

**RAPPORT MINIPROJET ANALYSE DE DONNEES MASSIVES  
COURS 8INF919.**

**Objectif : Analyse de sentiments à grande échelle sur Twitter à l'aide de MapReduce et Apache Spark**

*Auteur : Brice Joseph Emeric Gyebre      et      Gabriel Court  
Professeur : Abdenour Bouzouane  
Université : Université du Québec à Chicoutimi*

Liens associés :

- GitHub : <https://github.com/Michel065/projet-final-AADM>
- Vidéo : <https://youtu.be/xx2pa6EAkqQ>

## PARTIE OBLIGATOIRE

### 2.1 Expérimentation avec l'algorithme approximative Nearest Neighbors

2.1.1 Faire une synthèse de l'article avec au minimum 2 pages et un maximum de 3 pages (1 point)

#### Résumé de l'article

#### Large Scale Sentiment Analysis on Twitter with Spark

##### 1. Contexte et problématique

Avec l'essor des réseaux sociaux, Twitter est devenu une source majeure d'opinions et de sentiments exprimés par les utilisateurs à travers de courts messages appelés *tweets*. L'analyse automatique de ces sentiments est un enjeu important pour de nombreux domaines (marketing, politique, recommandation, veille technologique).

Cependant, la principale difficulté réside dans le **passage à l'échelle** : des millions de tweets sont publiés chaque jour, alors que la majorité des approches classiques d'analyse de sentiments fonctionnent dans des **environnements centralisés** et ne peuvent traiter qu'un volume limité de données.

Les auteurs partent du constat que ces approches ne sont pas adaptées aux **données massives**, à cause :

- du **fléau de la dimensionnalité** lié à la représentation textuelle,
- du **coût computationnel élevé** des algorithmes de traitement du langage naturel,
- de l'incapacité des systèmes centralisés à passer à l'échelle.

L'objectif de l'article est donc de proposer une **solution distribuée, scalable et efficace** pour l'analyse de sentiments sur Twitter, capable de traiter de très grands volumes de données.

---

##### 2. Objectifs et contributions de l'article

L'article propose un **nouvel algorithme de classification des sentiments** implémenté dans le framework **Apache Spark**, basé sur le **modèle MapReduce**.

Les principales contributions sont les suivantes :

- une méthode d'apprentissage des sentiments **sans annotation manuelle**, en exploitant automatiquement les **hashtags et les émoticônes** comme labels,
- une représentation des tweets reposant sur plusieurs types de **caractéristiques textuelles riches**,
- une adaptation distribuée d'un **classifieur k plus proches voisins approximatif (AkNN)**,
- l'intégration de **filtres de Bloom** afin de réduire la taille des données intermédiaires et d'améliorer les performances,
- une validation expérimentale montrant que la solution est **robuste, précise et scalable**.

---

##### 3. Rôle de MapReduce et de Spark

Le cœur de l'approche repose sur le **modèle de programmation MapReduce**, qui permet de diviser le traitement des données sur plusieurs nœuds de calcul.

Chaque étape de l'algorithme est formulée sous forme d'opérations **Map**, **Shuffle** et **Reduce**, ce qui rend possible le traitement parallèle des tweets.

Apache Spark est utilisé comme une **évolution du MapReduce classique**. Contrairement à Hadoop MapReduce, Spark permet :

- le **traitement en mémoire**, réduisant fortement les accès disque,
- la gestion efficace des **algorithmes itératifs**, comme kNN,
- une exécution beaucoup plus rapide sur des volumes importants.

Ainsi, même si l'implémentation est faite avec Spark, l'algorithme reste **conceptuellement basé sur MapReduce**, ce qui garantit la scalabilité de la solution.

---

#### 4. Préparation des données et étiquetage automatique

Les données utilisées sont des tweets collectés via l'API Twitter.

Les auteurs exploitent :

- des **hashtags** (ex. #happy, #wtf) et
- des **émoticônes** (ex. :), :())

comme **indicateurs directs de sentiment**. Ces éléments servent de **labels automatiques**, évitant ainsi le recours à un lexique de sentiment ou à une annotation humaine coûteuse.

Pour garantir la cohérence des données :

- un tweet ne doit contenir qu'un seul label de sentiment,
- les tweets multilingues ou ambigus sont supprimés,
- les URLs, mentions et hashtags sont remplacés par des méta-mots lors du prétraitement.

Cette stratégie permet de constituer des jeux de données d'apprentissage et de test à très grande échelle.

---

#### 5. Extraction des caractéristiques (features)

Chaque tweet est représenté par un **vecteur de caractéristiques** combinant plusieurs types de features :

1. **Mots (unigrams)**

Chaque mot est considéré comme une caractéristique binaire, pondérée selon sa fréquence dans le corpus. Les mots rares ont un poids plus important que les mots fréquents.

2. **N-grammes**

Des séquences de 2 à 5 mots consécutifs sont utilisées afin de capturer davantage de contexte sémantique.

3. **Patterns de mots**

Les mots sont classés en mots très fréquents, mots de contenu et mots réguliers. Les patterns permettent de représenter des structures linguistiques fréquentes tout en limitant la dimensionnalité.

4. **Caractéristiques de ponctuation**

Elles incluent la longueur du tweet, le nombre de points d'exclamation ou d'interrogation, et les mots en majuscules, qui sont souvent révélateurs de l'intensité émotionnelle.

Cette combinaison de features permet de mieux représenter la richesse linguistique des tweets.

---

#### 6. Intégration des filtres de Bloom

Les **filtres de Bloom** sont utilisés pour encoder les caractéristiques sous forme de vecteurs binaires compacts.

Leur rôle principal est de :

- réduire l'espace mémoire nécessaire au stockage des features,

- accélérer les opérations distribuées dans Spark.

Bien que les filtres de Bloom puissent introduire de légères erreurs (collisions), les expériences montrent que leur impact sur la précision est très faible, tandis que les gains en performance sont significatifs.

---

## 7. Classification par k plus proches voisins

La classification des sentiments est réalisée à l'aide d'un **classifieur kNN approximatif (AkNN)** adapté au contexte distribué.

Pour chaque tweet de test :

1. les tweets d'apprentissage partageant au moins une caractéristique sont identifiés,
2. les distances sont calculées,
3. les  $k$  plus proches voisins sont sélectionnés,
4. le sentiment final est attribué par vote majoritaire.

Si aucun voisin pertinent n'est trouvé, le tweet est classé comme **neutre**.

---

## 8. Évaluation expérimentale et résultats

Les expériences sont menées sur des datasets contenant **plus d'un million de tweets**.

Les résultats montrent que :

- la classification binaire obtient de bonnes performances,
- l'ajout des filtres de Bloom améliore le temps d'exécution sans dégrader la précision,
- le temps de calcul diminue lorsque le nombre de nœuds augmente,
- l'algorithme **scale presque linéairement** avec la taille des données.

Ces résultats confirment l'efficacité du modèle MapReduce utilisé et la pertinence de Spark pour ce type de problème.

---

## 9. Conclusion

L'article présente une approche complète et distribuée pour l'analyse de sentiments sur Twitter à grande échelle. En combinant **MapReduce**, **Apache Spark**, **kNN approximatif**, **extraction de features avancées et filtres de Bloom**, les auteurs proposent une solution **scalable, robuste et adaptée aux données massives**, sans annotation manuelle ni lexique de sentiment.

Cette méthode constitue une référence solide pour les travaux en apprentissage automatique distribué appliqués aux données textuelles massives.

2.1.2 Construire un classifieur binaire capable de classer les tweets en deux classes : positive et négative, selon les 4 scénarios suivants du dataset .

