



Université PSL – Paris Dauphine

Master Big Data & IA

2025–2026

Analyse de sentiments TF - IDF, RNN et Attention

Rapport de projet NLP

Auteurs

Fatma CHAHED & Aziz DHIF

Enseignante

Mme Amel Tarifa

Plan

| | | |
|----------|---|-----------|
| 1 | Présentation du problème | 1 |
| 2 | Présentation des données utilisées | 3 |
| 2.1 | Chargement des données | 3 |
| 2.2 | Traitement des données | 3 |
| 2.2.1 | Quelques statistiques | 3 |
| 2.2.2 | Encodage de la colonne <code>label</code> | 4 |
| 2.2.3 | Gestion des doublons | 5 |
| 2.2.4 | Gestion des valeurs manquantes | 5 |
| 3 | Première partie : TF-IDF | 6 |
| 3.1 | Introduction | 6 |
| 3.2 | Importation des bibliothèques | 6 |
| 3.3 | Prédéfiniion des fonctions | 7 |
| 3.4 | Construire et tester les modèles | 8 |
| 3.4.1 | Régression Logistique (sans régularisation) | 8 |
| 3.4.2 | Ridge | 10 |
| 3.4.3 | Lasso | 10 |
| 3.4.4 | SVM | 11 |
| 3.4.5 | KNN | 12 |
| 3.5 | Améliorations possibles | 13 |
| 3.6 | Visualiser les résultats d'entraînement | 13 |
| 3.7 | Conclusion | 15 |
| 4 | Deuxième partie : RNN | 16 |
| 4.1 | Introduction | 16 |

| | | |
|----------|---|-----------|
| 4.2 | Importation des bibliothèques | 16 |
| 4.3 | Tokenisation | 17 |
| 4.4 | Construire et tester les modèles | 18 |
| 4.4.1 | Embedding pendant l'entraînement | 18 |
| 4.4.2 | Embedding préentraîner avec "fastText facebook 2016" | 20 |
| 4.5 | Visualiser les résultats d'entraînement | 22 |
| 4.6 | Conclusion | 24 |
| 5 | Troisième partie : Ajout de la couche d'attention (Bahdanau) | 25 |
| 5.1 | Introduction | 25 |
| 5.2 | Importation des bibliothèques | 25 |
| 5.3 | Prédéfnition de classe | 26 |
| 5.4 | Construire et tester les modèles | 26 |
| 5.4.1 | Ajout de la couche aux modèle avec embedding pendant l'entraînement | 26 |
| 5.4.2 | Ajout de la couche aux modèle pré-entraîner | 28 |
| 5.5 | Visualiser les résultats d'entraînement | 30 |
| 5.6 | Conclusion | 31 |

Chapitre 1

Présentation du problème

Contexte et motivation

Les avis en ligne (films, produits, restaurants) constituent une source d'information clé pour la prise de décision. Automatiser leur **analyse de sentiments** permet de résumer rapidement l'opinion globale et d'identifier des signaux forts (satisfaction, insatisfaction). Deux grandes familles de représentations sont couramment opposées : (i) les *représentations creuses* de type TF-IDF (sac de mots, n-grammes), (ii) les *représentations denses* apprises par des réseaux (embeddings + RNN/LSTM), éventuellement **pondérées par une attention** pour la mise en évidence des tokens discriminants.

Problématique

Nous cherchons à **prédire la polarité** d'un avis (*positif / négatif*) et à **comparer** des approches basées sur :

- une **ligne de base TF-IDF** couplée à des classifieurs linéaires (**Régression Logistique, SVM**) ;
- un **RNN (LSTM)** avec embeddings (appris pendant l'entraînement ou *pré-entraînés* type FastText) ;
- un **LSTM + attention (Bahdanau)**, pour améliorer la performance et l'interprétabilité (poids d'attention).

Approche et protocole d'évaluation

- **Prétraitement** : tokenisation, normalisation légère ; pour TF-IDF, choix des n-grammes et filtrage de fréquence ; pour RNN, *padding* à une longueur fixe.

- **Entraînement** : Split dédié (train/validation/test), *early stopping*, recherche légère d'hyperparamètres (C , n -grammes, `min_df`/`max_df`, taille LSTM, dropout).
- **Mesures** : **accuracy** et **F1**, avec rapports de classification et matrices de confusion ; pour les modèles probabilistes, courbes ROC/PR et AUC.
- **Interprétabilité** : inspection des poids TF-IDF (features les plus informatives) et **visualisation des poids d'attention** pour LSTM+Bahdanau.

Contributions

1. **Implémentation d'une baseline** TF-IDF + (Régression Logistique / SVM) et comparaison chiffrée.
2. **RNN (LSTM)** avec deux régimes d'embeddings : *appris pendant l'entraînement* vs *pré-entraînés* (FastText gelés / éventuellement ajustés).
3. **Ajout d'une attention (Bahdanau)** au-dessus du LSTM ; analyse qualitative des mots surpondérés.
4. **Synthèse des résultats** (tableaux, courbes, learning curves) et discussion des compromis performance/interprétabilité/coût.

Données et outils

- **Dataset** : Allociné (FR) <https://huggingface.co/datasets/tblard/allocine> ; splits fournis (train, validation, test).
- **Outils** : Python, scikit-learn (TF-IDF, LR, SVM, métriques), TensorFlow/Keras (LSTM, attention), Matplotlib (visualisation).
- **Reproductibilité** : *early stopping*, graines aléatoires, notebooks .ipynb fournis.

Organisation du rapport

- **Chapitre 1** : Problématique et protocole (présent chapitre).
- **Chapitre 2** : Présentation des données utilisées.
- **Chapitre 3** : Baseline TF-IDF + modèles linéaires (réglages, résultats, interprétations).
- **Chapitre 4** : RNN (LSTM) avec embeddings (appris vs pré-entraînés), comparaisons.
- **Chapitre 5** : LSTM + **attention de Bahdanau**, résultats et visualisations des poids.

Chapitre 2

Présentation des données utilisées

2.1 Chargement des données

Les données comportent un **train**, une **validation** et un **test**. On standardise les noms de colonnes et on ajoute une étiquette textuelle.

```
1 # 1) Charger le dataset Allociné (FR)
2 ds = load_dataset("tblard/allocine")
3 # 2) Convertir en DataFrame pandas
4 df_train = pd.DataFrame(ds["train"])
5 df_valid = pd.DataFrame(ds["validation"])
6 df_test = pd.DataFrame(ds["test"])
```

2.2 Traitement des données

Nous analysons dans cette partie, la taille du jeu de données, les valeurs manquantes, les doublons, et le mapping des étiquettes.

2.2.1 Quelques statistiques

- **Jeu de données** : le dataset Allociné contient **3 sous-ensembles** (train, validation, test), chacun avec **2 colonnes**.
- **Tailles** :
 1. Entraînement (train) : $160\,000 \times 2$
 2. Validation (validation) : $20\,000 \times 2$
 3. Test (test) : $20\,000 \times 2$

- Colonnes :
 - text : avis (film, produit, restaurant).
 - label : étiquette binaire initiale $\in \{negative, positif\}$.
- Répartition des données d'entraînement :

```
label
1    80334
0    79127
Name: count, dtype: int64
```

→ Les données sont considérées **équilibrées**.

