



République Algérienne Démocratique et Populaire
Ministère de l'Enseignement Supérieur et de la Recherche
Scientifique

Université des Sciences et de la Technologie Houari Boumediene

Faculté d'Informatique
Département de Système Informatique
Master 2 : Big Data Analytics

Algorithme Apriori avec Automatisation du Support Minimal

Étude comparative : Apriori Optimisé vs Apriori Classique
Application au marché de l'emploi en Intelligence Artificielle

Réalisé par le binome 8 :

- MOKRANI Sarra
- BENDAHGANE Sami

Année académique : 2025/2026

31 octobre 2025

Table des matières

1	Introduction	1
1.1	Contexte du projet	1
1.2	Objectif principal	1
1.3	Jeu de données	1
2	Prétraitement des données	1
2.1	Chargement et sélection des attributs	1
2.2	Transformation en format transactionnel	1
2.3	Construction de la matrice de contexte	2
3	Automatisation du support minimal	2
3.1	Problématique	2
3.2	Méthode proposée : Support adaptatif hybride	2
3.2.1	Formule mathématique	3
3.2.2	Implémentation	3
3.3	Avantages de cette méthode	4
4	Algorithme Apriori implémenté	4
4.1	Principe de base	4
4.2	Fonction de calcul du support	4
4.3	Algorithme principal	4
4.4	Étapes de l'algorithme	5
5	Génération des règles d'association	5
5.1	Principe	5
5.2	Métriques d'intérêt	5
5.3	Implémentation	6
5.4	Interprétation	6
6	Exports et résultats	6
6.1	Fichiers générés	6
6.2	Code d'export	7
7	Conclusion et perspectives	7
8	Comparaison avec l'algorithme Apriori classique (Weka)	8
8.1	Objectif de comparaison	8
8.2	Critères de comparaison	8
8.3	Enjeux méthodologiques	8
9	Préparation des données pour Weka	8
9.1	Format ARFF	8
9.2	Problématique des co-absences	8
9.3	Solution adoptée : Format "true-only"	8
9.4	Script de conversion Python vers ARFF	9
9.5	Exemple de fichier ARFF généré	9

10 Exécution de l'algorithme Apriori dans Weka	10
10.1 Chargement du fichier	10
10.2 Configuration de l'algorithme Apriori	10
10.3 Commande ligne équivalente	10
10.4 Format des résultats Weka	10
11 Extraction et normalisation des résultats Weka	10
11.1 Parsing des règles	10
11.2 Normalisation des données	11
12 Fusion et analyse comparative	12
12.1 Création de clés canoniques	12
12.2 Fusion des datasets	12
12.3 Statistiques de recouvrement	13
13 Visualisations comparatives	13
13.1 Graphiques générés	13
13.2 Code de visualisation	15
13.3 Tableau statistique récapitulatif	15
14 Résultats et interprétation	16
14.1 Résultats attendus	16
14.2 Différences observables	16
14.3 Avantages de l'approche optimisée	16
15 Conclusion générale	16
15.1 Synthèse du projet	16
15.2 Apports méthodologiques	17
15.3 Perspectives d'amélioration	17
15.4 Applications pratiques	17

1 Introduction

1.1 Contexte du projet

Ce rapport présente l'implémentation d'un algorithme d'extraction de règles d'association (Apriori) appliqué au marché de l'emploi en intelligence artificielle. Le projet vise à découvrir des motifs de co-occurrence entre compétences techniques et outils professionnels demandés par les employeurs dans le secteur de l'IA.

1.2 Objectif principal

L'objectif est de développer une variante optimisée de l'algorithme Apriori qui **automatise le choix du support minimal** en fonction des caractéristiques intrinsèques du dataset, éliminant ainsi le besoin de fixer manuellement ce paramètre critique.