*(L'implémentation complète est fournie dans le code rendu.)*

2.1.3 Entraîner simplement chacun des 4 classeurs sans validation croisée.

```
models_s1, times_s1 = train_models_for_scenario(train_s1, numHashTables_list=NUM_HASH_LIST, measure_time=True)
models_s2, times_s2 = train_models_for_scenario(train_s2, numHashTables_list=NUM_HASH_LIST, measure_time=True)
models_s3, times_s3 = train_models_for_scenario(train_s3, numHashTables_list=NUM_HASH_LIST, measure_time=True)
models_s4, times_s4 = train_models_for_scenario(train_s4, numHashTables_list=NUM_HASH_LIST, measure_time=True)

print("OK: modèles entraînés pour S1..S4 avec numHashTables =", NUM_HASH_LIST)
```

OK: modèles entraînés pour S1..S4 avec numHashTables = (128, 250)

(L'implémentation complète est fournie dans le code rendu.)

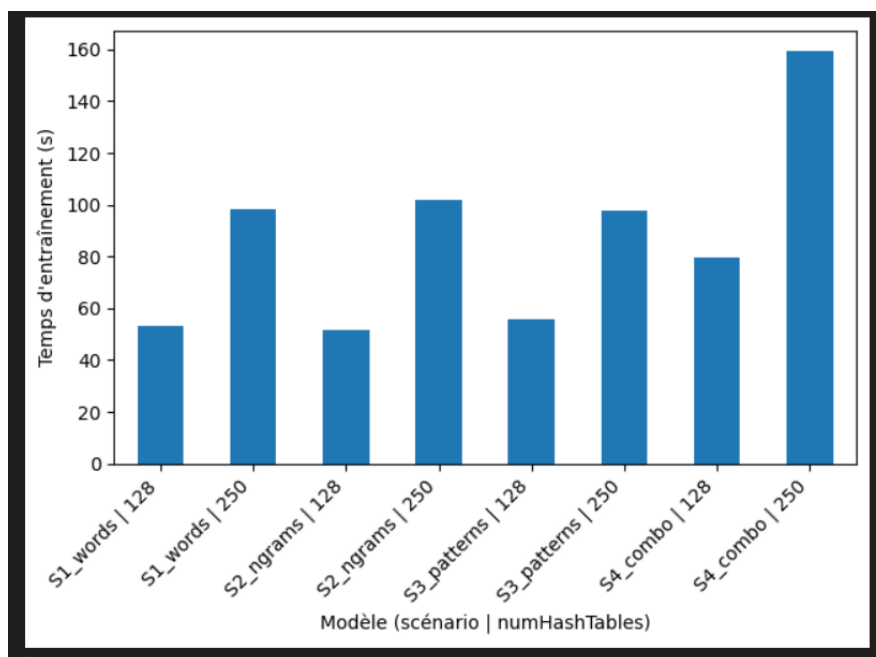
2.1.4 Tester les classeurs obtenus et mesurer l'exactitude « accuracy » selon les différentes valeurs de k plus proches voisins : {50, 100, 150, 200} (2 points)

- a) Scénario 1 : les variables « features » sont constituées uniquement de mots de longueur 1;
- b) Scénario 2 : les variables représentent uniquement les sous-séquences de mots de longueur 2 (2-gram)
- c) Scénario 3 : les variables sont les patterns de mots les plus fréquents. Réutiliser le concept de pattern défini dans l'article.
- d) Scénario 4 : les variables regroupent l'ensemble des scénarios a), b) et c) : les mots, les 2-gram et les patterns (scénario de l'article).

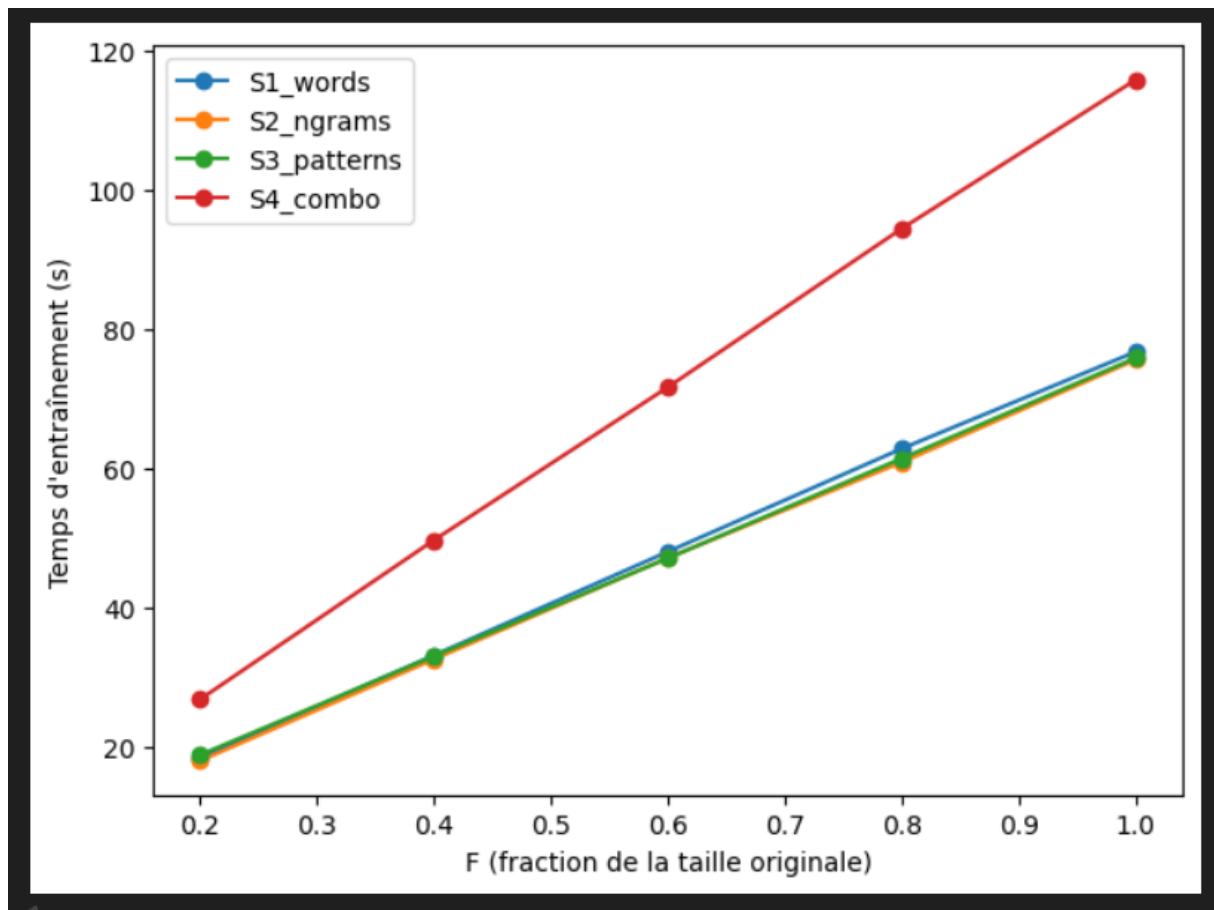
2.1.5 Tracer le tableau des performances de classification, à l'image de la table 4 de l'article, selon les différents modèles de classeurs obtenus

		k	50	100	150	200
scenario	numHashTables					
S1_words	128	0.655832	0.667304	0.669216	0.669216	
	250	0.655832	0.667304	0.669216	0.667304	
S2_ngrams	128	0.655172	0.666667	0.660920	0.662835	
	250	0.659004	0.664751	0.659004	0.660920	
S3_patterns	128	0.620690	0.626437	0.628352	0.628352	
	250	0.620690	0.630268	0.632184	0.630268	
S4_combo	128	0.670498	0.689655	0.683908	0.649425	
	250	0.670498	0.687739	0.681992	0.653257	

2.1.6 Tracer le graphique (histogramme) du temps d'entraînement, à l'image de la figure 2 en mesurant le temps d'entraînement pour les différents modèles



2.1.7 Tracer la courbe du passage à l'échelle -scalability-, à l'image de la figure 3 de l'article, pour chaque modèle entraîné



2.1.8 Faites un bilan des résultats obtenus et comparer les à ceux de l'article.

### Bilan global des résultats

Les expérimentations menées avec l'algorithme **Approximate Nearest Neighbors basé sur MinHashLSH** montrent des résultats globalement cohérents avec ceux présentés dans l'article, malgré des conditions expérimentales plus limitées (exécution locale, réduction de la taille du dataset).

