèle choisi : Linear SVC  
entraîné sur les trois bases des  
données

sur features : `['presence_monnaie',  
'presence_caratere_speciaux',  
'presence_lien']`

---

# Résultat sur test manuel

```
modelPipe.predict(pd.DataFrame(["Join us today ! Flexible work without constainst. Earn between 50 and 3600 euros per  
day payment sent daily. To know more add : https://wam.me/nawak5546431645",  
"Chronopost[]: your package 7d6595466533 is waiting clic here to comfirm the delivery of your package[]: https://  
chronoquost.fr",  
"I'm a nigerian prince i have 250 millions dollars ($250 000 000) in asset and gold and i need help to fight corruption  
charge levied against me in England . You will be rewarded for your help please reply.",  
"There was a hack on your amazon account clic on this link to resolve the issue https://www.amozon.com/reclamation . If  
you don't you will be fined $250 dollars a day",  
"Your energy provider want to talk to you about a reduction of your bill call this number 0648874598",  
"URGENT: Your account has been compromised. Call 1-800-123-4567 immediately to secure your account!"],columns=["text"])))
```

22]

✓ 0.0s

Python

```
array(['spam', 'spam', 'spam', 'spam', 'ham', 'spam'], dtype=object)
```

# Résultats comparatifs sur tests manuels

Premier Model

résult
ham
ham
spam
ham
ham
ham

1/6

Naif Bayesien Main

result
spam
spam
ham
ham
ham
spam

3/6

Model SVC

result
spam
spam
spam
spam
ham
spam

5/6





# Application

# Démonstration !

---

# Spam Detector

## To detect if the sms is a spam

Enter the sms here

Join us today ! Flexible work without constainst. Earn between 50 and 3600 euros per day payement sent daily. To know more add : <https://wam.me/nawak5546431645>

Detect

The sms is a spam

# Spam Detector

## To detect if the sms is a spam

Enter the sms here

Hello it's been a while ! What are you up to ? I'm coming to town this weekend , are you free for a beer?

Detect

The sms is not a spam

# Piste de continuation

- Appliquer l'automatisation aux différents modèles
- Implémenter un système de barre de temps de chargement pour long test
- Faire un système de gestion de l'entraînement avec de nouvelles données.  
Gérez la pipeline d'insertion des données. (Concaténation sensible taille base de données)
- Tester le modèle avec des data actuel et concrète (Récupération de Data)
- Sérialiser un modèle pour le mettre en production
- Créer des variable environnement (localisation des dataset)
- Data Analyse : implémenter graphique en courbes

*merci!*

**Est-ce que vous avez des questions ?**

# Data Extraction

```
class DataExtraction(BaseEstimator, TransformerMixin):  
  
    def __init__(self, features_list):  
        self.features_list = features_list  
    def fit(self, X, y=None):  
        return self  
    def transform(self, X):  
        retour = X  
        for feature_name, feature_function in self.features_list.items():  
            retour[feature_name] = X["text"].apply(feature_function)  
        retour = retour.drop(columns = "text")  
        return retour
```

# List extraction features

```
features = {  
    "taille_phrase": lambda x: len(x),  
    "Nombre_mot": lambda x: len(x.split()),  
    "email": lambda x: len(re.findall(r'[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}', x)) > 0,  
    "presence_monnaie": lambda x: 1 if re.search(r'[\$€£]', x) else 0 ,  
    "presence_telephone": lambda x: 1 if re.search(r'\b\d{10,}\b', x) else 0,  
    "presence_caractere_speciaux": lambda x: 1 if re.search(r'[!@#%^&*(),.?":{}|<>]', x) else 0,  
    "proportion_majuscule": lambda x: sum(1 for c in x if c.isupper()) / len(x) if len(x) > 0 else 0,  
    "presence_lien": lambda x: 1 if re.search(r'\b(http|www)\S+', x) else 0  
}
```