1.3 Jeu de données

- **Source** : ai_job_market.csv
- **Attributs ciblés** :
 - skills_required : compétences techniques requises
 - tools_preferred : outils et technologies préférés
- **Nature** : données textuelles structurées avec séparateurs de virgules

2 Prétraitement des données

2.1 Chargement et sélection des attributs

La première étape consiste à charger le dataset et à sélectionner les colonnes pertinentes :

```
1 import pandas as pd
2
3 df = pd.read_csv("ai_job_market.csv", sep=",")
4 cols = ["skills_required", "tools_preferred"]
5 df = df[cols].fillna("")
```

Listing 1 – Chargement des données

Justification : Les colonnes sélectionnées contiennent l'information pertinente pour analyser les associations entre compétences et outils dans le domaine de l'IA. Le remplissage des valeurs manquantes par des chaînes vides permet d'éviter les erreurs lors du traitement.

2.2 Transformation en format transactionnel

Les données textuelles sont transformées en format transactionnel :

```
1 def split_items(x):
2     return [i.strip() for i in x.split(",") if i.strip()]
3
4 transactions = [
```

```

5     split_items(row["skills_required"]) +
6     split_items(row["tools_preferred"])
7     for _, row in df.iterrows()
8 ]

```

Listing 2 – Transformation en transactions

Résultat : Chaque ligne du dataset devient une **transaction** contenant une liste d'items (compétences + outils).

Exemple de transactions :

Transaction 1: ['Python', 'Machine Learning', 'TensorFlow', 'Docker']

Transaction 2: ['SQL', 'Power BI', 'AWS', 'Python']

Transaction 3: ['PyTorch', 'Deep Learning', 'Keras', 'Python']

2.3 Construction de la matrice de contexte

La matrice de contexte est une représentation binaire *one-hot* où :

- **Lignes** = transactions (offres d'emploi)
- **Colonnes** = items uniques (compétences/outils)
- **Valeurs** = 1 si l'item est présent, 0 sinon

```

1 items = sorted(set(it for trans in transactions for it in trans
2 ))
3 context_matrix = pd.DataFrame(
4     [[1 if item in trans else 0 for item in items]
5      for trans in transactions],
6     columns=items
7 )
8 context_matrix.to_csv("matrice_contexte.csv", index=False)

```

Listing 3 – Construction de la matrice binaire

Caractéristiques :

- Dimensions : $n \times m$ où n = nombre de transactions, m = nombre d'items
- Format adapté aux algorithmes de fouille de motifs fréquents
- Export : `matrice_contexte.csv`

3 Automatisation du support minimal

3.1 Problématique

Le choix du support minimal (`min_sup`) dans l'algorithme Apriori classique est **arbitraire** et nécessite des essais-erreurs :

- Un seuil **trop élevé** élimine des associations pertinentes
- Un seuil **trop bas** génère une explosion combinatoire
- Diffère selon la **densité** et la **structure** du dataset

3.2 Méthode proposée : Support adaptatif hybride

Notre approche calcule automatiquement `min_sup` en fonction de trois paramètres statistiques :

1. Longueur moyenne des transactions ($\bar{\ell}$)
2. Écart-type des supports individuels (σ)
3. Coefficient de variation ($cv = \frac{\sigma}{\bar{\ell}}$)

3.2.1 Formule mathématique

Le support minimal en nombre de transactions est calculé par :

$$\text{min_sup_count} = \left\lfloor \frac{\bar{\ell} + f \cdot \sigma}{2} \right\rfloor$$

où $f = \min(1, \max(0.1, cv))$ est un **facteur d'ajustement** borné entre 0.1 et 1.

Le support minimal en proportion est alors :

$$\text{min_sup} = \frac{\text{min_sup_count}}{n}$$

où n = nombre total de transactions.