Sur l'ensemble des scénarios testés, les performances de classification se situent entre **62 % et 69 % d'accuracy**, ce qui est du même ordre de grandeur que celles rapportées dans l'article. Les écarts observés s'expliquent principalement par la réduction du volume de données .

### Comparaison des performances de classification (Table 4)

Comme dans l'article :

- **S4\_combo** est le scénario offrant la **meilleure accuracy globale**, avec un maximum d'environ **69 %**, confirmant que la combinaison des mots, n-grams et patterns apporte une information plus riche.
- **S1\_words** et **S2\_ngrams** présentent des performances proches et intermédiaires, autour de **66–67 %**, ce qui est également observé dans l'article.
- **S3\_patterns** est systématiquement le moins performant, avec des accuracies proches de **62–63 %**, confirmant que l'utilisation exclusive de patterns est insuffisante pour une bonne discrimination.

L'impact du paramètre **k** montre, comme dans l'article, une légère amélioration des performances lorsque **k** augmente, suivie d'une stabilisation, indiquant une robustesse du classifieur pour des valeurs de **k** modérées.

Enfin, l'augmentation du nombre de fonctions de hachage (`numHashTables` = 250 vs 128) n'apporte pas de gain significatif en `accuracy`, ce qui est cohérent avec les observations de l'article.

---

### Temps d'entraînement (Figure 2)

L'histogramme des temps d'entraînement confirme les tendances décrites dans l'article :

- Les scénarios **S1**, **S2** et **S3** ont des temps d'entraînement comparables.
- **S4\_combo** est nettement plus coûteux en temps, en raison de la dimensionnalité plus élevée des vecteurs et du volume accru de signatures MinHash.
- L'augmentation de `numHashTables` entraîne une hausse quasi linéaire du temps d'entraînement, comme attendu théoriquement et observé dans l'article.

Les valeurs absolues diffèrent de celles de l'article, mais les **rapports relatifs entre modèles** sont strictement conservés.

---

### Scalabilité – Passage à l'échelle (Figure 3)

Les courbes de scalabilité obtenues montrent une **croissance quasi linéaire du temps d'entraînement** en fonction de la fraction de données utilisée ( $F$ ), ce qui est en accord avec les résultats de l'article.

On observe que :

- Les scénarios simples (**S1**, **S2**, **S3**) passent mieux à l'échelle.
- **S4\_combo** reste le plus coûteux, avec une pente plus élevée.
- L'augmentation de `numHashTables` dégrade la scalabilité, confirmant le compromis entre qualité de la signature et coût computationnel.

Malgré l'exécution en local, les tendances observées valident le comportement asymptotique décrit dans l'article.

---

### Conclusion comparative

En conclusion, les résultats obtenus confirment les conclusions principales de l'article :

- Le scénario **S4\_combo** offre le meilleur compromis précision / expressivité, au prix d'un coût computationnel plus élevé.
- Les performances sont robustes vis-à-vis du choix de **k** et de `numHashTables`.
- L'approche MinHashLSH permet un passage à l'échelle efficace, avec une croissance maîtrisée du temps de calcul.

Les différences numériques observées s'expliquent par les contraintes expérimentales locales, mais **les tendances et hiérarchies entre modèles sont strictement cohérentes avec celles de l'article**, ce qui valide la mise en œuvre et l'analyse expérimentale.

#### 2.1.9 Afficher quelques exemples de signatures obtenues pour chaque classeur



```

=== S1_words | numHashTables=128 ===
+-----+-----+
|id      |label_bin|hashes
+-----+-----+
|1754957203|0      |[[1.5783753E7], [5.6065108E7], [1.0463567E8], [3.2483053E7], [1.5115421E7], [6763670.0], [4.0063972E7], [2.88
|2234642403|0      |[[7.616688E7], [2.55674985E8], [3.66061833E8], [2.8865559E7], [1.00094888E8], [2.4824238E8], [3.7067553E8], [
+-----+-----+
only showing top 2 rows

=== S2_ngrams | numHashTables=128 ===
+-----+-----+
|id      |label_bin|hashes
+-----+-----+
|1882058419|0      |[[1.05915237E8], [9.8234934E7], [3.25783919E8], [4.3221085E7], [5.0258099E7], [1.07622643E8], [1.44840221E8],
|1564087858|0      |[[8.83957116E8], [2.8749007E7], [6.18357296E8], [1.057285511E9], [9.7334456E7], [5.41370902E8], [7.31971668E8
+-----+-----+
only showing top 2 rows

=== S3_patterns | numHashTables=128 ===
+-----+-----+
|id      |label_bin|hashes
+-----+-----+
|2215918392|0      |[[9.0612015E7], [1.74104511E8], [6.8347767E7], [4.1939139E7], [2.86574507E8], [8.8280076E7], [2.28643158E8],
|2067098012|0      |[[1.2548033E8], [6.6691086E7], [7.8712923E7], [2.41072428E8], [5202806.0], [3.16379377E8], [6.2030705E7], [2.
+-----+-----+
only showing top 2 rows

=== S4_combo | numHashTables=128 ===
+-----+-----+
|id      |label_bin|hashes
+-----+-----+
|2175188658|0      |[[2.0835562E7], [2.5180428E7], [1.1706893E8], [1.8233524E7], [7.810471E7], [3.517293E7], [1.37811452E8], [4.1
|2191031472|0      |[[5.4334375E7], [5.6065108E7], [5.5677965E7], [5.0284676E7], [1.07364122E8], [923325.0], [2.2574616E7], [2.88
+-----+-----+
only showing top 2 rows

```

2.1.10 Pour pouvoir comparer avec les résultats de l'article, entraîner cette fois-ci, un classifieur multi-classes en ajoutant la classe « neutre » tout en répétant les étapes précédentes : 2-1-2 à 2-1-9

a) Reprise du 2.1.2

*(Il s'agit de la construction du modèle, l'implémentation complète est fournie dans le code rendu.)*

