

Université des Sciences et de la Technologie Houari Boumediene  
DÉPARTEMENT D'IA ET SD



Master Informatique visuelle

---

Apriori

**Apriori**

---

KASHI Thiziri  
MESSAOUDENE Lydia

## Contents

<b>1</b>	<b>Introduction à l'algorithme Apriori et ses applications</b>	<b>3</b>
1.1	Introduction . . . . .	3
1.2	Description de l'algorithme Apriori . . . . .	3
<b>2</b>	<b>Fonctionnement de l'algorithme Apriori</b>	<b>4</b>
2.1	Itemset . . . . .	4
2.2	Frequent Itemset . . . . .	4
2.3	Association Rules . . . . .	4
2.4	Support et confidence . . . . .	4
2.5	Etape : . . . . .	4
2.6	Generate Association Rules: . . . . .	6
<b>3</b>	<b>Limitations de l'algorithme Apriori</b>	<b>7</b>
<b>4</b>	<b>Présentation de l'algorithme Close et ses avantages</b>	<b>8</b>
4.1	groupes fermés d'éléments . . . . .	8
4.2	Etape . . . . .	8
4.3	Comparaison de l'algorithme close a apriori . . . . .	8
<b>5</b>	<b>Amélioration du calcul du support dans l'algorithme Apriori</b>	<b>9</b>
<b>6</b>	<b>Utilisation de la corrélation pour améliorer l'efficacité de l'algorithme Apriori</b>	<b>10</b>
<b>7</b>	<b>Utilisation de l'algorithme de clustering en conjonction avec l'algorithme Apriori</b>	<b>11</b>
<b>8</b>	<b>temps d'execution</b>	<b>11</b>

# 1 Introduction à l'algorithme Apriori et ses applications

## 1.1 Introduction

L'algorithme Apriori a été proposé en 1994 par Agrawal et Srikant comme une méthode d'exploration de données pour extraire des règles d'association à partir de données transactionnelles. Depuis lors, il est devenu un outil important dans de nombreux domaines tels que la gestion de la relation client, la recherche de marché et la publicité en ligne.

## 1.2 Description de l'algorithme Apriori

L'algorithme Apriori est un algorithme utilisé pour calculer les règles d'association entre des objets, En d'autres termes , comment deux ou plusieurs objets sont liés les uns aux autres. c'est un apprentissage de règles d'association qui analyse les personnes qui ont acheté le produit A ont également acheté le produit B.

L'objectif principal de l'algorithme Apriori est de créer la règle d'association entre différents objets. La règle d'association décrit comment deux ou plusieurs objets sont liés les uns aux autres. L'algorithme Apriori est également appelé fouille de motifs fréquents.

Prenons un exemple concret : imaginez que vous allez faire vos courses a UNO et que vous achetez différents produits. L'algorithme Apriori peut être utilisé pour déterminer les règles d'association entre ces produits, par exemple, que les personnes qui ont acheté du lait ont également acheté du pain. Cela peut aider les clients à acheter leurs produits plus facilement et augmenter les performances de vente de UNO.

Dans ce rapport, nous allons nous intéresser plus en détail à l'algorithme Apriori et à ses applications dans différents domaines, tels que la gestion de la relation client, la recherche de marché et la publicité en ligne. Nous examinerons également les limites de l'algorithme Apriori et les améliorations qui peuvent être apportées, comme l'utilisation des l'algorithmes Close , FP-Growth, pour surmonter certaines de ces limites. Enfin, nous discuterons également l'utilisation de l'amelioration du support , de l'utulisation de la correlation ou meme de fussionner apriori avec l'algorithme de clustering

## 2 Fonctionnement de l'algorithme Apriori

Pour mieux comprendre le fonctionnement d'Apriori, il faut d'abord comprendre quelques mots :

### 2.1 Itemset

: en anglais "a set of items together is called an itemset", Itemset est un ensemble de 2 ou plus items

### 2.2 Frequent Itemset

un ensemble d'items est dit fréquent si son support et sa confiance est supérieur à un certain seuil minimum

### 2.3 Association Rules

aussi appelées les règles d'association en français, elles sont utilisées pour déterminer la relation entre des attributs dans un dataset, exemple :  $A \Rightarrow B$ , cela veut dire que certaine valeur de A détermine la valeur B

### 2.4 Support et confiance

$$\text{Support (A)} = \frac{\text{Number of transaction in which A appears}}{\text{Total number of transactions}}$$
$$\text{Confidence (A} \rightarrow \text{B)} = \frac{\text{Support(A} \cup \text{B)}}{\text{Support(A)}}$$

### 2.5 Etape :

Prenons un exemple simple supposant que dans un supermarché nous vendons les produits suivants : Riz, huile, lait, poulet, pomme.