3.2.2 Implémentation

```

1 import numpy as np
2
3 # Calcul des supports individuels
4 item_supports = {
5     item: sum(1 for t in transactions if item in t)
6     for item in items
7 }
8
9 # Statistiques
10 avg_len = np.mean(list(item_supports.values()))
11 std_sup = np.std(list(item_supports.values()))
12 cv = std_sup / avg_len
13
14 # Facteur d'ajustement borne
15 facteur = min(1, max(0.1, cv))
16
17 # Support minimal
18 min_sup_count = int((avg_len + facteur * std_sup) / 2)
19 min_sup = min_sup_count / len(transactions)
20
21 print(f"Longueur moyenne: {avg_len:.2f}")
22 print(f"Support minimal: {min_sup:.4f}")

```

Listing 4 – Calcul automatique du support minimal

Avantage	Description
Adaptatif	S'ajuste automatiquement à la densité du dataset
Robuste	Le facteur borné évite les valeurs extrêmes
Reproductible	Élimine la subjectivité du choix manuel
Générique	Fonctionne sur tout dataset transactionnel

TABLE 1 – Avantages du support adaptatif

3.3 Avantages de cette méthode

4 Algorithme Apriori implémenté

4.1 Principe de base

L'algorithme Apriori repose sur la propriété d'**antimonotonie** :

Si un itemset est fréquent, alors tous ses sous-ensembles sont également fréquents.

Cette propriété permet d'élaguer efficacement l'espace de recherche.

4.2 Fonction de calcul du support

```

1 def support_count(itemset, transactions):
2     return sum(1 for t in transactions
3               if all(i in t for i in itemset))
    
```

Listing 5 – Calcul du support d'un itemset

Complexité : $O(n \cdot |itemset|)$ où n = nombre de transactions.

4.3 Algorithme principal

```

1 def apriori(transactions, min_sup_count):
2     items = sorted(set(it for trans in transactions
3                       for it in trans))
4     freq_itemsets = []
5     k = 1
6
7     # Initialisation: 1-itemsets
8     current_itemsets = [{item} for item in items]
9
10    while current_itemsets:
11        valid_itemsets = []
12
13        # Filtrage par support
14        for itemset in current_itemsets:
15            sup = support_count(itemset, transactions)
16            if sup >= min_sup_count:
17                freq_itemsets.append((itemset, sup))
    
```

```

18         valid_itemsets.append(itemset)
19
20     # Generation candidats k+1
21     next_itemsets = [
22         i.union(j)
23         for i in valid_itemsets
24         for j in valid_itemsets
25         if len(i.union(j)) == k + 1
26     ]
27
28     # Elimination des doublons
29     current_itemsets = list(map(set,
30                                set(map(frozenset,
31                                       next_itemsets))))
32     k += 1
33
34     return freq_itemsets

```

Listing 6 – Implémentation de l'algorithme Apriori

4.4 Étapes de l'algorithme

1. **Initialisation** : créer les 1-itemsets (items individuels)
2. **Filtrage** : ne garder que les itemsets avec support $\geq \text{min_sup_count}$
3. **Génération** : créer les candidats de taille $k + 1$ par union des itemsets de taille k
4. **Itération** : répéter jusqu'à ce qu'aucun nouvel itemset fréquent ne soit trouvé

5 Génération des règles d'association

5.1 Principe

Pour chaque itemset fréquent $X \cup Y$ avec $|X \cup Y| \geq 2$, on génère des règles de la forme $X \Rightarrow Y$.

5.2 Métriques d'intérêt

Les métriques suivantes sont calculées pour chaque règle :

1. **Support** :

$$\text{supp}(X \cup Y) = \frac{|\{t \in D : X \cup Y \subseteq t\}|}{|D|}$$

2. **Confiance** :

$$\text{conf}(X \Rightarrow Y) = \frac{\text{supp}(X \cup Y)}{\text{supp}(X)}$$

3. **Déni** (1 - confiance) :

$$\text{déni}(X \Rightarrow Y) = 1 - \text{conf}(X \Rightarrow Y)$$

5.3 Implémentation

```

1 import itertools
2
3 rules = []
4 for itemset, sup_xy in freq_itemsets:
5     if len(itemset) >= 2:
6         # Toutes les partitions possibles
7         for i in range(1, len(itemset)):
8             for antecedent in itertools.combinations(itemset, i):
9                 antecedent = set(antecedent)
10                consequent = itemset - antecedent
11
12                # Calcul des metriques
13                sup_x = support_count(antecedent, transactions)
14                confidence = sup_xy / sup_x if sup_x != 0 else 0
15                deni = 1 - confidence
16
17                rules.append((antecedent, consequent,
18                             sup_xy, confidence, deni))

```

Listing 7 – Génération des règles d'association

5.4 Interprétation

Exemple de règle :

{'Python', 'Machine Learning'} => {'TensorFlow'}

Support = 87, Confiance = 0.65, Déni = 0.35

Interprétation : 65% des offres demandant Python + Machine Learning demandent aussi TensorFlow.