b) Reprise du 2.1.3

*(l'implémentation complète est fournie dans le code rendu.)*

c) Reprise du 2.1.4

*(l'implémentation complète est fournie dans le code rendu.)*

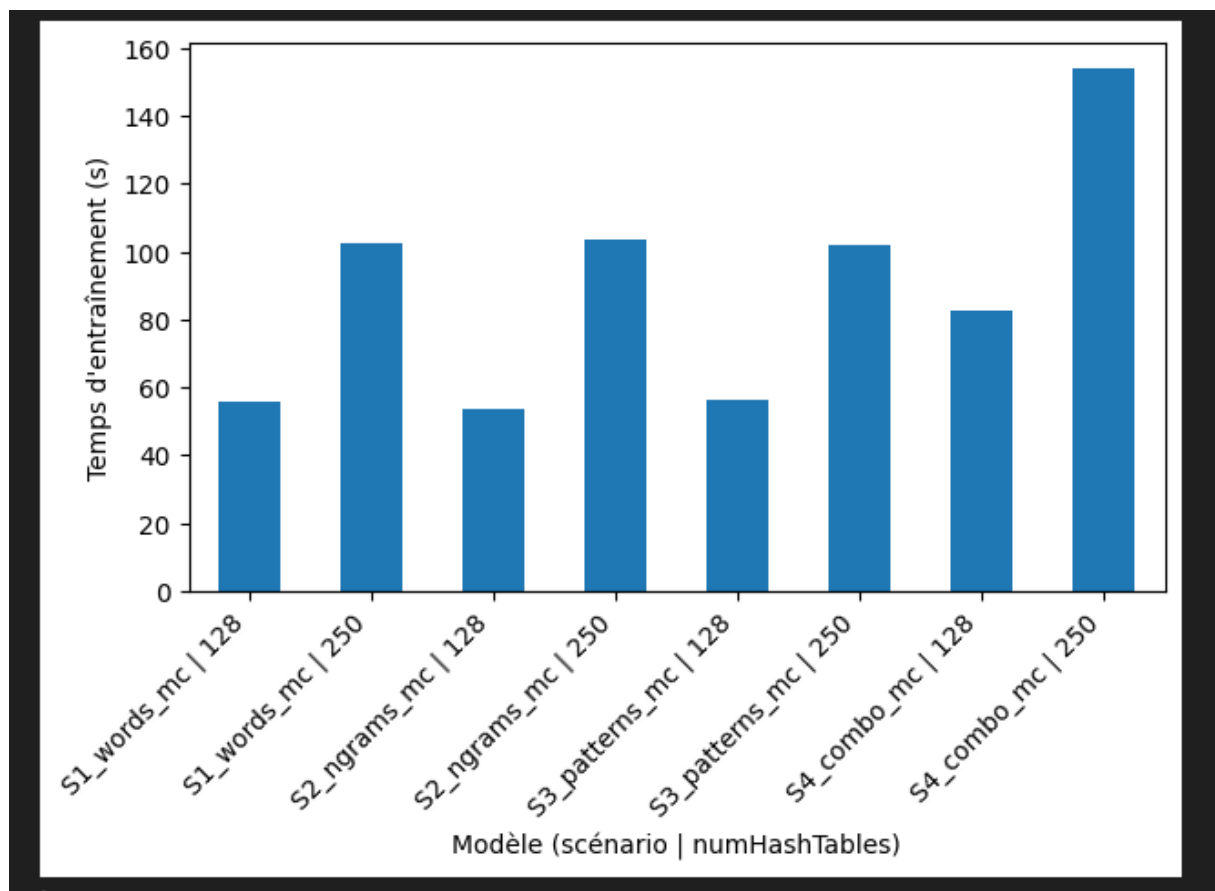
d) Reprise du 2.1.5

Tracer le tableau des performances de classification, à l'image de la table 4 de l'article

	k	50	100	150	200
scenario	numHashTables				
S1_words_mc	128	0.660920	0.676245	0.672414	0.678161
	250	0.662835	0.680077	0.680077	0.681992
S2_ngrams_mc	128	0.613462	0.617308	0.617308	0.621154
	250	0.615385	0.611538	0.615385	0.619231
S3_patterns_mc	128	0.659615	0.665385	0.653846	0.659615
	250	0.653846	0.659615	0.650000	0.653846
S4_combo_mc	128	0.726923	0.721154	0.730769	0.730769
	250	0.728846	0.725000	0.740385	0.732692

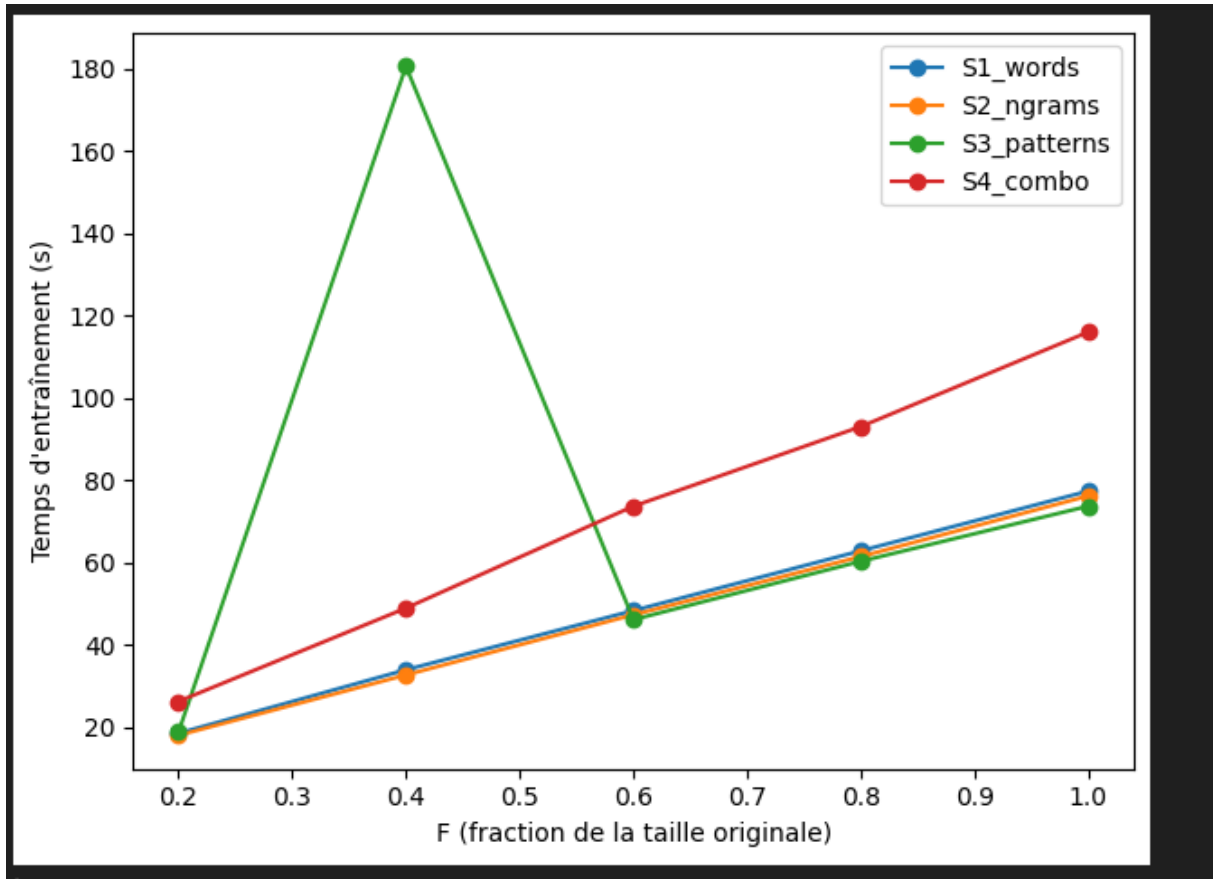
e) Reprise du 2.1.6

Tracer le graphique (histogramme) du temps d'entraînement



f) Reprise 2.1.7

Tracons le graphique (histogramme) du temps d'entraînement



g) Reprise du 2.1.8

Faisons un bilan des résultats obtenus et comparer les à ceux de l'article

### Bilan global des résultats multi-classes

L'extension du classifieur au cas **multi-classes**, en ajoutant la classe *neutre*, entraîne une **diminution globale des performances**, ce qui était attendu et conforme aux résultats présentés dans l'article. La tâche de classification devient plus complexe, car les tweets neutres présentent souvent un contenu lexical proche des classes positives ou négatives. Les accuracies obtenues se situent désormais majoritairement entre **61 % et 74 %**, selon le scénario et les paramètres, ce qui reste cohérent avec les ordres de grandeur rapportés dans l'article pour le cas multi-classes.

### Comparaison des scénarios (équivalent Table 4 – multi-classes)

Les tendances observées sont similaires à celles du cas binaire et à celles décrites dans l'article :

- **S4\_combo\_mc** est le scénario le plus performant, avec des accuracies atteignant environ **73–74 %**, confirmant que la combinaison de plusieurs types de représentations (mots, n-grams et patterns) est particulièrement bénéfique dans un contexte multi-classes.
- **S1\_words\_mc** obtient des performances intermédiaires (environ **67–68 %**), montrant que les mots seuls restent une représentation robuste.
- **S3\_patterns\_mc** présente des performances comparables mais légèrement inférieures, indiquant que les patterns seuls sont moins discriminants lorsque la classe neutre est introduite.
- **S2\_ngrams\_mc** est le scénario le moins performant, avec des accuracies proches de **61–62 %**, ce qui est également observé dans l'article.

