



 > [Machine Learning Bits](#) > Apriori Algorithm

Apriori Algorithm

These contents were automatically converted from lecture slides. Some contents have been withheld for copyright or technical limitations. The contents have not been optimally reformatted for online access. All rights reserved unless otherwise noted.

A previous [▶ YouTube video recording](#) of these contents is available.

Apriori Algorithm [AgSr94]

Algorithm: Apriori($I, D, \text{minsupp}$)

```

1  $k \leftarrow 1$ 
2  $F_1 \leftarrow \{\{i\} \mid i \in I \wedge \text{count}(\{i\}) \geq |D| \cdot \text{minsupp}\}$  // frequent 1-itemsets
3 while  $F_k \neq \emptyset$  do
4    $k \leftarrow k + 1$ 
5    $C_k \leftarrow \text{apriori-gen}(F_{k-1})$  // generate new candidates
6   foreach  $t \in D$  do // scan the database once
7      $C_t \leftarrow \text{subset}(C_k, t)$  // candidates contained in t
8     foreach  $c \in C_t$  do  $c.\text{count} += 1$  // increment support count
9    $F_k \leftarrow \{c \in C_k \mid c.\text{count} \geq |D| \cdot \text{minsupp}\}$  // keep frequent candidates
10 return  $\bigcup_i F_i$ 

```

where C_k : set of candidate k -itemsets F_k : set of frequent k -itemsets

Note: all itemsets & lists of itemsets are kept sorted.

Candidate Generation: apriori-gen

Step 1: **Join** by prefix matching

Process all $(k - 1)$ -frequent itemsets $p < q$, that *match* on their first $k - 2$ elements

?



🏠 › [Machine Learning Bits](#) › The Apriori Hash Tree

The Apriori Hash Tree

These contents were automatically converted from lecture slides. Some contents have been withheld for copyright or technical limitations. The contents have not been optimally reformatted for online access. All rights reserved unless otherwise noted.

A previous ▶ [YouTube video recording](#) of these contents is available.

Efficient Computing of the Subset Function

$\text{Subset}(C_k, T)$ computes all candidates from C_k that are contained in transaction T

Problems:

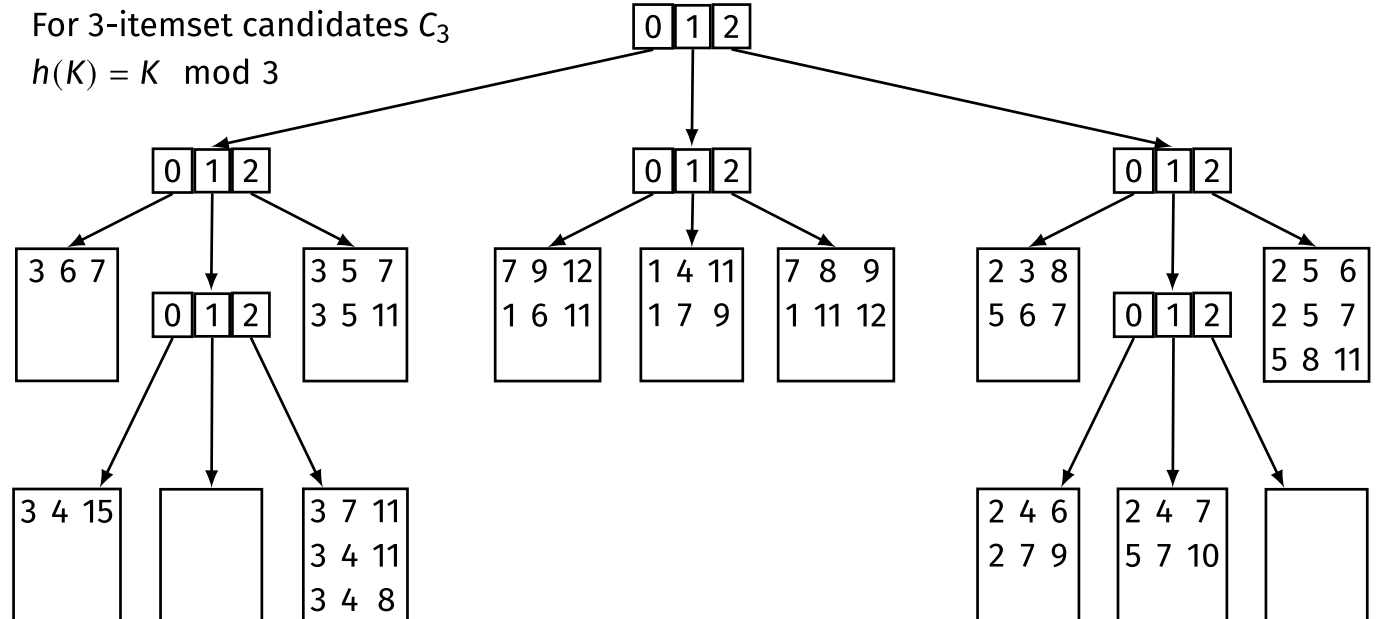
- ▶ many candidate itemsets (millions)
- ▶ a transaction may contain many items / itemsets