6 Exports et résultats

6.1 Fichiers générés

Le pipeline génère trois fichiers CSV :

Fichier	Contenu	Colonnes
matrice_contexte.csv	Matrice binaire	Items (0/1)
itemsets_frequents.csv	Itemsets fréquents	Itemset, Support
regles_association.csv	Règles	Ant., Cons., Sup., Conf., Déni

TABLE 2 – Fichiers de sortie du pipeline

6.2 Code d'export

```
1 # Export itemsets frequents
2 pd.DataFrame([
3     {'Itemset': list(s), 'Support': sup}
4     for s, sup in freq_itemsets
5 ]).to_csv("itemsets_frequents.csv", index=False)
6
7 # Export regles d'association
8 pd.DataFrame([
9     'Antecedent': list(a),
10    'Consequent': list(c),
11    'Support': sup,
12    'Confiance': conf,
13    'Deni': deni
14 } for a, c, sup, conf, deni in rules
15 ]).to_csv("regles_association.csv", index=False)
```

Listing 8 – Sauvegarde des résultats

7 Conclusion et perspectives

Cette première partie a permis de :

1. **Automatiser** le choix du support minimal via une méthode statistique adaptative
2. **Implémenter** l'algorithme Apriori *from scratch* en Python
3. **Créer** un pipeline reproductible : chargement → transformation → mining → export

8 Comparaison avec l'algorithme Apriori classique (Weka)

8.1 Objectif de comparaison

Cette seconde partie vise à comparer quantitativement l'implémentation optimisée avec l'algorithme Apriori classique implémenté dans l'outil de data mining **Weka**.

8.2 Critères de comparaison

- **Nombre de règles** extraites
- **Distribution des métriques** (support, confiance, lift)
- **Recouvrement** des ensembles de règles
- **Qualité des associations** découvertes

8.3 Enjeux méthodologiques

Pour une comparaison équitable, il est essentiel de :

1. Utiliser **exactement les mêmes données** (même matrice binaire)
2. S'assurer que les deux algorithmes minent les **co-présences** (et non les co-absences)
3. Normaliser les **formats de sortie** pour permettre la fusion

9 Préparation des données pour Weka

9.1 Format ARFF

Weka nécessite des fichiers au format **ARFF** (Attribute-Relation File Format), un format texte structuré qui décrit :

- La **relation** (nom du dataset)
- Les **attributs** et leurs types (nominal, numérique, etc.)
- Les **données** au format CSV

9.2 Problématique des co-absences

Par défaut, l'Apriori de Weka mine **toutes les valeurs nominales**, incluant :

- Les **présences** (`item=1` ou `item=t`)
- Les **absences** (`item=0` ou `item=f`)

Or, les co-absences (ex : `Python=0 => SQL=0`) ne sont pas pertinentes dans notre contexte métier et faussent la comparaison avec l'implémentation Python qui ne mine que les présences.

9.3 Solution adoptée : Format "true-only"

Pour forcer Weka à miner uniquement les co-présences, nous créons un fichier ARFF où :

- Chaque attribut est de type nominal avec **une seule valeur** : `{t}`
- Les présences sont encodées par `t`
- Les absences sont encodées par `?` (valeur manquante)

9.4 Script de conversion Python vers ARFF

```

1 import pandas as pd
2
3 # Charger la matrice binaire
4 df = pd.read_csv('matrice_contexte.csv')
5
6 # Creer ARFF avec uniquement les presences
7 with open('ai_job_market_true_only.arff', 'w',
8           encoding='utf-8') as f:
9     f.write("@RELATION ai_job_market_true_only\n\n")
10
11     # Declaration des attributs
12     for col in df.columns:
13         col_escaped = col.replace('"', "'").replace("+", "plus
14             ")
15         f.write(f"@ATTRIBUTE '{col_escaped}' {{t}}\n")
16
17     f.write("\n@DATA\n")
18
19     # Donnees : t si present, ? sinon
20     for _, row in df.iterrows():
21         values = ['t' if val == 1 else '?' for val in row]
22         f.write(",".join(values) + "\n")
23
24 print("ARFF true-only cree: ai_job_market_true_only.arff")