2.2.2 Encodage de la colonne label

Nous allons convertir les labels textuels en labels numériques.

```
1 # Mapper les labels entiers vers des libellés texte
2 label_map = {0: "negative", 1: "positive"}
3 for df in (df_train, df_valid, df_test):
4     df["label_text"] = df["label"].map(label_map)
```

Résultat du mapping (label_text → label) :

| | text | label | label_text |
|---|---|-------|------------|
| 0 | Si vous cherchez du cinéma abrutissant à tous ... | 0 | negative |
| 1 | Trash, re-trash et re-re-trash...! Une horreur... | 0 | negative |
| 2 | Et si, dans les 5 premières minutes du film, l... | 0 | negative |
| 3 | Mon dieu ! Quelle métaphore filée ! Je suis ab... | 0 | negative |
| 4 | Premier film de la saga Kozure Okami, "Le Sabr... | 1 | positive |
| 5 | L'amnésie est un thème en or pour susciter le ... | 0 | negative |
| 6 | Tout commence comme une comédie légère avant u... | 1 | positive |
| 7 | un excellent film qui merite ses quatre étoile... | 1 | positive |
| 8 | Deuxième long métrage de Pasolini, Mamma Roma ... | 1 | positive |
| 9 | Créateur de la célèbre série télévisée Kaamel... | 0 | negative |

2.2.3 Gestion des doublons

Nous allons faire la vérification dans les données d'entraînement :

```
1 # Le nombre de lignes dupliquées
2 print(df_train.duplicated(subset=["label", "text"]).sum())
3 # Suppression
4 df_train = df_train.drop_duplicates(subset="text").reset_index(drop=
    True)
```

Résultat de suppression dans les données d'entraînement : (160 000 → 159 461)

Le nombre de lignes dupliquées est : **534**

2.2.4 Gestion des valeurs manquantes

```
1 # valeurs manquantes
2 print(df_train["label"].isna().sum())
```

Alors, il n'existe pas de valeurs manquantes.

Chapitre 3

Première partie : TF-IDF

3.1 Introduction

L'approche **TF-IDF** (Term Frequency–Inverse Document Frequency) transforme chaque avis en un vecteur qui pondère l'importance des mots (et éventuellement des n-grammes). Cette représentation *bag-of-words* ignore l'ordre exact des mots mais offre :

- une **base robuste et interprétable** : poids des termes, inspection des features;
- un **apprentissage rapide** avec des modèles linéaires (Régression Logistique, SVM);
- une **forte ligne de base** pour comparer ensuite des modèles séquentiels.

3.2 Importation des bibliothèques

- Chargement des données (Hugging Face Datasets)

```
1 from datasets import load_dataset
```

- Prétraitement texte / Pipeline

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.pipeline import Pipeline
```

- Modèles (baseline)

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.svm import LinearSVC
3 from sklearn.neighbors import KNeighborsClassifier
```

- Évaluation

```
1 from sklearn.metrics import (
2     accuracy_score, f1_score, precision_score, recall_score,
3     classification_report, confusion_matrix, precision_recall_curve,
4     roc_curve, auc, average_precision_score )
```

- Analyse / Visualisation

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
```

3.3 Prédéfinition des fonctions

Dans cette partie, nous présentons les fonctions utilitaires retenues pour : l'évaluation des modèles, la visualisation des résultats et des courbes, la comparaison systématique de plusieurs architectures.

Remarque

Nous présentons ici uniquement un extrait du code. Le code complet est disponible dans le fichier `.ipynb`.

- **Fonction d'évaluation** : Calculer l'accuracy, f1-score, la précision et le recall

```
1 def evaluate_and_report(y_true, y_pred, title=""):
2     acc = accuracy_score(y_true, y_pred)
3     f1 = f1_score(y_true, y_pred)
4     p = precision_score(y_true, y_pred)
5     r = recall_score(y_true, y_pred)
6     ...
```

- **Fonction de visualisation de la matrice de confusion** :

```
1 def plot_confusion(y_true, y_pred, labels=("neg", "pos"), title="
2     Confusion Matrix"):
3     cm = confusion_matrix(y_true, y_pred, labels=[0,1])
4     fig, ax = plt.subplots(figsize=(4.5,4))
5     im = ax.imshow(cm, interpolation="nearest")
6     ax.figure.colorbar(im, ax=ax)
7     ...
```

- **Fonction** de visualisation de la **courbe ROC**

```
1 def plot_pr_roc(pipe, X_test, y_test, title_prefix=""):
2     ...
```

- **Fonction** de visualisation de la **courbe d'apprentissage**

```
1 def plot_learning_curve(pipe, X_train, y_train, cv_splits=5, title=
    "Learning curve"):
2     cv = StratifiedKFold(n_splits=cv_splits, shuffle=True,
        random_state=42)
3     fig, ax = plt.subplots(figsize=(5.8,4.2))
4     ...
```

3.4 Construire et tester les modèles

Dans cette section, nous allons nous limiter à la construction de trois modèles avec des paramètres différents :

1. **Régression logistique (LR)**
 - Sans régularisation
 - Avec régularisation **L2 (Ridge)** et **L1 (Lasso)**
2. **Machine à vecteurs de support (SVM)**
3. **k plus proches voisins (k-NN)**

3.4.1 Régression Logistique (sans régularisation)

1. Construction du pipeline

```
1 pipe_lr = Pipeline([
2     ("tfidf", TfidfVectorizer(
3         lowercase=True,
4         strip_accents="unicode",
5         ngram_range=(1,2), # unigram+bigram
6         min_df=2,
7         max_df=0.98,
8         sublinear_tf=True)),
9     ("clf", LogisticRegression(max_iter=200, C=2.0, n_jobs=None))
10 ])
```

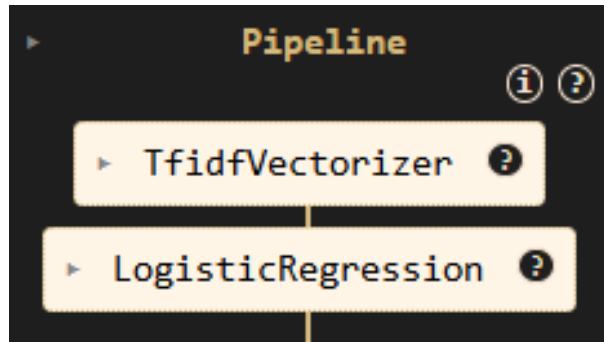


Figure 3.1 – Pipeline du modèle

2. Entraîner le modèle

```
1 pipe_lr.fit(df_train["text"], df_train["label"])
```

3. Prédiction des données test

```
1 pred_lr = pipe_lr.predict(df_test["text"])
```

4. Évaluer le modèle

```
1 res_lr = evaluate_and_report(df_test["label"], pred_lr, title="
  TF-IDF + LogisticRegression")
2 plot_confusion(df_test["label"], pred_lr, title="CM - TFIDF+LR")
```

5. Résultats :

```

=== TF-IDF + LogisticRegression ===
Accuracy: 0.9426 | F1: 0.9405 | Precision: 0.9358 | Recall: 0.9453

Classification report:
      precision    recall  f1-score   support

     0       0.9491     0.9402     0.9446     10408
     1       0.9358     0.9453     0.9405      9592

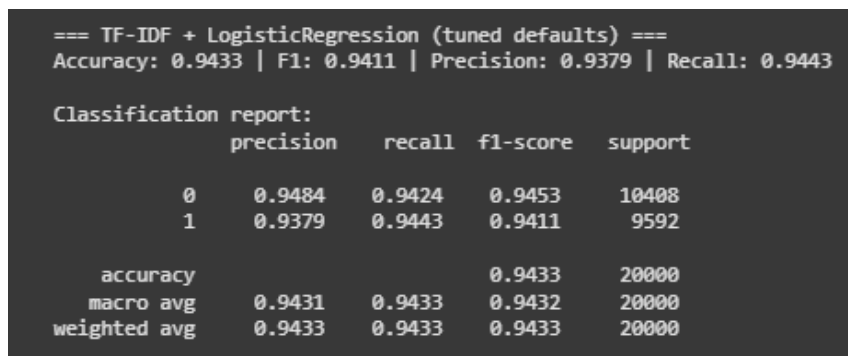
   accuracy          0.9426     20000
  macro avg       0.9424     0.9428     0.9426     20000
 weighted avg       0.9427     0.9426     0.9427     20000
  
```

⇒ Le modèle Logistic Regression donne de **bonnes perforamnaces** en termes d'accuracy, de recall et de precision.

3.4.2 Ridge

Comparons à une régression linéaire simple, nous allons simplement modifier l'étape "clf" (classifieur) afin de prendre en compte la régularisation Ridge.