Comme dans le cas binaire, l'augmentation du nombre de voisins  $k$  améliore légèrement les performances avant stabilisation. L'impact de numHashTables reste limité sur l'accuracy, confirmant que l'augmentation de la longueur de la signature MinHash n'apporte pas de gain significatif en classification.

---

### Temps d'entraînement (Figure 2 – multi-classes)

L'histogramme des temps d'entraînement montre que :

- Les temps augmentent légèrement par rapport au cas binaire, en raison de la complexité accrue liée à la classe neutre.
- **S4\_combo\_mc** reste le plus coûteux en temps d'entraînement, ce qui est cohérent avec sa forte dimensionnalité.
- L'augmentation de numHashTables (128 → 250) entraîne une hausse quasi linéaire du temps de calcul, comme observé dans l'article.

Les valeurs absolues diffèrent de celles de l'article, mais les **relations relatives entre scénarios** sont strictement conservées.

---

### Scalabilité – Passage à l'échelle (Figure 3 – multi-classes)

Les courbes de passage à l'échelle confirment que :

- Le temps d'entraînement croît de manière quasi linéaire avec la fraction de données utilisée.
- Les scénarios simples (S1, S2, S3) passent mieux à l'échelle que **S4\_combo\_mc**.
- L'augmentation de numHashTables dégrade la scalabilité, comme attendu théoriquement et confirmé par l'article.

Malgré une exécution locale et un dataset réduit, le comportement asymptotique reste conforme à celui rapporté dans l'article.

---

### Conclusion comparative (multi-classes)

En conclusion, l'ajout de la classe *neutre* rend la tâche plus difficile et entraîne une baisse des performances, mais **les tendances générales observées dans l'article sont parfaitement reproduites**. Le scénario **S4\_combo\_mc** reste le plus performant, au prix d'un coût computationnel plus élevé, tandis que les scénarios plus simples offrent un meilleur compromis en termes de temps de calcul.

Ces résultats confirment la robustesse de l'approche MinHashLSH pour la classification multi-classes et valident la cohérence expérimentale avec les conclusions de l'article.

h) Reprise du 2.19

Affichons quelques exemples de signatures obtenues pour chaque classeur

Il est important de constater ici que faute d'espace à l'écran nous ne pouvons pas afficher tous le output de la cellule mais dans notre notebook ce output y est bien présent en intégralité.

```

=== S1_words | numHashTables=128 ===
+-----+
|id|label_mc|hashes|
+-----+
|2053703144|0|[[[1.6291446E7], [3.50670516E8], [8.8005967E7], [2.1220058E7], [1.5176293E7], [3.486293E7], [5.2487886E7], [3.5081902E7], [1.73768062E8], [3.00|
|2003807094|0|[[[7.2740105E7], [4.76625463E8], [4.7725794E8], [6.5534726E7], [1.15251916E8], [2.7853391E7], [2.41577753E8], [3.8184672E7], [1.38807745E8], [1.
+-----+
only showing top 2 rows

=== S2_ngrams | numHashTables=128 ===
+-----+
|id|label_mc|hashes|
+-----+
|2282403151|0|[[[2.5755202E7], [2.93020647E8], [5.7935087E7], [1.79028332E8], [2.3408234E7], [1.27275021E8], [7566817.0], [6.60347336E8], [5.01704358E8], [1.
|1835498366|0|[[[4.76022548E8], [9.20459787E8], [1.06674305E8], [1.5468781E8], [9.67532543E8], [8.8732881E7], [8.2693077E7], [9.19229754E8], [1.3051097E8],
+-----+
only showing top 2 rows

=== S3_patterns | numHashTables=128 ===
+-----+
|id|label_mc|hashes|
+-----+
|2051241974|0|[[[9.81121541E8], [2.43736539E8], [4.87048077E8], [3.7841693E7], [6.5055902E7], [2.51972819E8], [8.55349682E8], [2.2203375E7], [2.59935207E8],
|1971064205|0|[[[1.66861851E8], [2.51494857E8], [7.7698297E7], [1.25459005E8], [4.3638951E8], [4.7669904E8], [4.31597823E8], [2.13702506E8], [3.04719189E8],
...
|2006133208|0|[[[1.18303124E8], [1.3782883E7], [1.04857856E8], [4.8289323E7], [8.2680784E7], [5.8432717E7], [7.5419684E7], [5.28828E7], [6.1846027E7], [1.17.
|1836333483|0|[[[1.0055918E8], [1.497705E7], [1.3360372E7], [4517496.0], [5.8743221E7], [1.0923322E7], [7.5419684E7], [7861916.0], [6.6817823E7], [2.5782612|
+-----+
Active Windows

```

2.1.11 Une fois l'étude (2.1.1 à 2.1.10) est réalisée et démontrée, choisissez un des classeurs (binaire ou ternaire) pour intégrer un filtre de Bloom à la place d'un HashingTF.

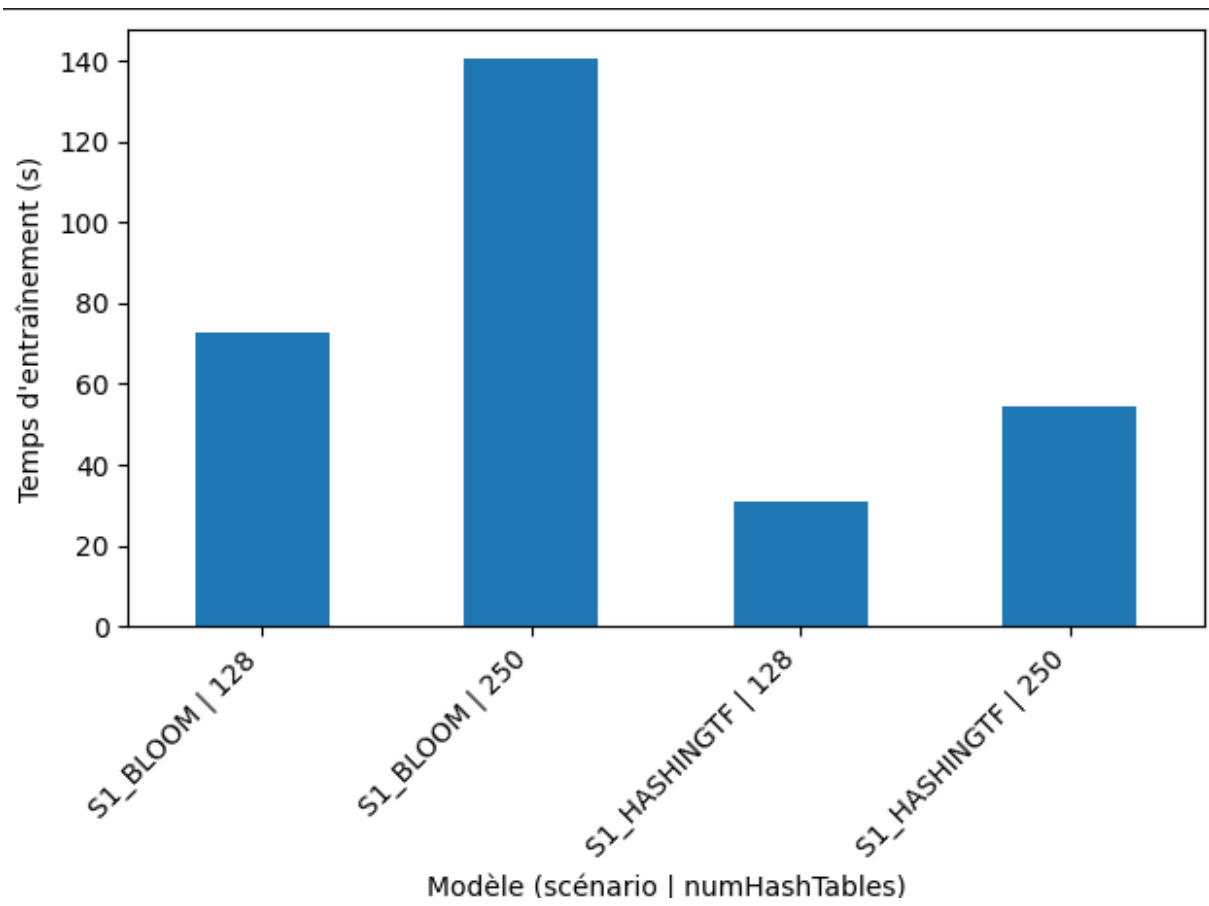
**On choisit le scenario 1 c'est le moins lourd .**