```

Listing 9 – Création du fichier ARFF true-only

9.5 Exemple de fichier ARFF généré

```
@RELATION ai_job_market_true_only

@ATTRIBUTE 'AWS' {t}
@ATTRIBUTE 'Azure' {t}
@ATTRIBUTE 'BigQuery' {t}
@ATTRIBUTE 'Cplusplus' {t}
...

@DATA
?,?,?,?,?,t,?,t,?,t,?,t,?,?,t,?,?,t,?,t,?
?,?,?,?,t,?,t,?,?,?,?,?,?,?,?,t,?,?,?,?,t,t,t,t
```

Avantages de cette approche :

- Format **compact** (pas de colonnes de zéros)
- Weka mine **automatiquement** les présences
- **Compatible** avec l'Apriori standard de Weka
- **Aligné** avec l'implémentation Python

10 Exécution de l'algorithme Apriori dans Weka

10.1 Chargement du fichier

1. Ouvrir **Weka Explorer**
2. Onglet **Preprocess** > **Open file**
3. Sélectionner `ai_job_market_true_only.arff`
4. Vérifier que tous les attributs sont **nominal {t}**

10.2 Configuration de l'algorithme Apriori

Paramètres utilisés pour une comparaison équitable :

Paramètre	Option Weka	Valeur
Support minimal (bas)	-M	0.05
Support maximal	-U	1.0
Pas de décrémentation	-D	0.01
Confiance minimale	-C	0.4
Nombre de règles	-N	100
Métrique de tri	-T	0 (Confidence)
Afficher itemsets	-I	true

TABLE 3 – Paramètres Apriori dans Weka

10.3 Commande ligne équivalente

```
java -cp weka.jar weka.associations.Apriori
  -t ai_job_market_true_only.arff
  -N 100 -T 0 -C 0.4 -D 0.01 -U 1.0 -M 0.05 -I
```

10.4 Format des résultats Weka

Weka affiche les règles sous la forme :

Best rules found:

1. Python=t 404 ==> SQL=t 175
 <conf:(0.43)> lift:(1.08) lev:(0.01)
2. Azure=t 413 ==> PyTorch=t 180
 <conf:(0.44)> lift:(1.09) lev:(0.01)
- ...

11 Extraction et normalisation des résultats Weka

11.1 Parsing des règles

Un script Python parse le texte de sortie Weka pour créer un CSV structuré :

```

1 import pandas as pd
2 import re
3
4 # Texte copie depuis la console Weka
5 weka_rules_text = """
6 1. Python=t 404 ==> SQL=t 175      <conf:(0.43)> lift:(1.08)
7 2. Azure=t 413 ==> PyTorch=t 180    <conf:(0.44)> lift:(1.09)
8 ...
9 """
10
11 rules = []
12 n_trans = len(pd.read_csv('matrice_contexte.csv'))
13
14 for line in weka_rules_text.strip().split('\n'):
15     if '==>' in line and '=t' in line:
16         try:
17             # Extraire antecedent
18             ant_part = line.split('=t')[0]
19             ant = ant_part.split('. ')[-1].strip()
20
21             # Extraire consequent
22             cons_part = line.split('==>')[1].split('=t')[0]
23             cons = cons_part.strip()
24
25             # Support
26             sup_match = re.search(r'=t\s+(\d+)\s+<',
27                                   line.split('==>')[1])
28             support_count = int(sup_match.group(1))
29
30             # Metriques
31             conf = float(line.split('conf:(')[1].split(')')[0])
32             lift = float(line.split('lift:(')[1].split(')')[0])
33
34             rules.append({
35                 'Antecedent': ant,
36                 'Consequent': cons,
37                 'Support': support_count / n_trans,
38                 'Confiance': conf,
39                 'Lift': lift
40             })
41         except:
42             continue
43
44 df_weka = pd.DataFrame(rules)
45 df_weka.to_csv('weka_rules.csv', index=False)