Nous allons analyser 6 transactions de différents clients, nous considérerons le 0 comme produit non-acheté et le 1 pour produit acheté

Transaction ID	Riz	poulet	huile	lait	pomme
t1	1	1	1	0	0
t2	0	1	1	1	0
t3	0	0	0	1	1
t4	1	1	0	1	0
t5	1	1	1	0	1
t6	1	1	1	1	1

L'algorithme Apriori repose sur les bases suivantes :

- Si  $P(I)$  est inférieur au seuil de support minimum, alors  $I$  n'est pas considéré comme fréquent.
- Si  $P(I+A)$  est inférieur au seuil de support minimum, alors  $I+A$  n'est pas considéré comme fréquent, où  $A$  appartient également à l'ensemble d'éléments.
- Si un ensemble d'éléments a une valeur inférieure au support minimum, alors tous ses ensembles supérieurs seront également inférieurs au support minimum et pourront donc être ignorés. Cette propriété est appelée **propriété d'antimonotonie**.

#### Étape 1 : Génération des ensembles de fréquences :

Dans cette étape, l'algorithme scanne la base de données transactionnelle pour identifier les ensembles d'éléments fréquents. Un ensemble d'éléments est considéré comme fréquent s'il apparaît dans un nombre suffisant de transactions, selon un seuil de support défini à l'avance.

dans notre cas nous avons la table suivante qui indique les produits les plus achetés par le client (support > 50)

Support threshold = 50  $\Rightarrow 0.5 * 6 = 3 \Rightarrow \text{minsup} = 3$

Produit	fréquence
Riz	4
poulet	5
huile	4
lait	4

#### Étape 2: Création d'un itemset à partir des items restants de l'étape précédente

Produit	fréquence
RP	4
RO	3
RM	2
PO	4
PM	3
OM	2

Nous nous débarrassons des items en dessous du seuil, il nous reste RP, RO, PO, PM

Création d'un itemset à partir des items restants de l'étape précédente

Produit	fréquence
RPO	4
POM	2

et donc RPO sont les items les plus fréquents

## 2.6 Generate Association Rules:

dans cette etapes nous devons genere toute les regles possible pour cette exemple ca sera :

$R, P \Rightarrow O, R, O \Rightarrow P, R \Rightarrow P, O, O \Rightarrow R, P, P, O \Rightarrow R, P \Rightarrow R, P$

Nous calculerons "confidence" avec la formule site au paravant et nous eliminerons les regles nons segnifiante grace au seuil de confidentialite .

### 3 Limitations de l'algorithme Apriori

Bien que l'algorithme Apriori soit une technique très populaire pour l'exploration des données, il présente certains inconvénients, tels que :

- **Coût de calcul :** La technique Apriori prend beaucoup de temps à calculer car elle doit générer et examiner de nombreux ensembles de données pour dériver les règles d'association.
- **Dimensions de la base de données :** La méthode Apriori peut rencontrer des problèmes lorsqu'elle est appliquée à des ensembles de données comportant un volume élevé de transactions, car elle doit examiner chaque transaction pour trouver des ensembles d'éléments récurrents.
- **Dépendance aux paramètres :** L'algorithme Apriori dépend fortement de facteurs tels que les seuils de soutien et de confiance. Le choix de ces paramètres peut avoir un impact significatif sur les résultats de l'algorithme
- **Incapacité à gérer des données complexes :** L'algorithme Apriori est conçu pour traiter des données transactionnelles simples, telles que les achats de produits dans un magasin. Il peut rencontrer des difficultés lorsqu'il traite des données plus complexes, telles que des données textuelles brutes ou des données multidimensionnelles.
- **Problèmes liés à la dimensionnalité :** L'algorithme Apriori peut souffrir d'un mauvais calcul de la dimensionnalité, ce qui signifie que le coût du calcul augmente de manière exponentielle avec le nombre de dimensions, rendant l'approche inefficace pour les grands ensembles de données.

Ces limites doivent être prises en compte lors de l'utilisation de l'algorithme d'Apriori, et d'autres approches peuvent être envisagées pour des ensembles de données plus importants ou plus compliqués et c'est ce que nous allons voir dans les prochains chapitre

## 4 Présentation de l'algorithme Close et ses avantages

La méthode close est une variante de l'algorithme Apriori qui identifie les groupes d'éléments fréquents et fixes dans une base de données. Les ensembles d'éléments fermés sont des groupes d'éléments qui ne peuvent pas être étendus tout en conservant leur support et leur fermeture. En d'autres termes, ils ne peuvent pas être améliorés sans changer leur signification.

### 4.1 groupes fermés d'éléments

Un groupe d'éléments est considéré comme fermé s'il ne peut être étendu sans modifier son support et sa fermeture.

### 4.2 Etape

L'algorithme Close suit les mêmes étapes de base que l'algorithme Apriori, mais il ajoute des étapes supplémentaires pour identifier les groupes fermés d'éléments. Les étapes supplémentaires sont les suivantes :

- Générer tous les groupes d'éléments fréquents, comme dans la méthode Apriori.
- Vérifier si chaque groupe d'éléments fréquents est fermé ou non.
- Supprimez tous les groupes d'éléments qui ne sont pas fermés
- Trouvez tous les sous-ensembles fermés de chaque ensemble fermé d'éléments trouvé.
- Supprimer les sous-ensembles fermés qui partagent le même support que leur ensemble parent.
- Renvoyer les groupes d'éléments fermés restants.