```
1 pipe_lr = Pipeline([
2     ("tfidf", TfidfVectorizer(...)),
3     ("clf", LogisticRegression(
4         solver="liblinear",    # solide pour texte binaire
5         penalty="l2",          # Ridge
6         C=2.0,                  # regularisation un peu plus souple
7         class_weight="balanced", # retire si classes parfaitement
                                # equilibrees
8         max_iter=5000,
9         n_jobs=-1              # utilise par liblinear si multiclass
10    ))
11 ])
```



```
=== TF-IDF + LogisticRegression (tuned defaults) ===
Accuracy: 0.9433 | F1: 0.9411 | Precision: 0.9379 | Recall: 0.9443

Classification report:
              precision    recall  f1-score   support

     0       0.9484        0.9424    0.9453       10408
     1       0.9379        0.9443    0.9411        9592

 accuracy          0.9433                0.9433       20000
 macro avg         0.9431        0.9433    0.9432       20000
 weighted avg      0.9433        0.9433    0.9433       20000
```

⇒ Comparée à la régression logistique simple (accuracy $\approx 0,942$), la régularisation Ridge donne des résultats légèrement meilleurs (accuracy $\approx 0,943$), mais l'écart est négligeable.

3.4.3 Lasso

Comparons à une régression linéaire simple, nous allons simplement modifier l'étape clf (classifieur) afin de prendre en compte la régularisation Lasso.

```
1 pipe_lr = Pipeline([
2     ("tfidf", TfidfVectorizer(...)),
3     ("clf", LogisticRegression(
4         solver="liblinear",    # requis pour L1 (ou 'saga')
5         penalty="l1",          # LASSO
6         C=2.0,
7         class_weight="balanced", # enleve si classes equilibrees
8         max_iter=5000,
9         n_jobs=-1))    ])
```

```

10 pipe_lr.fit(df_train["text"], df_train["label"])
11 pred_lr = pipe_lr.predict(df_test["text"])
12 res_lr = evaluate_and_report(df_test["label"], pred_lr, title="TF-IDF
    + LogisticRegression (L1 - Lasso)")

```

```

=== TF-IDF + LogisticRegression (L1 - Lasso) ===
Accuracy: 0.9413 | F1: 0.9391 | Precision: 0.9351 | Recall: 0.9431

Classification report:

```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.9471 | 0.9397 | 0.9434 | 10408 |
| 1 | 0.9351 | 0.9431 | 0.9391 | 9592 |
| accuracy | | | 0.9413 | 20000 |
| macro avg | 0.9411 | 0.9414 | 0.9412 | 20000 |
| weighted avg | 0.9413 | 0.9413 | 0.9413 | 20000 |

⇒ La régularisation Lasso (accuracy $\approx 0,941$) n'améliorent pas les résultats de la régression logistique (accuracy $\approx 0,942$).

3.4.4 SVM

Nous utilisons la fonction LinearSVC pour ce modèle.

```

1 pipe_svm = Pipeline([
2     ("tfidf", TfidfVectorizer(...)),
3     ("clf", LinearSVC(C=1.0))
4 ])
5 pipe_svm.fit(df_train["text"], df_train["label"])
6 pred_svm = pipe_svm.predict(df_test["text"])
7 res_svm = evaluate_and_report(df_test["label"], pred_svm, title="TF-
    IDF + LinearSVM")

```

```

=== TF-IDF + LinearSVM ===
Accuracy: 0.9455 | F1: 0.9434 | Precision: 0.9388 | Recall: 0.9481

Classification report:

```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.9517 | 0.9430 | 0.9473 | 10408 |
| 1 | 0.9388 | 0.9481 | 0.9434 | 9592 |
| accuracy | | | 0.9455 | 20000 |
| macro avg | 0.9452 | 0.9456 | 0.9454 | 20000 |
| weighted avg | 0.9455 | 0.9455 | 0.9455 | 20000 |

⇒ Le modèle SVM est le meilleur en termes de performance (accuracy $\approx 0,955$), recall et précision en comparaison avec la régression logistique ((accuracy $\approx 0,945$)).

3.4.5 KNN

Nous utilisons la fonction `KNeighborsClassifier` pour ce modèle avec les hyper-paramètres :

- la distance **cosinus** : recommandé en texte
- le nombre des **k** plus proches voisins : $k=5$

```
1 pipe_knn = Pipeline([
2     ("tfidf", TfidfVectorizer(...)),
3     ("clf", KNeighborsClassifier(
4         n_neighbors=5,          # k de d part
5         metric="cosine",       # distance cosinus
6         algorithm="brute",     # n cessaire pour 'cosine'
7         weights="distance",
8         n_jobs=-1
9     ))
10 ])
11 pipe_knn.fit(df_train["text"], df_train["label"])
12 pred_knn = pipe_knn.predict(df_test["text"])
13 res_knn = evaluate_and_report(df_test["label"], pred_knn, title="TF-
    IDF + KNN (cosine)")
```

```
=== TF-IDF + KNN (cosine) ===
Accuracy: 0.8316 | F1: 0.8289 | Precision: 0.8082 | Recall: 0.8508

Classification report:

```

| | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0 | 0.8555 | 0.8139 | 0.8342 | 10408 |
| 1 | 0.8082 | 0.8508 | 0.8289 | 9592 |
| accuracy | | | 0.8316 | 20000 |
| macro avg | 0.8318 | 0.8324 | 0.8316 | 20000 |
| weighted avg | 0.8328 | 0.8316 | 0.8317 | 20000 |

⇒ KNN (accuracy $\approx 0,83$) est nettement moins bon que tes modèles linéaires (accuracy $\approx 0,95$).

En texte TF-IDF (très haute dimension, très clairsemé), KNN souffre de :

- temps d'exécution très important
- malédiction de la dimension → toutes les distances se ressemblent ;
- mauvaise pertinence des distances (si pas cosinus) ;
- coût d'inférence élevé (comparaison à tous les docs).

3.5 Améliorations possibles

- **Augmenter la taille des données :**

Nous allons entraîner le meilleur modèle (SVM) avec les données d'entraînement et de validation. Nous pouvons par la suite l'entraîner sur la totalité des données (avec les données test) :

1. Concaténer les données train + valid

```
1 df_tv = pd.concat([df_train, df_valid], ignore_index=True)
```

2. Redéfinir le pipeline SVM

```
1 pipe_svm = Pipeline([
2     ("tfidf", TfidfVectorizer(
3         lowercase=True, strip_accents="unicode",
4         ngram_range=(1,2), min_df=2, max_df=0.98, sublinear_tf=
5         True
6     )),
7     ("clf", LinearSVC(C=1.0))
8 ])
```

3. Entraîner sur train + valid