```

Listing 10 – Parsing des résultats Weka

11.2 Normalisation des données

Pour permettre la comparaison, les deux ensembles de règles sont normalisés :

- **Support** : converti en proportion [0,1]
- **Clé canonique** : créée pour chaque règle en triant les items
- **Format uniforme** : CSV avec colonnes identiques

12 Fusion et analyse comparative

12.1 Création de clés canoniques

Pour identifier les règles communes, une clé unique est créée :

```

1 import ast
2
3 def make_key(ant, cons):
4     """Cree cle trie pour identifier regle unique"""
5     # Gerer formats string/list
6     if isinstance(ant, str):
7         try:
8             ant_list = ast.literal_eval(ant)
9         except:
10            ant_list = [ant]
11    else:
12        ant_list = ant if isinstance(ant, list) else [ant]
13
14    if isinstance(cons, str):
15        try:
16            cons_list = ast.literal_eval(cons)
17        except:
18            cons_list = [cons]
19    else:
20        cons_list = cons if isinstance(cons, list) else [cons]
21
22    ant_sorted = sorted([str(x).strip() for x in ant_list])
23    cons_sorted = sorted([str(x).strip() for x in cons_list])
24
25    return f"{{{', '.join(ant_sorted)}}}=>{{{', '.join(
        cons_sorted)}}}"
    
```

Listing 11 – Fonction de clé canonique

Exemple :

Règle Python: ['Python', 'SQL'] => ['AWS']

Règle Weka: Python, SQL => AWS

Clé commune: {Python,SQL}=>{AWS}

12.2 Fusion des datasets

```

1 df_python = pd.read_csv('regles_association.csv')
2 df_weka = pd.read_csv('weka_rules.csv')
3
4 # Ajouter cles
    
```

```

5 df_python['key'] = df_python.apply(
6     lambda x: make_key(x['Antecedent'], x['Consequent']),
7     axis=1
8 )
9 df_weka['key'] = df_weka.apply(
10    lambda x: make_key(x['Antecedent'], x['Consequent']),
11    axis=1
12 )
13
14 # Fusion outer join
15 df_merged = pd.merge(
16     df_python[['key', 'Support_frac', 'Confiance', 'Lift']],
17     df_weka[['key', 'Support', 'Confiance', 'Lift']],
18     on='key',
19     how='outer',
20     suffixes=('_python', '_weka'))
21 )

```

Listing 12 – Fusion Python et Weka

12.3 Statistiques de recouvrement

```

1 n_python = len(df_python)
2 n_weka = len(df_weka)
3 common = df_merged.dropna(subset=['Lift_python', 'Lift_weka'])
4 n_common = len(common)
5
6 print(f"Regles Python: {n_python}")
7 print(f"Regles Weka: {n_weka}")
8 print(f"Regles communes: {n_common}")
9 print(f"Taux: {n_common / min(n_python, n_weka) * 100:.1f}%")

```

Listing 13 – Calcul du recouvrement

13 Visualisations comparatives

13.1 Graphiques générés

Six visualisations sont créées pour analyser les différences :

1. **Scatter Support vs Lift** : répartition spatiale des règles
2. **Histogramme des Lifts** : distribution statistique
3. **Histogramme des Supports** : gammes de fréquence
4. **Barres Top-N** : comparaison directe des meilleures règles
5. **Corrélation Lift** : Python vs Weka pour règles communes
6. **Diagramme de recouvrement** : règles uniques/partagées