### 4.3 Comparaison de l'algorithme close a apriori

La méthode Close et l'algorithme Apriori sont tous deux des algorithmes de fouille de données qui sont utilisés pour identifier des modèles dans un ensemble de données. Néanmoins, il existe des différences significatives entre les deux algorithmes.

L'algorithme Apriori est souvent considéré comme plus simple et plus rapide que l'algorithme Close car il utilise une approche de recherche **large** qui implique l'exploration de toutes les combinaisons d'éléments potentielles dans un premier temps, suivie de l'exclusion des combinaisons peu fréquentes. Close, quant à lui, utilise une approche de recherche **approfondie** qui examine chaque combinaison d'éléments en profondeur et détermine si elle est fermée ou non.

Cependant, l'approche Close présente l'avantage de générer des ensembles fermés d'éléments, qui peuvent être plus utiles pour l'analyse des données que les ensembles ouverts d'éléments de l'algorithme Apriori. Les groupes d'éléments fermés sont des groupes d'éléments qui ne peuvent être étendus sans modifier leur support et leur fermeture, ce qui implique qu'ils représentent des modèles de données plus significatifs et précis



## 5 Amélioration du calcul du support dans l'algorithme Apriori

Le choix du calcul du support dans la méthode Apriori est principalement déterminé par la taille de la base de données et la mémoire disponible.

Si la base de données est relativement petite, il est possible de stocker toutes les transactions en mémoire et de parcourir la base de données plusieurs fois pour calculer le support de chaque itemet. Cette méthode est plutôt basique et ne nécessite pas de techniques avancées de gestion de la mémoire ou de parallélisation.

Néanmoins, pour les bases de données plus importantes, le stockage de toutes les transactions en mémoire peut s'avérer peu pratique. Dans ce cas, des structures de données telles que les tables de hachage peuvent être utilisées pour stocker les transactions de manière efficace et compacte. Cette méthode permet de réduire l'utilisation de la mémoire et d'augmenter le temps de calcul du support.

## 6 Utilisation de la corrélation pour améliorer l'efficacité de l'algorithme Apriori

L'utilisation de la corrélation peut être une méthode efficace pour améliorer l'efficacité de l'algorithme Apriori. En effet, la corrélation peut aider à identifier les itemsets qui ont tendance à apparaître ensemble, même s'ils ne répondent pas toujours à l'exigence de support minimal.

L'idée est de calculer une mesure de corrélation entre chaque paire d'itemsets et d'utiliser cette mesure pour identifier les itemsets ayant une forte corrélation. Ces ensembles peuvent ensuite être combinés pour former des ensembles plus importants sans qu'il soit nécessaire de vérifier leur support individuel.

Plusieurs mesures de corrélation peuvent être utilisées, notamment le coefficient de corrélation de Pearson, le coefficient de corrélation de Spearman et le coefficient de corrélation de Kendal-Tau. Ces mesures peuvent être calculées en utilisant la fréquence d'apparition de chaque paire d'itemsets dans la base de données.

L'utilisation de la corrélation peut contribuer à réduire le temps de calcul et le nombre d'ensembles d'éléments candidats générés par la méthode Apriori tout en conservant une grande précision. Cependant, cette méthode peut ne pas être applicable à toutes les applications car elle suppose que les itemsets sont fortement corrélés.

## 7 Utilisation de l'algorithme de clustering en conjonction avec l'algorithme Apriori

L'utilisation de l'algorithme de regroupement en conjonction avec l'algorithme Apriori peut être une méthode efficace pour améliorer l'efficacité et la pertinence des résultats obtenus.

L'idée est d'utiliser un algorithme de clustering pour regrouper les transactions similaires en grappes, puis d'appliquer la méthode Apriori à chaque grappe individuellement plutôt qu'à l'ensemble des transactions. Cela permet de réduire le nombre de transactions à considérer, et donc le temps nécessaire au calcul de la méthode Apriori.

En outre, en regroupant des transactions similaires, la méthode de regroupement peut aider à identifier des corrélations plus significatives entre les articles. Par exemple, les articles fréquemment achetés ensemble dans un certain groupe peuvent présenter une corrélation plus forte que s'ils étaient considérés dans l'ensemble des transactions.

Toutefois, l'utilisation d'un algorithme de regroupement peut poser des problèmes si les regroupements ne sont pas représentatifs de l'ensemble des transactions. En outre, la technique de regroupement et ses paramètres doivent être adaptés à chaque application en fonction de la nature des données.

## 8 temps d'exécution

Nous avons testé sur 2 datasets :

- la dataset titanic de 891 lignes nous avons obtenu un temps d'exécution de 22s pour apriori
- la dataset wether qui contient 20 lignes nous avons obtenu un temps de 0.00001s pour l'algorithme apriori