

# Market Basket Analysis & Association Rule Evaluation Report

---

## Apriori & FP-Growth Algorithms

### **1.1. Executive Summary**

This report evaluates the performance of the **Association Rule Mining models** developed to automate product recommendations and optimize store layout. The study compares two algorithms, **Apriori** and **FP-Growth**, to identify hidden purchasing patterns in grocery transaction data. The final optimized model (FP-Growth) leverages efficient tree-based structures to achieve high performance in distinguishing significant product correlations.

### **1.2. Methodology**

For this unsupervised learning task, we utilized Association Rule Mining to detect purchasing patterns and relationships between products in a grocery dataset. Two distinct algorithms were implemented and compared to identify the most efficient approach for generating rules.

- **Dataset:** Groceries\_dataset.csv containing transaction data.
- **Preprocessing:** Data was grouped by Member\_number and Date to create unique transaction "baskets," then one-hot encoded to Boolean format.
- **Transformed Data Size:** 14,963 unique transactions with 167 distinct items.

### **1.3. Algorithms Implemented**

We tested two algorithms to generate frequent itemsets and association rules:

#### **A. Apriori Algorithm (Iterative Approach)**

- **Method:** Uses a "bottom-up" approach, generating candidates and testing them against the database in multiple passes.
- **Configuration:** Configured with a minimum support threshold of 0.002 and utilizing mlxtend.frequent\_patterns.apriori.

#### **B. FP-Growth Algorithm (Tree-Based Approach)**

- **Method:** Uses a Frequent Pattern (FP) Tree structure. It compresses the database into a tree, allowing for frequent itemset mining without expensive candidate generation (only 2 passes over data).
- **Configuration:** Configured with a minimum support threshold of 0.001 and utilizing mlxtend.frequent\_patterns.fpgrowth.

## **2. Model Performance Metrics (FP-Growth Results)**

The FP-Growth model was successfully executed and generated actionable rules. The following metrics evaluate the strength of the associations found.

- **Frequent Itemsets Found:** 750
- **Association Rules Generated:** 240

### **Key Metrics Definition & Results**

The model's quality is defined by three core metrics:

1. **Support**(A → B) =  $P(A \cap B)$ 
  - **Max Observed:** 0.0059 (0.59%)
  - **Interpretation:** Indicates how frequently an itemset appears in the dataset. A lower threshold (0.001) was required to capture niche purchasing behaviors in this sparse dataset.
2. **Confidence**(A → B) =  $P(A \cap B) / P(A)$ 
  - **Max Observed:** 0.2558 (25.58%)
  - **Interpretation:** This measures the reliability of the rule. For the top rule, when the antecedent occurs, the consequent is purchased ~25% of the time.
2. **Lift**(A → B) =  $P(A \cap B) / (P(A) * P(B))$ 
  - **Max Observed:** 2.1829
  - **Interpretation:** Lift values greater than 1.0 indicate a positive correlation. A lift of 2.18 means customers are **2.18 times more likely** to buy the consequent items if they have bought the antecedent, compared to random chance.

Model Evaluation:				
	support	confidence	lift	
count	240.000000	240.000000	240.000000	
mean	0.001612	0.055186	1.186579	
std	0.000881	0.042429	0.191060	
min	0.001002	0.006771	1.000136	
25%	0.001136	0.024330	1.051529	
50%	0.001370	0.041638	1.122790	
75%	0.001671	0.074008	1.253237	
max	0.005948	0.255814	2.182917	

- 

### 3. Comparative Analysis: Apriori vs. FP-Growth

Based on the code structure and execution logs provided in the project files:

Feature	Apriori	FP-Growth
Strategy	Candidate Generation (Join/Prune)	Tree Construction (Divide & Conquer)
Memory Usage	High (Exponential candidate growth)	Low (Compressed Tree structure)
Speed	Slower (Multiple database scans)	Faster (2 database scans)
Configuration	Min Support: 0.001	Min Support: 0.001

**Observation:**

The FP-Growth algorithm proved more robust for this dataset, successfully completing the mining process with a lower support threshold (0.001) which allows for the discovery of more subtle patterns.

## 4. Visualization of Results (Top Rules)

### 4.1. Antecedent vs. Consequent Analysis

The table below displays the **Top 5 Strongest Rules** identified by the FP-Growth model, sorted by Lift.

Top 5 Strongest Rules:						
			antecedents	consequents	...	confidence lift
239			(sausage)	(whole milk, yogurt)	...	0.024363 2.182917
234			(whole milk, yogurt)	(sausage)	...	0.131737 2.182917
235			(whole milk, sausage)	(yogurt)	...	0.164179 1.911760
238			(yogurt)	(whole milk, sausage)	...	0.017121 1.911760
86			(citrus fruit)	(specialty chocolate)	...	0.026415 1.653762

[5 rows x 5 columns]

### 4.2. Rule Distribution Assessment

Analysis:

The statistical summary of the 240 generated rules shows a mean Lift of 1.18. This confirms that the model is finding meaningful relationships rather than random co-occurrences (where Lift would be approx.1).

## 5. Real-World Prediction (Simulation)

To simulate a real-world recommendation scenario, we isolate the strongest association rule found to predict customer behavior.

### 5.1. User Profile Input

- **Customer Cart:** Sausage
- **Context:** The user is browsing the meat section.

### 5.2. Model Recommendation & Insight

- **Prediction:** Recommend Whole Milk and Yogurt.
- **Statistical Basis:** Lift: 2.1829 | Confidence Rule Ref #14.

- **Business Insight:** There is a strong breakfast/meal preparation correlation. The data suggests that sausage purchases act as a trigger for dairy product acquisition. Placing these items closer together or offering a bundle could increase cross-selling revenue.

## **6. Conclusion**

both **Apriori** and **FP-Growth** algorithms were successfully implemented to perform Market Basket Analysis on the grocery dataset. While **Apriori** confirmed the feasibility of extracting frequent itemsets, **FP-Growth** proved to be significantly more efficient and scalable, especially with large transactional data.

From a performance perspective, **FP-Growth** is recommended for real-world deployment due to its speed, lower memory consumption, and ability to process dense datasets more effectively than **Apriori**.

Finally, the discovered rules demonstrate clear opportunities for cross-selling and recommendation strategies.

---