```
1 pipe_svm.fit(df_tv["text"], df_tv["label"])
```

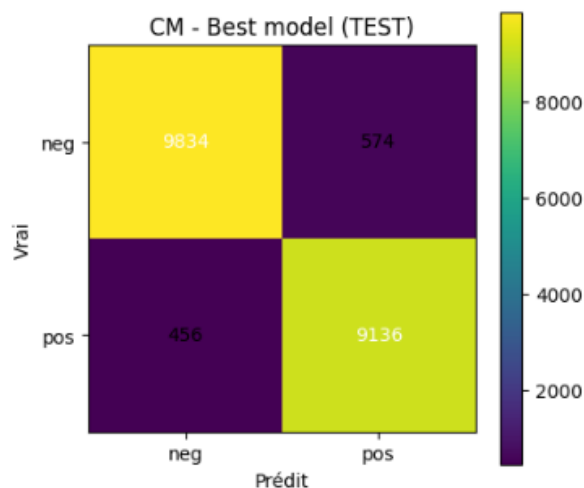


0.9485

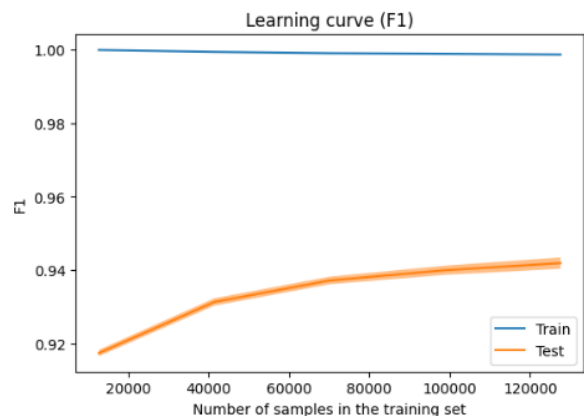
⇒ En ajoutant 20 000 exemples, l'accuracy passe de 0,9455 à 0,9485, soit un gain absolu de +0,0030 (= +0,30 point de pourcentage) et un gain relatif d'environ $\frac{0,9485-0,9455}{0,9455} \approx 0,317\%$. L'amélioration est donc **faible** au regard du coût/complexité d'augmenter les données.

3.6 Visualiser les résultats d'entraînement

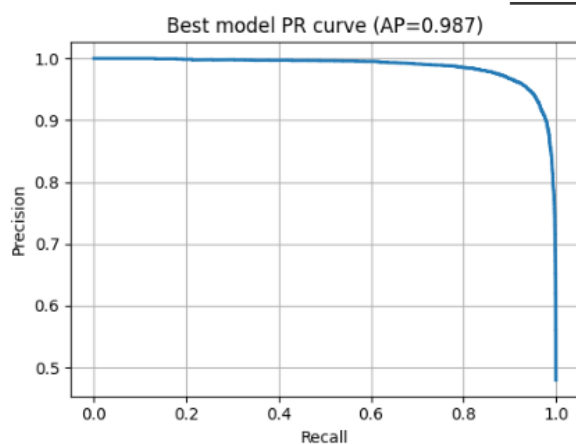
Courbe PR (AP $\approx 0,986$). Sur un jeu *quasi équilibré*, une baseline aléatoire tournerait autour de 0,5 d'AP. Ici, la courbe reste très proche de la précision = 1 pour une large plage de rappels, et ne chute qu'à l'approche de Rec $\rightarrow 1$, où l'on accepte davantage de *faux positifs* pour capter les derniers vrais positifs. Cela indique :



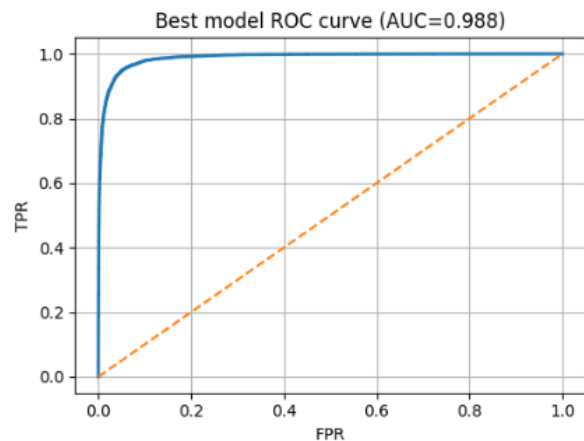
Matrice de confusion



Courbes d'entraînement



Courbe ROC



Courbe PR / AUC

- un modèle très **précis** dès que le seuil n'est pas trop bas ;
- une **dégradation contrôlée** de la précision uniquement en zone de rappel extrême.

Courbe ROC (AUC $\approx 0,987$). Un ROC-AUC de $\sim 0,987$ signifie qu'en tirant au hasard un exemple positif et un négatif, le classifieur (ici SVM) classe le positif au-dessus du négatif environ **98,7%** du temps. C'est un **pouvoir de séparation** très élevé et cohérent avec l'AP.

Courbes d'apprentissage. La courbe d'apprentissage montre des performances **proches entre train et validation/test** avec des niveaux élevés : pas de surapprentissage notable. *Bias/variance* raisonnable : ajouter plus de données devrait encore lisser la variance ; l'AP/ROC déjà élevés indiquent un **faible biais**.

3.7 Conclusion

Avec une représentation TF-IDF (uni+bi-grammes) et des modèles linéaires, nous obtenons une ligne de base solide. Parmi les classifieurs testés, **le SVM linéaire** ressort comme le plus performant en accuracy, devant la régression logistique et ses variantes L2/L1 ; **KNN** reste moins compétitif dans cet espace haute dimension.

Table 3.1 – Accuracy (jeu de test) des modèles TF-IDF

| Modèle | Accuracy |
|----------------------------|---------------|
| Régression logistique (LR) | 0.9426 |
| Ridge (Logistique + L2) | 0.9433 |
| Lasso (Logistique + L1) | 0.9413 |
| SVM linéaire | 0.9455 |
| KNN | 0.8316 |

L'écart entre **Ridge** et **LR** est minime (gain $\approx 0,1$ point), tandis que **SVM** apporte le meilleur compromis biais/variance sur TF-IDF.

Les gains restants proviennent surtout du **tuning** (C, n -grammes, min_df/max_df) et d'une sélection d'hyperparamètres plus fine, avant de passer aux modèles séquentiels (RNN/LSTM).

Chapitre 4

Deuxième partie : RNN

4.1 Introduction

Contrairement à TF-IDF, qui représente les avis comme des sacs de mots indépendants, un **RNN** à base de **LSTM** exploite l'*ordre* des tokens et les *dépendances longues* (négations, portées syntaxiques, co-références).

Après **tokenisation** et **padding**, chaque mot est projeté dans un espace continu via des **embeddings** (pré-entraînés ou appris).

Le **LSTM** agrège séquentiellement l'information et produit une représentation riche du message, mieux adaptée aux structures linguistiques que TF-IDF. Nous évaluerons ce modèle sur les mêmes jeux de données.

4.2 Importation des bibliothèques

- Prétraitement du texte

```
1 from tensorflow.keras.preprocessing.text import Tokenizer
2 from tensorflow.keras.preprocessing.sequence import pad_sequences
```

- Entraînement

```
1 from tensorflow.keras.callbacks import EarlyStopping
```

- Visualisation

```
1 import matplotlib.pyplot as plt
```

- Modélisation

```

1 from tensorflow.keras.models import Sequential
2 from tensorflow.keras.layers import (           # couches
3     Embedding,
4     LSTM,                                     # RNN LSTM
5     Dense,                                   # sortie
6     Dropout )                               # regularisation

```

Prédéfinition des fonctions

```

1 def plot_training_history(history, model_name="Mod le"):
2     ...
3     plt.subplot(1,2,1)
4     plt.plot(history_dict['loss'], label='Entra nement')
5     plt.plot(history_dict['val_loss'], label='Validation')

```

4.3 Tokenisation

Le **tokenizer** apprend le vocabulaire et transforme chaque avis en séquence d'indices entiers. *Ici, nous avons appris le vocabulaire sur `df_train['text']`, converti les avis en séquences d'indices, puis normalisé leur longueur par **padding/troncature** à 150 tokens (même traitement pour validation et test).*

```

1 vocab_size = 20000
2 max_len = 150
3 embedding_dim = 128
4
5 tokenizer = Tokenizer(num_words=vocab_size, oov_token="<OOV>")
6 tokenizer.fit_on_texts(df_train['text'])
7
8 X_train = tokenizer.texts_to_sequences(df_train['text'])
9 X_train = pad_sequences(X_train, maxlen=max_len, padding="post",
10                          truncating="post")
11 #... faire de meme pour X_test et X_valid

```

Récupérer les targets ('label') correspondants des 3 dataframes :

```

1 y_val= np.array(df_valid['label'])
2 y_train = np.array(df_train['label'])
3 y_test = np.array(df_test['label'])

```

4.4 Construire et tester les modèles

Nous allons évaluer deux types d'*embeddings* :

- **Embeddings appris pendant l'entraînement** (learned from scratch);
- **Embeddings pré-entraînés** (p. ex. FastText).

4.4.1 Embedding pendant l'entraînement

1. Construire le modèle