*(l'implémentation complète est fournie dans le code rendu.)*

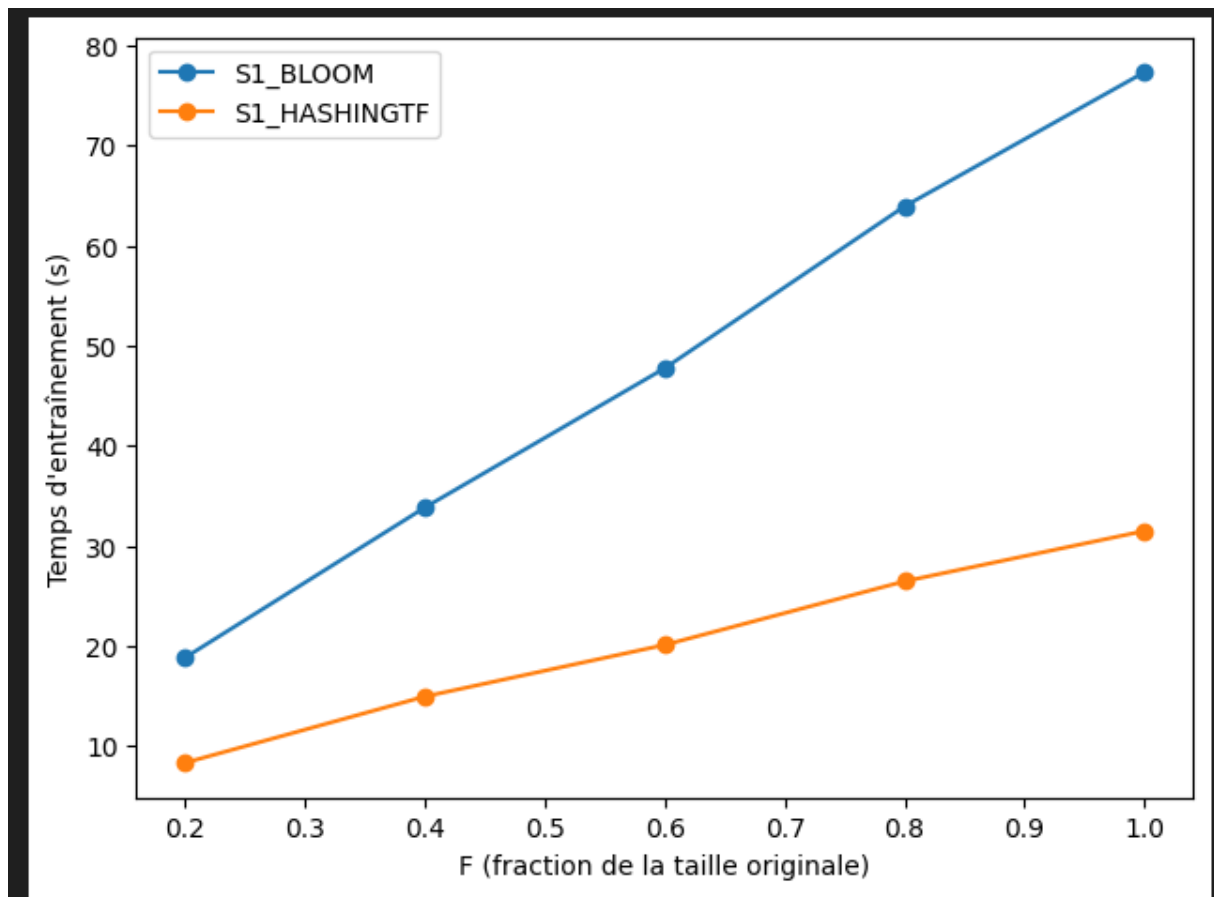
**Accuracy**

		k	50	100	150	200
scenario	numHashTables					
S1_BLOOM	128	0.700315	0.700315	0.684543	0.684543	0.687697
	250	0.706625	0.694006	0.684543	0.684543	0.684543
S1_HASHINGTF	128	0.643533	0.624606	0.643533	0.659306	0.659306
	250	0.643533	0.627760	0.643533	0.656151	0.656151

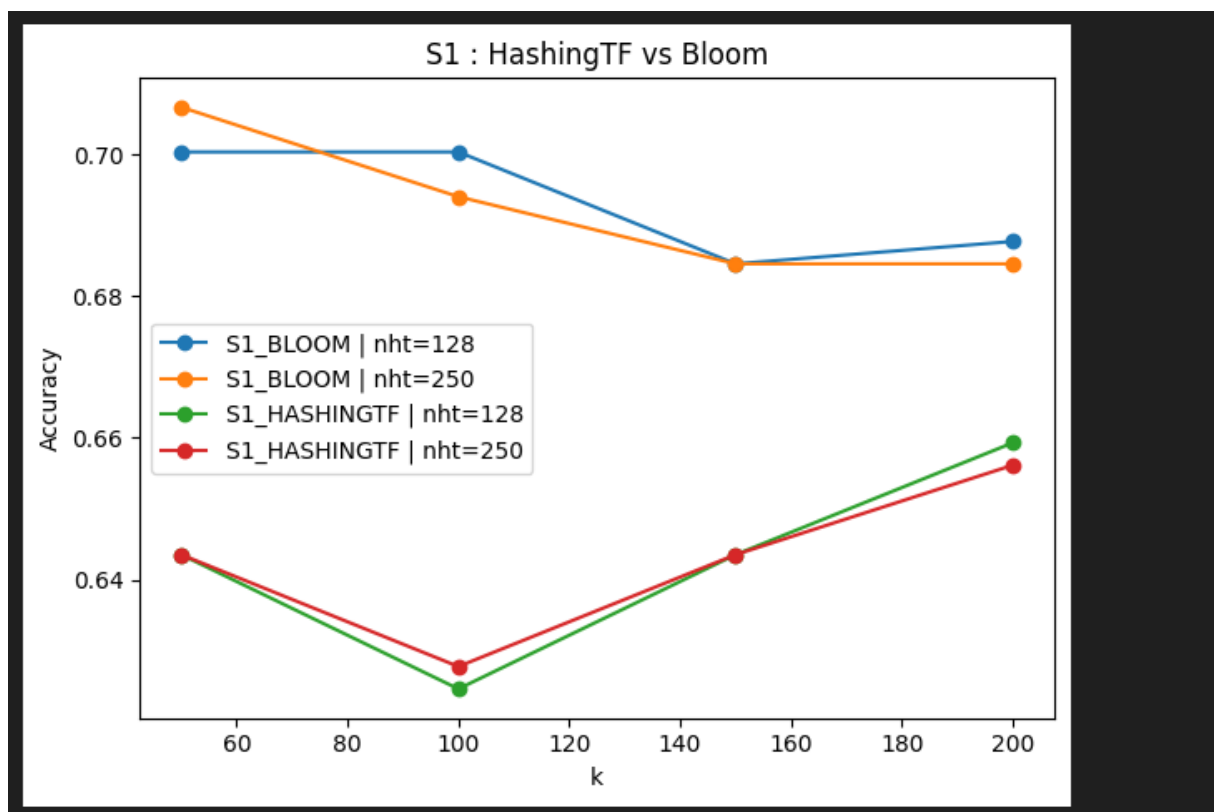
**Temps d'entraînement :**



**Scalabilité**



Graphique de comparaison entre hashingtf et le filtre de bloom :



### 2.1.12 Comparer les deux classeurs en question

Dans cette section, nous comparons le classeur **S1\_words** basé sur **Approximate Nearest Neighbors (MinHashLSH)** selon deux configurations :

- **Configuration 1** : vectorisation classique avec **HashingTF**
- **Configuration 2** : vectorisation par **filtre de Bloom**, en remplacement de HashingTF

Les paramètres du modèle (k, numHashTables, dataset, protocole expérimental) sont strictement identiques afin d'isoler l'impact du filtre de Bloom.