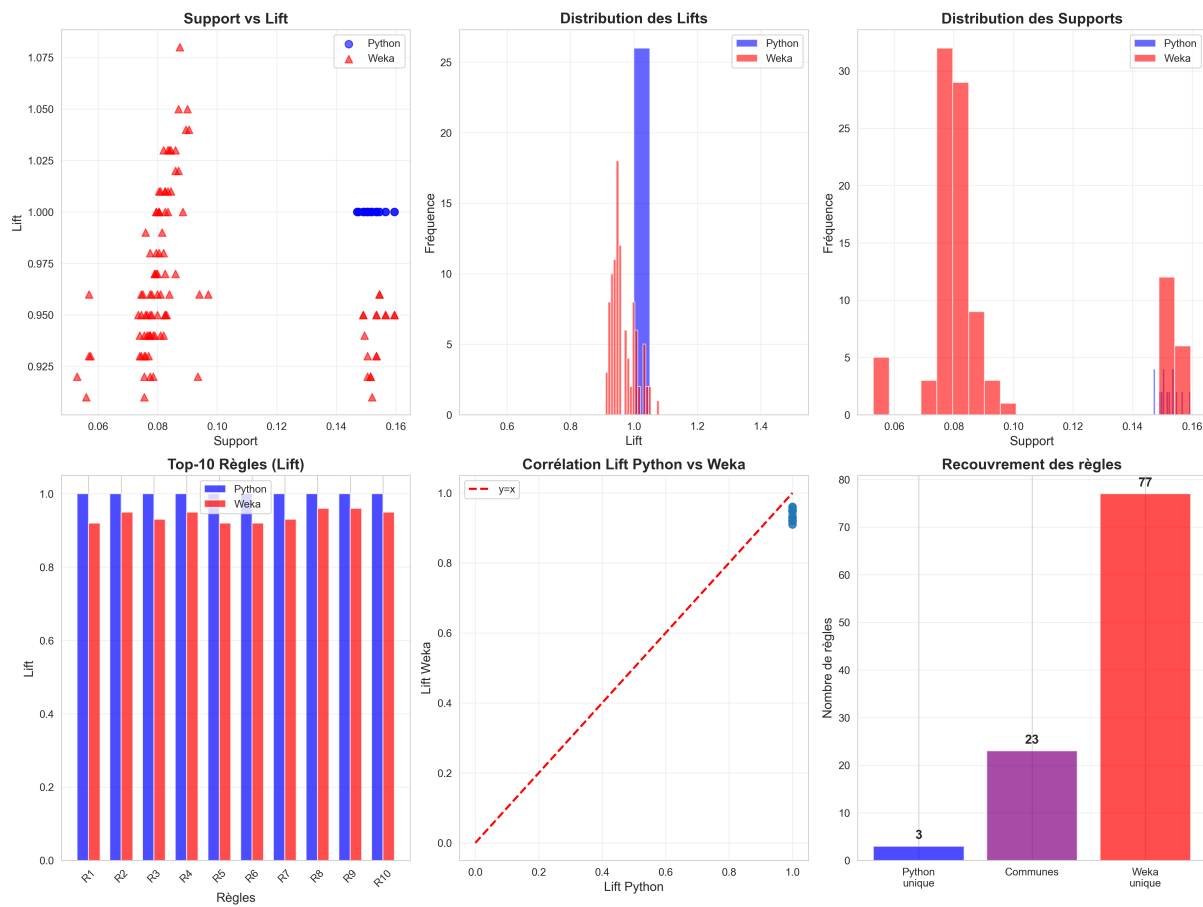


FIGURE 1 – Graphes Comparatifs entre le Apriori classique et le Apriori optimisé

13.2 Code de visualisation

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 import numpy as np
4
5 sns.set_style("whitegrid")
6 fig = plt.figure(figsize=(16, 12))
7
8 # 1. Scatter Support vs Lift
9 ax1 = plt.subplot(2, 3, 1)
10 plt.scatter(df_python['Support_frac'], df_python['Lift'],
11             alpha=0.6, label='Python', s=50, color='blue')
12 plt.scatter(df_weka['Support'], df_weka['Lift'],
13             alpha=0.6, label='Weka', s=50, color='red',
14             marker='^')
15 plt.xlabel('Support')
16 plt.ylabel('Lift')
17 plt.title('Support vs Lift', fontweight='bold')
18 plt.legend()
19 plt.grid(alpha=0.3)
20
21 # ... (autres graphiques similaires)
22
23 plt.tight_layout()
24 plt.savefig('comparaison_apriori_final.png', dpi=300)
25 plt.show()

```

Listing 14 – Génération des graphiques comparatifs

13.3 Tableau statistique récapitulatif