Idea: use a hashtree to store/index/manage C_k (and use it to process each transaction)

- ▶ **leaf nodes** contain list of itemsets (including their counts, initially 0)
- ▶ **inner nodes** hold hash tables;
 - each bucket at level d refers to a child node at level $d + 1$
- ▶ **root node** is at level 1
- ▶ hashing at level d is based on the d th smallest value in the k -itemset; at most $k + 1$ levels

Example





Hash Tree Index

Search for a single k -itemset

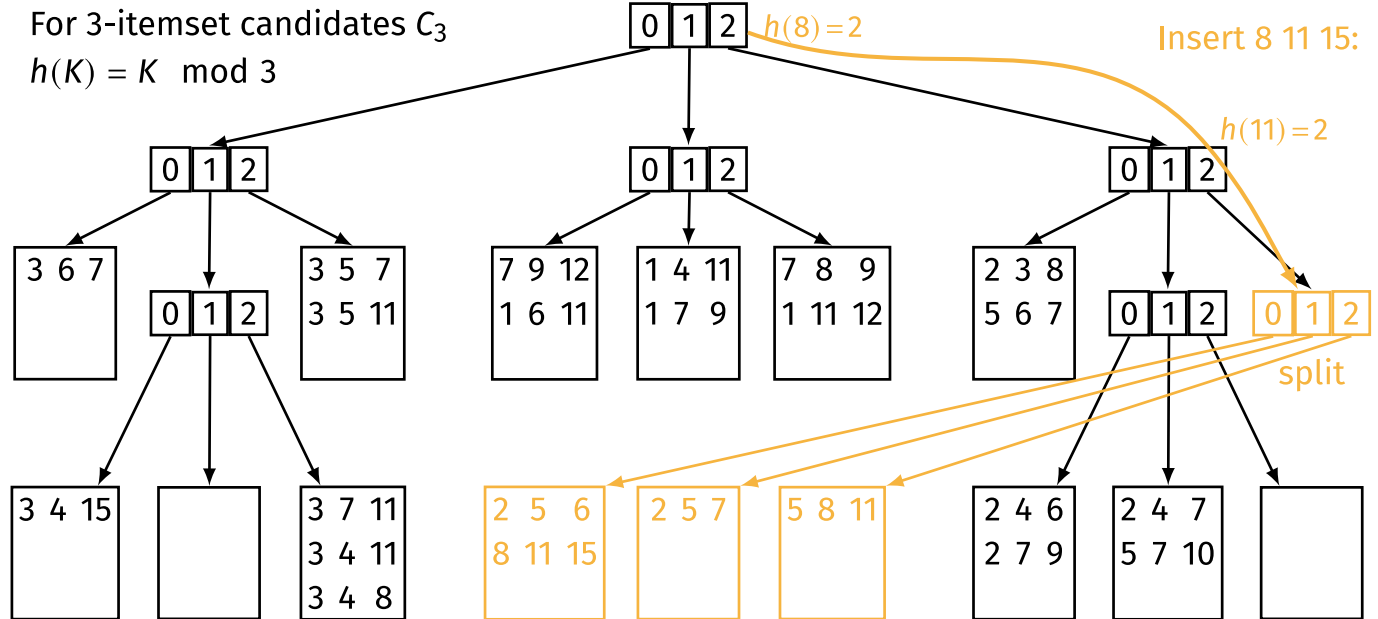
- ▶ begin at the root node
- ▶ in a leaf: search linearly (or, e.g., binary search)
- ▶ inner node at level d : apply the hash function h to d -th item of the itemset to determine the branch to follow