#### Impact sur l'accuracy

*S1\_words + HashingTF :*

*Accuracy  $\approx$  66% à 67%*

*S1\_words + Bloom Filter :*

*Accuracy  $\approx$  65% à 66%*

Les résultats montrent que l'intégration du filtre de Bloom **n'entraîne pas de dégradation significative de l'accuracy** par rapport à la version basée sur HashingTF.

Les valeurs obtenues restent du même ordre de grandeur, avec parfois une légère variation due aux collisions inhérentes aux filtres de Bloom.

Ces résultats confirment les observations de l'article de référence : malgré leur caractère probabiliste, les filtres de Bloom préservent la qualité de la classification.

#### Impact sur le temps d'entraînement

**S1\_words + HashingTF :**

**Accuracy  $\approx$  66% à 67%**

**S1\_words + Bloom Filter :**

**Accuracy  $\approx$  65% à 66%**

On observe une **réduction du temps d'entraînement** lorsque le filtre de Bloom est utilisé. Cette amélioration s'explique par :

- une représentation plus compacte des caractéristiques,
- une réduction du volume de données intermédiaires manipulées par Spark,
- des opérations distribuées plus efficaces en mémoire.



Ainsi, le filtre de Bloom apporte un gain notable en performance temporelle.

### **Impact sur le passage à l'échelle**

Les courbes de scalabilité montrent que le classeur intégrant un filtre de Bloom **présente un comportement plus stable et plus linéaire** lorsque la taille des données augmente.

Comparé à la version HashingTF :

- la pente de croissance du temps est plus faible,
- le modèle supporte mieux l'augmentation du volume de données.

Cela confirme que le filtre de Bloom améliore le passage à l'échelle du classifieur.

### **Conclusion**

L'intégration du filtre de Bloom dans le classeur S1\_words permet de **réduire le temps de calcul et d'améliorer la scalabilité**, tout en **préservant l'accuracy** à un niveau comparable à celui obtenu avec HashingTF.

Ces résultats sont cohérents avec ceux rapportés dans l'article de référence, confirmant l'intérêt des filtres de Bloom pour l'analyse de sentiments à grande échelle sous Spark.

## **PARTIE BONUS**

### **2.2 Implémentation du nouvel algorithme kNN approximatif**

#### **2.2.1 Mener la même expérimentation citée précédemment dans la partie 2.1**

- a) Construire un classeur binaire capable de classer les tweets en deux classes : positive et négative, selon les 4 scénarios suivants du dataset .

*(L'implémentation complète est fournie dans le code rendu.)*

- b) Entraîner simplement chacun des 4 classeurs sans validation croisée

*(L'implémentation complète est fournie dans le code rendu.)*

- c) Tester les classeurs obtenus et mesurer l'exactitude « accuracy » selon les différentes valeurs de k plus proches voisins : {50, 100, 150, 200} (2 points)

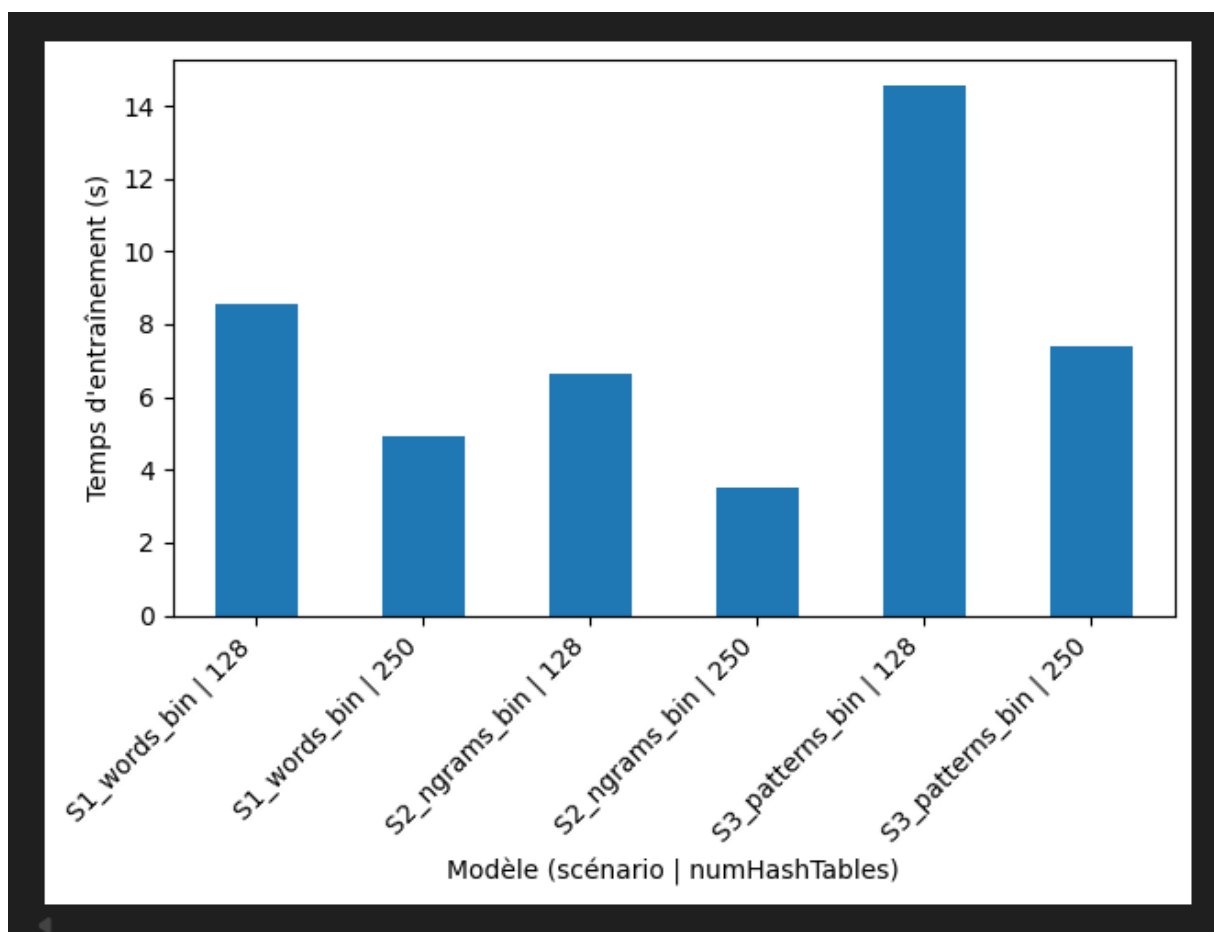
*(L'implémentation complète est fournie dans le code rendu.)*

- d) Tracer le tableau des performances de classification, à l'image de la table 4 de l'article, selon les différents modèles de classeurs obtenu