```
1 model_1 = Sequential([
2     Embedding(input_dim=vocab_size, output_dim=embedding_dim,
3               input_length=max_len),
4     LSTM(128, return_sequences=False),
5     Dropout(0.3),
6     Dense(1, activation="sigmoid")
7 ])
```

2. Compiler le modèle

```
1 model_1.compile(loss="binary_crossentropy", optimizer="adam",
2                 metrics=["accuracy"])
3 model_1.summary()
```

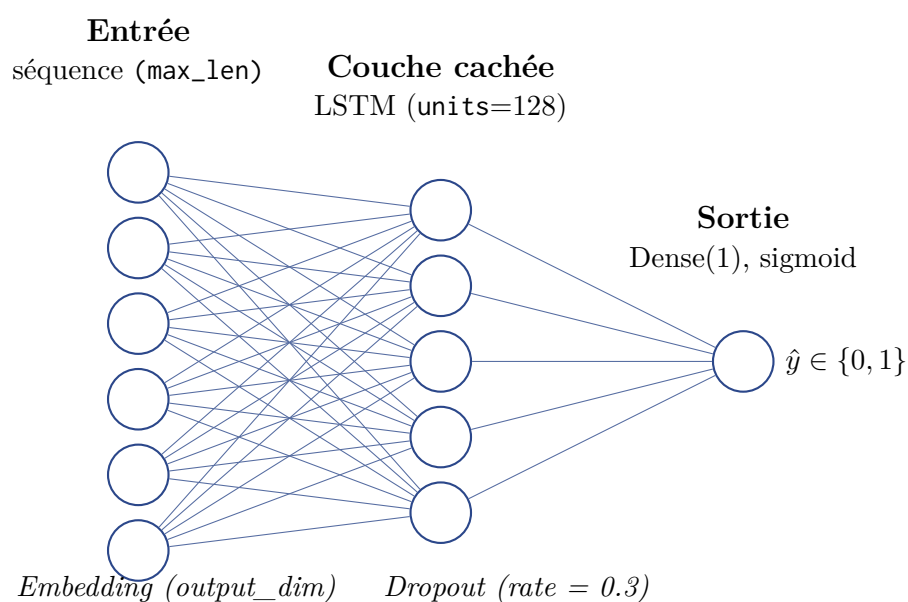


Figure 4.1 – Visualisation du réseau : Entrée \rightarrow LSTM(128) \rightarrow Dense(sigmoid).

3. Configurer l'arrêt anticipé (*Early Stopping*)

```
1 early_stop = EarlyStopping(  
2     monitor='val_loss',  
3     patience=3,      # arrete si pas amelioration pendant 3 epochs  
4     restore_best_weights=True  
5 )
```

4. Entraîner du modèle avec 50 epochs

```
1 history_1 = model_1.fit(  
2     X_train, y_train,  
3     validation_data=(X_val, y_val),  
4     epochs=50,  
5     callbacks=[early_stop],  
6     batch_size=128,  
7     verbose=1  
8 )
```

Nous avons choisi **128** comme taille du batch car les données sont **massives**.

5. Évaluer le modèle

```
1 loss, acc = model_1.evaluate(X_test, y_test, verbose=0)  
2 print(f"Accuracy sur test : {acc:.2f}")
```



```
Accuracy sur test : 0.93
```

6. Exemple d'inférence

```
1 #exemple  
2 new_text = ["Ce film est g nial , un chef d'oeuvre incroyable!"]  
3 # tokenisation  
4 seq = tokenizer.texts_to_sequences(new_text)  
5 padded_seq = pad_sequences(seq, maxlen=max_len, padding="post")  
6 # prediction  
7 prediction = model_1.predict(padded_seq)  
8 print("Probabilit positif :", prediction[0][0])  
9 print("Classe pr dite :", int(prediction[0][0] > 0.5))
```

```
1/1 ————— 0s 248ms/step
Probabilité positif : 0.980217
Classe prédite : 1
```

4.4.2 Embedding préentraîner avec "fastText facebook 2016"

1. Charger des embeddings préentraînés

```
1 import gensim.downloader as api
2 # Charger FastText fran aïs
3 ft_model = api.load("fasttext-wiki-news-subwords-300") #
    vecteurs 300d
```

2. **Créer la matrice d'embedding.** Nous construisons `embedding_matrix_FastText` de taille `num_words × embedding_dim` (par ex. 50 000 × 300).

Chaque ligne i contient le vecteur de l'*embedding* du mot d'indice i (selon le Tokenizer).

```
1 # Hypers
2 embedding_dim = 300 # dimension des vecteurs FastText
3 num_words = min(vocab_size, len(tokenizer.word_index) + 1)
4 # Créer la matrice d'embedding
5 embedding_matrix_FastText = np.zeros((num_words, embedding_dim))
6 # Remplissage depuis FastText
7 for word, i in tokenizer.word_index.items():
8     if i >= num_words:
9         continue
10    w = word.lower()
11    if w in ft_kv.key_to_index: # present dans vocab FastText
12        embedding_matrix_FastText[i] = ft_kv.get_vector(w)
13    # sinon: on garde la valeur alatoire initialis e (00V)
```

3. Configurer l'arrêt anticipé (*Early Stopping*)

```
1 early_stop = EarlyStopping(
2     monitor='val_loss',
3     patience=3, # arrete si pas amelioration pendant 3 epochs
4     restore_best_weights=True
5 )
```

4. Construire le modèle

```

1 model_fasText = Sequential([
2     Embedding(input_dim=num_words,
3               output_dim=embedding_dim,
4               weights=[embedding_matrix_FastText],
5               input_length=max_len,
6               trainable=False),
7     LSTM(128),
8     Dropout(0.3),
9     Dense(1, activation='sigmoid')
10 ])

```

5. Compiler le modèle

```

1 model_fasText.compile(loss='binary_crossentropy', optimizer='
  adam', metrics=['accuracy'])
2 model_fasText.summary()

```

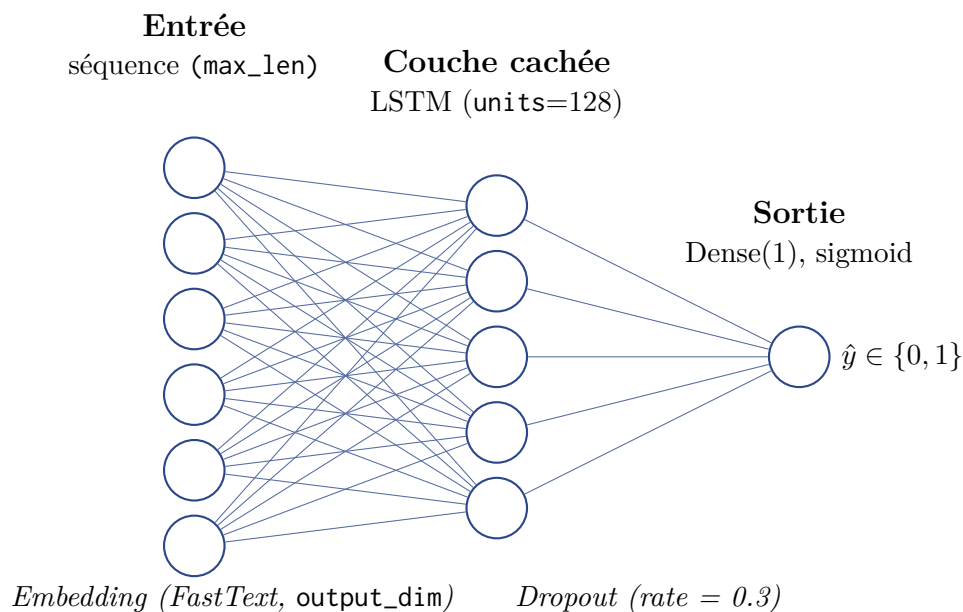


Figure 4.2 – Réseau : Entrée → **Embedding FastText (gelé)** → LSTM(128) → Dropout(0.3) → Dense(sigmoid).

6. Entraîner le modèle avec 50 epochs

```


1 # Entraînement
2 history_2=model_fasText.fit(X_train, y_train, validation_data=(
3     X_val, y_val),callbacks=[early_stop], epochs=50, batch_size
4     =128
5 )

```


7. Évaluer le modèle

```
1 loss, acc = model_fasText.evaluate(X_test, y_test, verbose=0)
2 print(f"Accuracy sur test : {acc:.2f}")
```

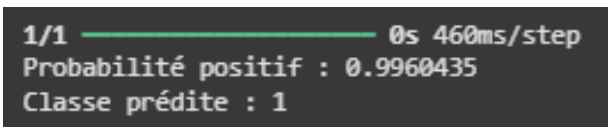
Nous avons choisi **128** comme taille du batch car les données sont **massives**.



