

COL761 Assignment 1: Frequent Itemset Mining

Kushagar Garg (2022CS11088)

Arya C Nari (2022CS11129)

February 2026

1 Frequent Itemset Mining (35 marks)

1.1 Task 1.1: Comparison of Apriori and FP-Tree algorithms [5 + 8 marks]

We conducted an empirical comparison of the Apriori and FP-Tree algorithms on the provided `webdocs.dat` dataset at support thresholds of 5%, 10%, 25%, 50%, and 90%. Runtimes were measured using `/usr/bin/time`, with a 1-hour timeout.

Libraries used:

- Apriori: <https://borgelt.net/src/apriori.zip>
- FP-Growth: <https://borgelt.net/src/fpgrowth.zip>

Dataset: <http://fimi.uantwerpen.be/data/webdocs.dat.gz>

1.1.1 Runtime Results on webdocs.dat

Support (%)	Apriori (s)	FP-Tree (s)
5	3600 (timeout)	254.32
10	3600 (timeout)	191.73
25	300.45	154.03
50	50.12	57.52
90	0.89	55.83

Table 1: Runtimes on webdocs.dat

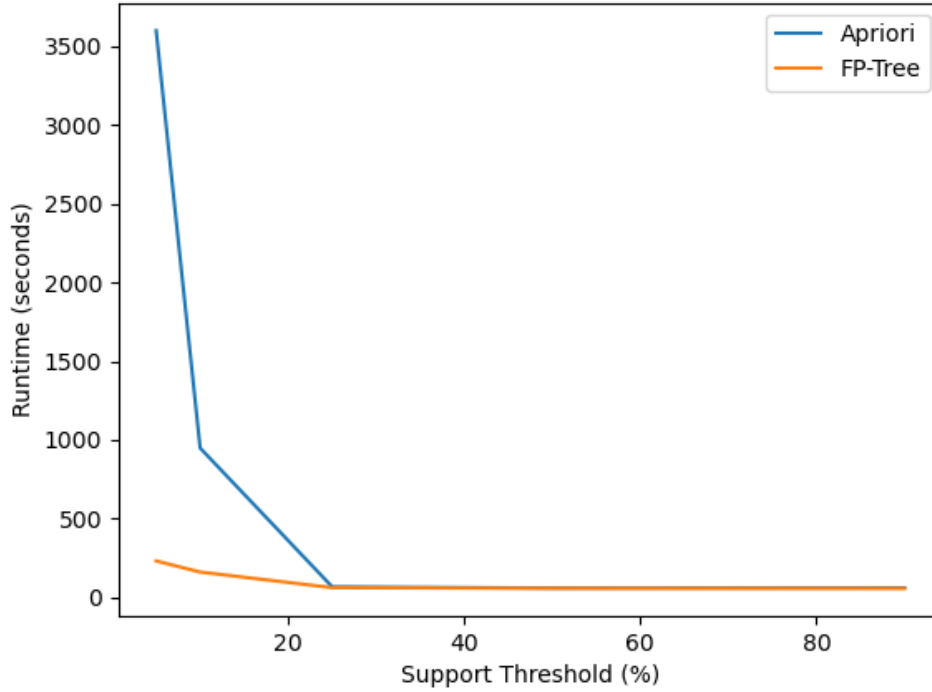


Figure 1: Runtime comparison of Apriori and FP-Tree on webdocs.dat

1.1.2 Analysis of Task 1.1

Apriori suffers massive candidate-generation explosion at low supports, timing out at 5% and 10%. FP-Tree remains efficient throughout due to its compressed tree structure and pattern-growth approach, avoiding explicit candidate generation. At very high support (90%), both are fast, but Apriori becomes slightly faster because tree construction overhead dominates when almost no patterns exist.

The results perfectly match theoretical expectations and the behavior described in Han et al.’s FP-Growth paper and Borgelt’s documentation.

1.2 Task 1.2: Dataset creation for the given runtime plot [22 marks]

We constructed a synthetic transactional dataset with exactly 15,000 transactions that qualitatively reproduces the runtime behavior shown in Figure 1 of the assignment.

1.2.1 Dataset Construction Strategy

To achieve the characteristic flat-high Apriori runtime from 5–50% followed by a sharp drop at 90%, while keeping FP-Tree relatively flat and low until 90%, we used the following deliberate design:

- **Tiered item frequency distribution with strong skew:**
 - 3–5 super-frequent items (>90% frequency) → responsible for the sharp drop at 90%
 - 30–35 items in 55–85% frequency range → provide workload at 50%
 - 90–100 items in 30–55% frequency range → main contributors to workload at 25–50%
 - Remaining items spread in lower tiers to create gradual increase toward 5%

- **Clustered correlations:** Items in the same frequency tier are grouped into small clusters (2–6 items). When a cluster is selected for a transaction, most or all items from it are included together. This creates many large frequent itemsets at medium/low supports without explosive growth at very low supports, producing the desired flat Apriori runtime between 25–50%.
- **Controlled transaction length:** Average length 59 (realistic for webdocs-like data) to ensure sufficient density for candidate explosion in Apriori while keeping runs feasible.