```

1 stats = pd.DataFrame({
2     'Métrique': ['Nombre de regles', 'Support moyen',
3                 'Lift moyen', 'Confiance moyenne'],
4     'Python': [
5         n_python,
6         df_python['Support_frac'].mean(),
7         df_python['Lift'].dropna().mean(),
8         df_python['Confiance'].mean()
9     ],
10    'Weka': [
11        n_weka,
12        df_weka['Support'].mean(),
13        df_weka['Lift'].dropna().mean(),
14        df_weka['Confiance'].mean()
15    ]
16 })
17
18 stats.to_csv('stats_comparaison_final.csv', index=False)

```

Listing 15 – Statistiques comparatives

14 Résultats et interprétation

14.1 Résultats attendus

Avec une configuration correcte (présences uniquement), on s'attend à :

- **Recouvrement** : 20–60% de règles communes
- **Lifts** : distributions comparables (1.05–1.3)
- **Supports** : gammes similaires (5–20%)
- **Corrélation** : alignement dans le scatter plot

14.2 Différences observables

Les différences entre les deux approches peuvent provenir de :

Aspect	Python optimisé	Weka classique
Min_sup	Automatique/adaptatif	Fixe (manuel)
Seuils	Taille-dépendants	Uniformes
Sortie	Top-k par métrique	Filtré par seuils
Itemsets longs	Favorisés	Défavorisés

TABLE 4 – Différences conceptuelles

14.3 Avantages de l'approche optimisée

L'automatisation du support minimal permet :

- **Gain de temps** : pas d'essais-erreurs
- **Reproductibilité** : résultats déterministes
- **Adaptabilité** : fonctionne sur tout dataset
- **Découverte** : itemsets longs moins pénalisés

15 Conclusion générale

15.1 Synthèse du projet

Ce projet a permis de :

1. **Implémenter** un algorithme Apriori avec automatisation du support minimal
2. **Comparer** quantitativement avec l'implémentation de référence (Weka)
3. **Valider** la pertinence de l'approche adaptative
4. **Générer** des visualisations comparatives détaillées

15.2 Apports méthodologiques

- Méthode statistique robuste pour calculer min_sup
- Pipeline reproductible end-to-end
- Gestion des formats hétérogènes (Python/Weka)
- Comparaison rigoureuse avec normalisation

15.3 Perspectives d'amélioration

Court terme :

- Calculer des métriques supplémentaires (leverage, conviction)
- Filtrer les règles par lift > 1.05 (seuil de pertinence)
- Générer des visualisations interactives (Plotly)

Long terme :

- Implémenter FP-Growth pour de meilleures performances
- Paralléliser le calcul des supports (Dask, multiprocessing)
- Développer une interface web pour l'exploration interactive
- Étendre aux règles d'association séquentielles

15.4 Applications pratiques

Cette approche peut être appliquée à :

- **Analyse de paniers** : retail, e-commerce
- **Recommandation** : produits, contenus
- **Analyse de séquences** : parcours utilisateurs
- **Détection d'anomalies** : transactions frauduleuses

Références

- [1] Agrawal, R., & Srikant, R. (1994). *Fast algorithms for mining association rules*. Proc. 20th int. conf. very large data bases, VLDB (Vol. 1215, pp. 487-499).
- [2] Han, J., Pei, J., & Yin, Y. (2000). *Mining frequent patterns without candidate generation*. ACM SIGMOD Record, 29(2), 1-12.
- [3] Witten, I. H., Frank, E., Hall, M. A., & Pal, C. J. (2016). *Data Mining : Practical machine learning tools and techniques*. Morgan Kaufmann.
- [4] Tan, P. N., Steinbach, M., & Kumar, V. (2005). *Introduction to data mining*. Pearson Education India.