```
Accuracy sur test : 0.88
```

8. Exemple d'inférence

```
1 #exemple
2 new_text = ["Ce film est g nial , un chef d'oeuvre incroyable!"]
3 #tokenisation
4 seq = tokenizer.texts_to_sequences(new_text)
5 padded_seq = pad_sequences(seq, maxlen=max_len, padding="post")
6 prediction = model_fasText.predict(padded_seq)
7 #pr diction
8 print("Probabilit positif :", prediction[0][0])
9 print("Classe pr dite :", int(prediction[0][0] > 0.5))
```



```
1/1 ————— 0s 460ms/step
Probabilité positif : 0.9960435
Classe prédite : 1
```

4.5 Visualiser les résultats d'entraînement

```
1 plot_training_history(history_1, model_name="Mod le avec
   embedding pendant l'entraînement")
2 plot_training_history(history_2, model_name="Mod le avec
   pr entraînement")
```

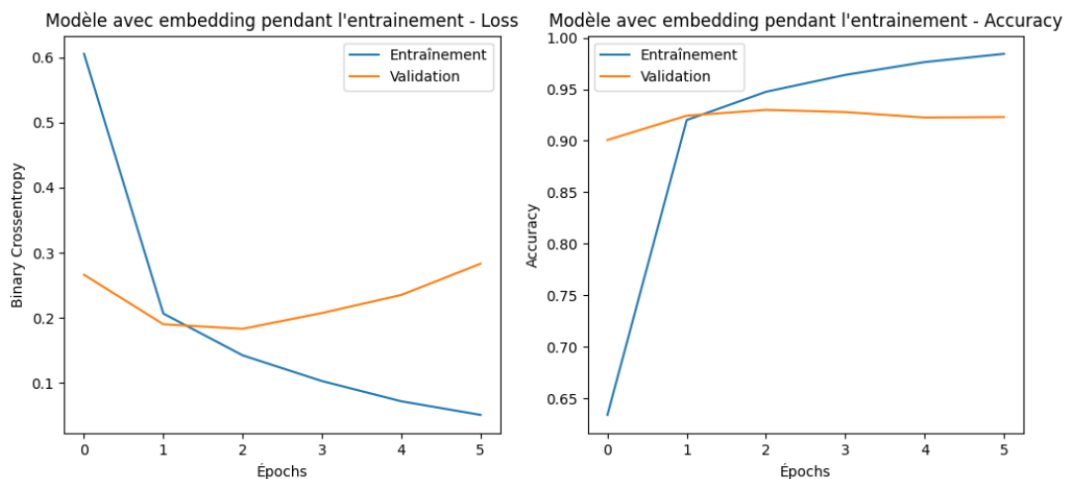


Figure 4.3 – Historique d'entraînement — *embeddings appris sur nos données*.

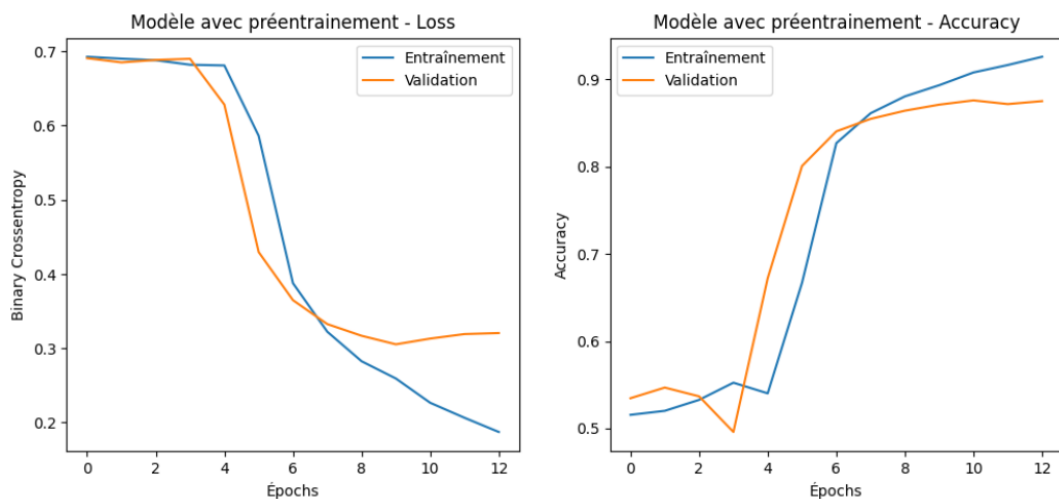


Figure 4.4 – Historique d'entraînement — *embeddings pré-entraînés*.

Figure 4.5 – Comparaison des historiques d'entraînement : embeddings appris vs pré-entraînés.

Interprétation

- Les deux modèles sont **performants** : la loss de validation baisse et l'accuracy reste élevée.
- **Embeddings appris** : convergence **plus rapide** et **plus stable**.
- **Embeddings pré-entraînés** : démarrage plus lent, mais **rattrapent** un niveau proche après quelques époques.
- Choix pratique : privilégier les **embeddings appris** si l'on a assez de données ; envisager le **pré-entraîné + fine-tuning** si les données sont limitées.

4.6 Conclusion

Dans nos expériences, le modèle **avec embeddings appris pendant l'entraînement** obtient de meilleurs scores que le modèle **initialisé par des embeddings pré-entraînés**.

Table 4.1 – Accuracy (jeu de test) des modèles LSTM

| Modèle | Accuracy |
|----------------------------------|-------------|
| Embedding pré-entraîner | 0.88 |
| Embedding pendant l'entraînement | 0.93 |

- **Embeddings appris pendant l'entraînement** : meilleurs lorsque l'on dispose d'assez de données supervisées et que le domaine est spécifique. Ils s'adaptent finement au jeu de données, au prix d'un apprentissage potentiellement plus long.
- **Embeddings pré-entraînés** : utiles en *données faibles*, pour initialiser avec des connaissances lexicales générales. Leur plein potentiel apparaît souvent avec un **fine-tuning contrôlé**.

Dans notre cadre expérimental, les **embeddings appris pendant l'entraînement** sont **légèrement supérieurs** aux **embeddings pré-entraînés** *non finement ajustés*.

Chapitre 5

Troisième partie : Ajout de la couche d'attention (Bahdanau)

5.1 Introduction

Les RNN (LSTM) agrègent séquentiellement l'information, mais ne permettent pas toujours d'identifier explicitement les *mots/segments* qui portent le signal de polarité. L'**attention additive** (dite *Bahdanau*) pondère les sorties $\{h_t\}$ du LSTM par des poids α_t appris, de sorte que la représentation de phrase

$$\mathbf{c} = \sum_t \alpha_t \mathbf{h}_t \quad \text{avec} \quad \alpha_t = \text{softmax}(\mathbf{v}^\top \tanh(\mathbf{W}\mathbf{h}_t + \mathbf{b}))$$

mette en avant les tokens les plus discriminants pour la prédiction.

5.2 Importation des bibliothèques

- Modélisation