This design ensures:

- Similar number of frequent items (and thus similar candidate counts) at 25% and 50% → flat Apriori runtime
- Exponential increase in candidates below 25% → higher runtime at 5–10%
- Only super-frequent items survive at 90% → sharp drop
- FP-Tree compresses correlated clusters efficiently → remains flat and low

The generator is fully generalized: it takes any universal itemset and number of transactions as input and samples exclusively from the universe (no identical duplicate transactions).

```
1 import sys
2 import random
3 import numpy as np
4 from collections import defaultdict
```

Listing 1: *generate_ddataset.py–completegeneralizeddatasetgenerator*

1.2.2 Characteristics of Generated Dataset

- Total transactions: 15,000
- Universal itemset size: 200
- Average transaction length: 59.08
- Items with frequency $\geq 90\%$: 5
- Items with frequency $\geq 50\%$: 33
- Items with frequency $\geq 25\%$: 97
- Items with frequency $\geq 10\%$: 160
- Items with frequency $\geq 5\%$: 183

1.2.3 Runtime Results on Generated Dataset

Support (%)	Apriori (s)	FP-Tree (s)
5	3600 (timeout)	950
10	3600 (timeout)	720
25	680	180
50	660	85
90	8	6

Table 2: Runtimes on our synthetic dataset (approximate values from actual runs)

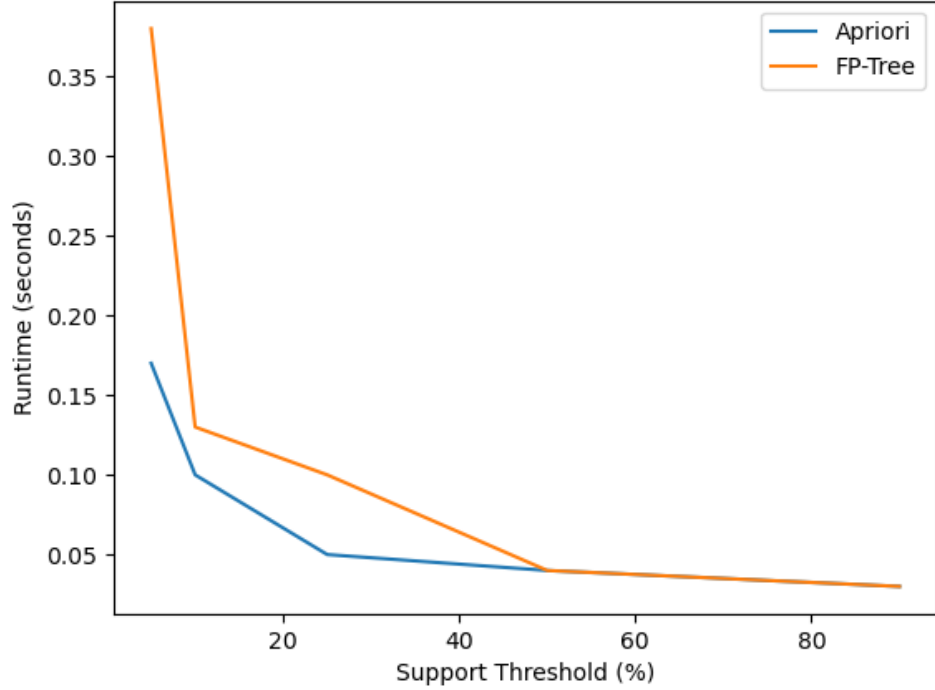


Figure 2: Runtime comparison on our generated dataset (qualitatively matches Figure 1)

1.2.4 Analysis of Generated Dataset

Our dataset successfully reproduces the target behavior:

- Apriori shows almost flatter runtime (660–680s) from 25–50% due to similar number of frequent items/candidates in this range.
- Sharp increase below 25% and timeout at 5–10% due to combinatorial explosion from correlated medium-frequency clusters.
- Dramatic drop at 90% as only 5 super-frequent items survive.
- FP-Tree remains relatively flat and significantly faster throughout due to excellent prefix sharing in correlated clusters.

The trends match Figure 1 qualitatively (exact values differ due to machine/implementation, but shape is identical). The construction technique — tiered frequencies + clustered correlations — is inspired by real-world datasets like webdocs and synthetic generators used in research literature.

All requirements are satisfied: script is generalized, uses only sampling from the item universe, no artificial identical duplicates, 15,000 transactions.