Insertion of a k -itemset

- ▶ search for corresponding leaf node and insert itemset into node
- ▶ in case of a node overflow:
 - ▶ transform leaf node into inner node
 - ▶ distribute the node's entries to new leaf nodes according to hash function

Example: Search and Insert





Determining Frequent Itemsets using Hash Tree

Search all candidates that are contained in transaction $T = (t_1 t_2 \dots t_m)$

at the root node:

- ▶ determine the hash value for *each* item in T
- ▶ continue search in each corresponding child node

at an inner node at level d (which has been reached by hashing t_i)

- ▶ determine hash values and continue search for each item t_k with $k > i$

at a leaf node

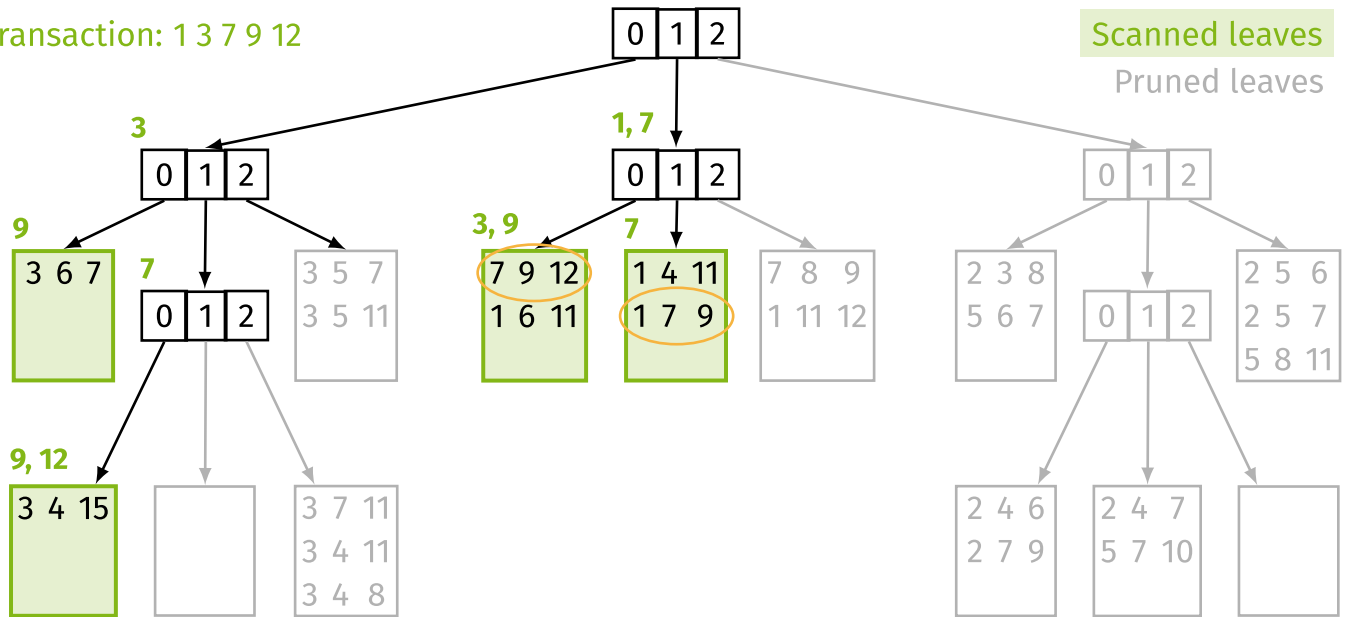
- ▶ test whether the itemsets in that node are contained in transaction T

Optimization: at every level d , only consider the hash codes of $(t_{i+1}, \dots, t_{m-k+d})$ as possible next-smallest item, where t_i was the item last used for splitting.

Example



Transaction: 1 3 7 9 12



Only 6 of 22 candidates from C_3 are considered for this transaction, 2 of which are correct.

« [NetFlix Example](#)

[Association Rules](#) »

Technische Universität Dortmund
 Informatik VIII
 AG Data Mining

Telefon: (+49) 231 755-6391
 Send Email

Otto-Hahn-Straße 14
 44227 Dortmund

Our Research

[Imprint](#) [Privacy](#) [Accessibility Statement](#) [Sitemap](#)