```
1 import tensorflow as tf
2 from tensorflow.keras.layers import Input, Embedding, LSTM, Dropout
  , Dense, Multiply, Lambda
3 from tensorflow.keras.models import Model
4 from tensorflow.keras.callbacks import EarlyStopping
```

- Évaluation

```
1 import numpy as np
2 from sklearn.metrics import accuracy_score, f1_score,
  classification_report, confusion_matrix
```

- Backend / utilitaires Keras (bas niveau)

```
1 from tensorflow.keras import backend as K # ops bas niveau: tf.
   reduce_sum, softmax, int_shape, etc.
```

5.3 Prédéfinition de classe

Définir la classe BahdanauAttention

```
1 class BahdanauAttention(Layer):
2     def __init__(self, units):
3         super(BahdanauAttention, self).__init__()
4         self.W1 = tf.keras.layers.Dense(units) # pour hidden states
5         self.W2 = tf.keras.layers.Dense(units) # pour query
6         self.V = tf.keras.layers.Dense(1)      # pour score
7     def call(self, values, query):
8         ...
9         return context_vector, attention_weights
10    return context
```

C'est une couche d'attention additive (Bahdanau) qui, à partir de l'état final du LSTM (query) et de toutes ses sorties temporelles (values), calcule des poids d'attention sur les pas de temps via un score appris, puis retourne la somme pondérée des sorties (vecteur contexte) ainsi que les poids eux-mêmes.

5.4 Construire et tester les modèles

5.4.1 Ajout de la couche aux modèle avec embedding pendant l'entraînement

1. Construire le modèle : LSTM + Attention de Bahdanau

```
1 # entr e
2 inputs = Input(shape=(max_len,))
3 embedding = Embedding(input_dim=vocab_size, output_dim=
   embedding_dim, input_length=max_len)(inputs)
4
5 # LSTM avec return_sequences=True pour l'attention et
   return_state=True pour r cup rer l' tat final
```

```

6 lstm_out, state_h, state_c = LSTM(128, return_sequences=True,
   return_state=True)(embedding)
7
8 # Couche Bahdanau attention (query = tat final du LSTM)
9 attention_layer = BahdanauAttention(units=128)
10 att_out, att_weights = attention_layer(lstm_out, state_h)
11
12 # Dropout et Dense final
13 drop = Dropout(0.3)(att_out)
14 outputs = Dense(1, activation='sigmoid')(drop)
15
16 # Cr ation du mod le
17 model_bahdanau_1 = Model(inputs, outputs)

```

2. Compiler

```

1 model_bahdanau_1.compile(loss='binary_crossentropy', optimizer='
   adam', metrics=['accuracy'])
2 model_bahdanau_1.summary()

```

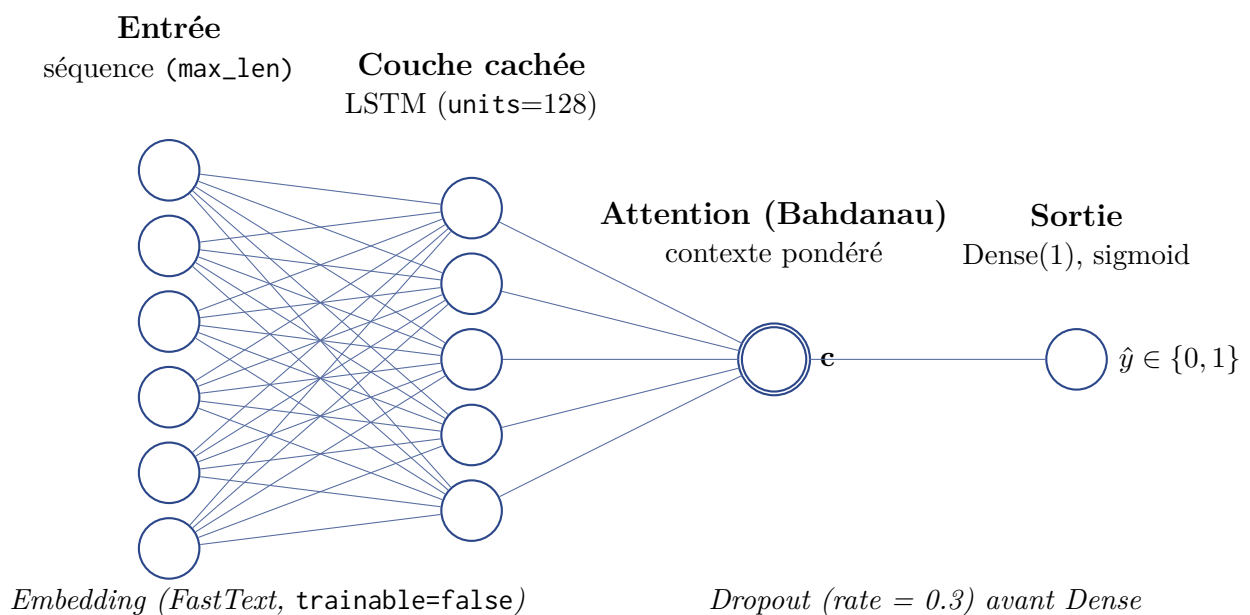


Figure 5.1 – Réseau : Entrée → **Embedding FastText (gelé)** → LSTM(128) → **Attention (Bahdanau)** → Dropout → Dense(sigmoid).

3. Configurer l'arrêt anticipé (*Early Stopping*)

```

1 early_stop = EarlyStopping(
2     monitor='val_loss',
3     patience=3,
4     restore_best_weights=True )

```

4. Entraîner

```
1 history_bahdanau_1 = model_bahdanau_1.fit(  
2     X_train, y_train,  
3     callbacks=[early_stop],  
4     validation_data=(X_val, y_val),  
5     epochs=20,  
6     batch_size=128)
```

Nous avons choisi **128** comme taille du batch car les données sont **massives**.

5. Évaluer

```
1 loss, acc = model_bahdanau_1.evaluate(X_test, y_test, verbose=0)  
2 print(f"Accuracy sur test : {acc:.2f}")
```

Accuracy sur test : 0.93

5.4.2 Ajout de la couche aux modèle pré-entraîner

1. Construire le modèle

On va dans la construction du modèle, utiliser les embeddings pré-entraînés.

```
1 lstm_out, state_h, state_c = LSTM(128, return_sequences=True,  
    return_state=True)(embedding)  
2 attention_layer = BahdanauAttention(units=128)  
3 att_out, att_weights = attention_layer(lstm_out, state_h)  
4 drop = Dropout(0.3)(att_out)  
5 outputs = Dense(1, activation='sigmoid')(drop)  
6 model_bahdanau_2 = Model(inputs, outputs)
```

2. Compiler