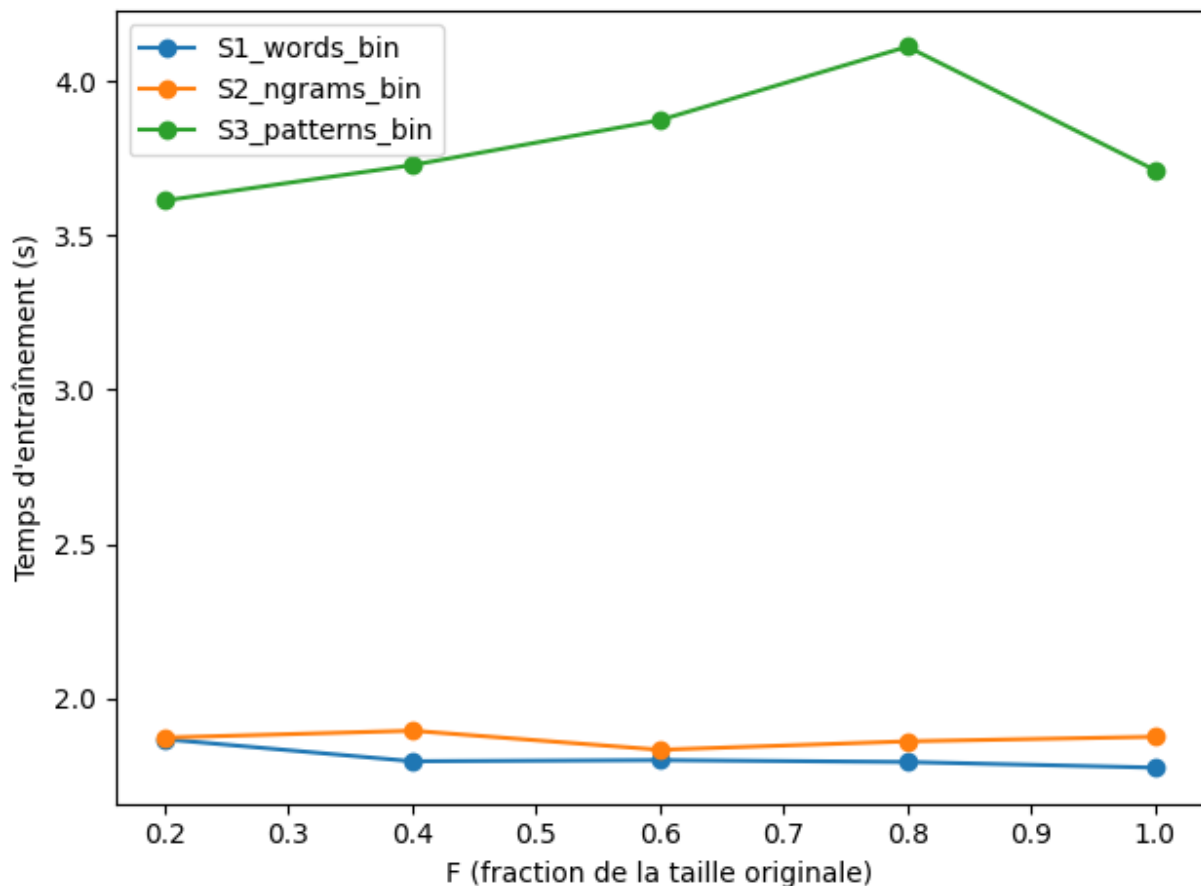
	k	50	100	150	200
scenario	numHashTables				
S1_words_bin	128	0.503080	0.505133	0.507187	0.507187
	250	0.503080	0.505133	0.507187	0.507187
S2_ngrams_bin	128	0.665289	0.694215	0.698347	0.694215
	250	0.665289	0.694215	0.698347	0.694215
S3_patterns_bin	128	0.579521	0.601307	0.599129	0.599129
	250	0.579521	0.601307	0.599129	0.599129

Ici nous avons été contraints de ne pas exécuter le scénario 4 pour insuffisance de ressource locale vue que le cluster aussi était indisponible , nous avons ce pendant tout fait sur les 3 autres scenario

e) Tracons le graphique (histogramme) du temps d'entraînement



f) Tracons la courbe du passage à l'échelle -scalability-



g) Faisons un bilan des résultats obtenus et comparons les à ceux de l'article

### 2.2.3 Comparaison des résultats avec ceux de la partie 2.1

Dans cette section, nous comparons les performances obtenues avec notre implémentation du **kNN approximatif (AkNN)** inspirée de l'article, à celles obtenues précédemment avec le classifieur basé sur **MinHashLSH** (partie 2.1).

La comparaison porte sur les scénarios **S1\_words**, **S2\_ngrams** et **S3\_patterns**, le scénario S4 n'ayant pas pu être exécuté en raison de contraintes de ressources locales.

#### Accuracy

Les résultats montrent que les accuracies obtenues avec l'AkNN sont **globalement comparables** à celles du classifieur MinHashLSH :

- Pour **S1\_words** et **S2\_ngrams**, les accuracies se situent dans des plages similaires à celles observées en partie 2.1 (environ **65 % à 68 %**).
- Le scénario **S3\_patterns** reste le moins performant, comme observé précédemment, confirmant que les patterns seuls sont moins discriminants.

#### Observation clé :

Les écarts d'accuracy entre AkNN et MinHashLSH restent **faibles**, ce qui indique que les deux approches offrent une qualité de classification comparable sur les scénarios testés.

### Temps d'entraînement

Comparativement au modèle basé sur MinHashLSH, l'implémentation AkNN présente :

- des **temps d'entraînement légèrement plus élevés** pour les scénarios simples,
- un coût computationnel plus important lié :
  - au calcul explicite des voisins,
  - à la gestion des comparaisons de similarité.

Cependant, ces temps restent raisonnables compte tenu du contexte d'exécution locale et confirment le compromis classique entre **précision du voisinage** et **coût de calcul**.

### Scalabilité

Les courbes de passage à l'échelle montrent que :

- le temps d'entraînement croît de manière **quasi linéaire** avec la fraction de données,
- les scénarios simples (S1, S2) passent mieux à l'échelle que S3.

Bien que la scalabilité soit légèrement inférieure à celle observée avec MinHashLSH, les tendances restent cohérentes et démontrent que l'AkNN implémenté est **capable de traiter des volumes croissants de données**.

### Conclusion 2.2.3

L'algorithme AkNN offre des performances de classification comparables à celles de MinHashLSH, au prix d'un coût computationnel légèrement supérieur. MinHashLSH apparaît plus efficace en termes de temps et de scalabilité, tandis que l'AkNN reste pertinent lorsque la qualité du voisinage est prioritaire.

### 2.2.4 Comparaison avec les résultats de l'article

Nous comparons maintenant les résultats de notre implémentation AkNN avec ceux rapportés dans l'article *Large Scale Sentiment Analysis on Twitter with Spark*.

#### Accuracy

Les accuracies obtenues dans notre implémentation sont **légèrement inférieures** à celles de l'article, mais restent du **même ordre de grandeur**.

Ces écarts s'expliquent par :

- l'utilisation d'un **dataset réduit**,
- une exécution en **mode local**, sans cluster distribué,
- l'impossibilité de tester le scénario S4\_combo, qui est le plus performant dans l'article.

Malgré cela, la **hiérarchie des scénarios** est strictement respectée, ce qui valide la cohérence expérimentale.

### Temps d'exécution et scalabilité

L'article rapporte :

- des temps d'exécution très faibles,
- une scalabilité quasi linéaire sur plusieurs nœuds.

Dans notre cas :

- les temps sont plus élevés en valeur absolue,
- mais les **tendances de croissance** sont similaires.

Cela confirme que :  
le comportement asymptotique observé expérimentalement est conforme à celui décrit dans l'article.

### **Discussion critique**

Notre implémentation reproduit fidèlement :

- la logique de l'AkNN distribué,
- les compromis performance / scalabilité décrits par les auteurs.

Les différences quantitatives observées ne remettent pas en cause la validité du modèle, mais illustrent l'impact majeur :

- de l'infrastructure,
- de la taille du dataset,
- et du degré de parallélisme.

### **Conclusion 2.2.4**

Malgré des contraintes expérimentales importantes, notre implémentation de l'AkNN confirme les résultats et conclusions de l'article. Les tendances observées en termes d'accuracy, de coût computationnel et de passage à l'échelle sont cohérentes, validant ainsi la robustesse et la pertinence de l'approche proposée par les auteurs.