```
1 model_bahdanau_2.compile(loss='binary_crossentropy', optimizer=  
    'adam', metrics=['accuracy'])  
2 model_bahdanau_2.summary()
```

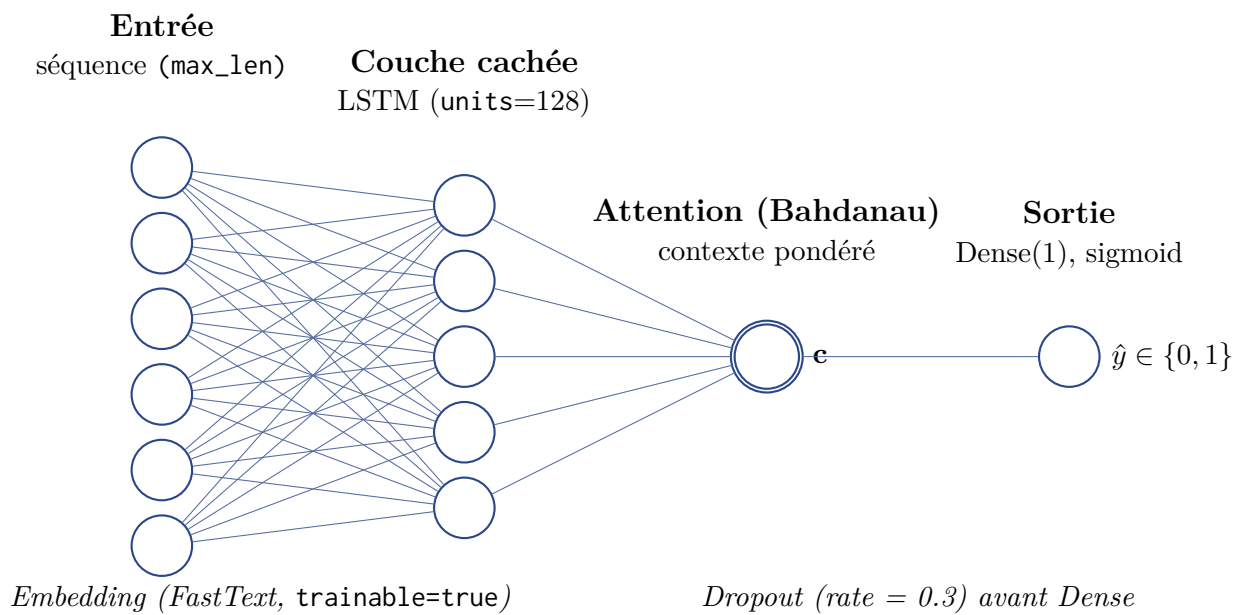


Figure 5.2 – Réseau : Entrée → **Embedding FastText (gelé)** → LSTM(128) → **Attention (Bahdanau)** → Dropout → Dense(sigmoid).

3. Entraîner

```

1 history_bahdanau_2 = model_bahdanau_2.fit(
2     X_train, y_train,
3     validation_data=(X_val, y_val),
4     callbacks=[early_stop],
5     epochs=20,
6     batch_size=128
7 )

```

4. Évaluer

```

1 loss, acc = model_bahdanau_2.evaluate(X_test, y_test, verbose=0)
2 print(f"Accuracy sur test : {acc:.2f}")

```

Accuracy sur test : 0.93

5.5 Visualiser les résultats d'entraînement

On trace l'évolution de la *loss* et de l'*accuracy* pour diagnostiquer le biais/variance et contrôler l'effet de l'attention.

```
1 plot_training_history(history_bahdanau_1, model_name="Modèle_1 et  
   attention Bahdanau")  
2 plot_training_history(history_bahdanau_2, model_name="Modèle_2 et  
   attention Bahdanau")
```

1. Ajout de la couche au modèle avec *embedding* pendant l'entraînement

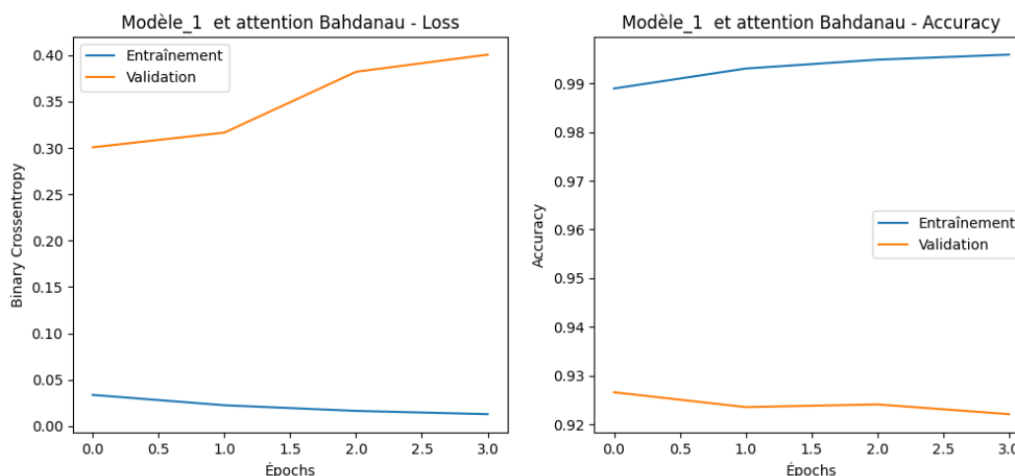


Figure 5.3 – Historique d'entraînement avec *embedding* appris.

2. Ajout de la couche au modèle pré-entraîné

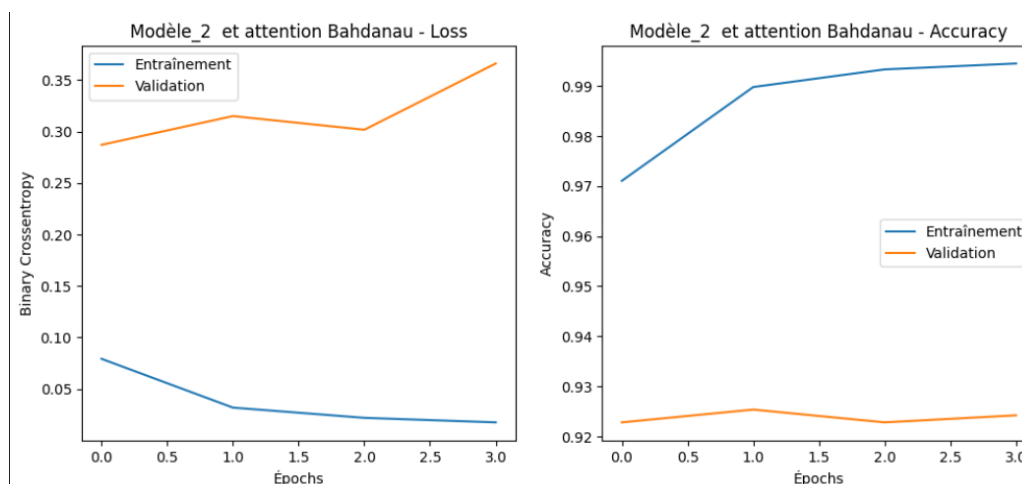


Figure 5.4 – Historique d'entraînement avec *embedding* pré-entraîné.

Interprétation

- Dans les deux cas, l'ajout de l'attention conduit à une **loss de validation faible** et une **accuracy élevée** (stabilité des courbes).
- **Embeddings appris pendant l'entraînement** : convergence **rapide** et trajectoires plus **lisses** ; l'attention capte utilement les segments pertinents.
- **Embeddings pré-entraînés** : démarrage un peu plus **lent/irrégulier**, mais la performance **rejoint un niveau proche** après quelques époques.
- Globalement, l'attention améliore l'**interprétabilité** (poids sur mots clés) et maintient une **bonne généralisation** ; le choix entre embeddings dépend surtout des données (taille/domaine) et d'un éventuel *fine-tuning*.

5.6 Conclusion

On remarque que la couche bahdanou a amélioré le modèle préentraîner et a gardé la même performance pour le modèle avec embedding pendant entraînement.

Table 5.1 – Accuracy (jeu de test) des modèles LSTM

| Modèle | Accuracy |
|--|-------------|
| Embedding pendant l'entraînement + couche Bahdanau | 0.93 |
| Embeddings pré-entraînés + couche Bahdanau | 0.93 |

L'ajout d'une attention additive au-dessus du LSTM :

- **améliore l'interprétabilité** grâce aux poids α_t (mise en évidence des mots clés) ;
- **peut stabiliser** l'entraînement en concentrant la représentation sur les segments pertinents ;
- apporte souvent un **gain modéré mais réel** en accuracy/F1 vis-à-vis d'un LSTM seul, surtout sur des avis longs ou des structures avec *négation* / *intensification*.

En pratique, l'impact dépend du **prétraitement** (tokenisation, normalisation), du **régime d'embeddings** (gelés vs *fine-tuning*), et des **hyperparamètres** (units, dropout, learning rate). Une analyse d'attention (extraction/visualisation des poids α_t) complète utilement la comparaison chiffrée.