

On the Complexity of Mining Itemsets from the Crowd Using Taxonomies

Antoine Amarilli^{1,2} Yael Amsterdamer¹ Tova Milo¹

¹Tel Aviv University, Tel Aviv, Israel

²École normale supérieure, Paris, France



Data mining

Data mining – discovering interesting patterns in large **databases**

Database – a (multi)set of **transactions**

Transaction – a set of **items** (aka. an **itemset**)

A simple kind of pattern to identify are **frequent itemsets**.

$$D = \left\{ \begin{array}{l} \{\text{beer, diapers}\}, \\ \{\text{beer, bread, butter}\}, \\ \{\text{beer, bread, diapers}\}, \\ \{\text{salad, tomato}\} \end{array} \right\}$$

- An itemset is **frequent** if it occurs in at least $\Theta = 50\%$ of transactions.
- $\{\text{salad}\}$ is not frequent.
- $\{\text{beer, diapers}\}$ is frequent. Thus, $\{\text{beer}\}$ is also frequent.

Human knowledge mining

- What if the database **doesn't really exist**?

Things to do in Athens:

$$D = \left\{ \begin{array}{l} \{\text{icdt, monday, laptop}\}, \\ \{\text{acropolis, sunglasses}\}, \\ \dots \\ \end{array} \right\}$$

Traditional medicine:

$$D = \left\{ \begin{array}{l} \{\text{hangover, coffee}\}, \\ \{\text{cough, honey}\}, \\ \dots \\ \end{array} \right\}$$

This data only exists in the **minds** of people!

Harvesting this data

- We **cannot** collect such data in a centralized database:
 - ① It's **impractical** to ask all users to surrender their data.

“Everyone please tell us all that you did the last three months.”
 - ② People do not **remember** the information.

“What were you doing on August 23th, 2013?”
- However, people remember **summaries** that we could access.

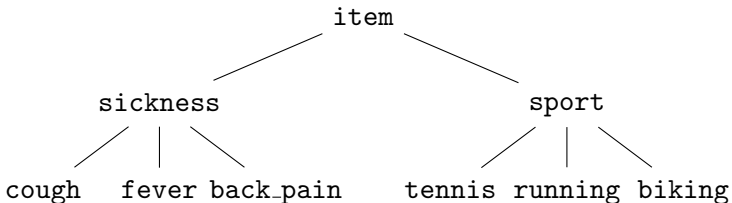
“Do you often play tennis on weekends?”
- We can just **ask** people if an itemset is frequent.

Crowdsourcing

- **Crowdsourcing** – solving hard problems through elementary queries to a crowd of users.
 - Find out if an itemset is **frequent** with the crowd:
 - ① **Draw** a sample of users from the crowd. (*black box*)
 - ② **Ask**: is this itemset frequent? (*“Do you often play tennis?”*)
 - ③ **Corroborate** the answers to eliminate bad answers. (*black box*)
 - ④ **Reward** the users. (*e.g., monetary incentive*)
- ⇒ An **oracle** that takes an itemset and finds out if it is frequent or not by asking crowd queries.

Taxonomies

Having a **taxonomy** over the items can save us work!



- If {sickness, sport} is **infrequent** then all itemsets such as {cough, biking} are **also infrequent**.
- Without the taxonomy, we need to test **all combinations!**
- Also avoids **redundant itemsets** like {sport, tennis}.

Cost

How to evaluate the **performance** of a strategy to identify the frequent itemsets?

Crowd complexity: The number of itemsets we ask about
(monetary cost, latency...)

Computational complexity: The complexity of computing the next question to ask

There is a **tradeoff** between the two:

- Asking **random** questions is computationally inexpensive but the crowd complexity is bad.
- Asking **clever** questions to obtain optimal crowd complexity is computationally expensive.

The problem

We can now describe the **problem**:

- We have:
 - A known **item domain** \mathcal{I} (set of items).
 - A known **taxonomy** Ψ on \mathcal{I} (is-a relation, partial order).
 - A crowd **oracle** freq to decide if an itemset is frequent or not.
 - Choose **interactively** questions based on past answers.
 - **Balance** crowd complexity and computational complexity.
- ⇒ Find out the status of **all** itemsets (learn freq exactly).

What is a good algorithm to solve this problem?

Table of contents

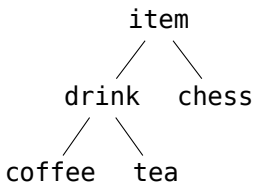
- 1 Background
- 2 Preliminaries**
- 3 Crowd complexity
- 4 Computational complexity
- 5 Conclusion

Itemset taxonomy

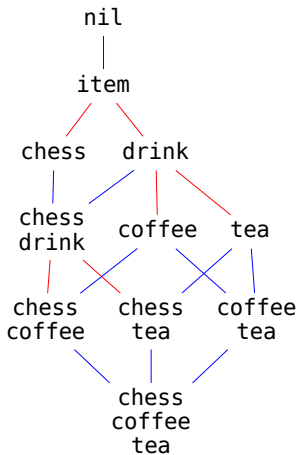
- **Itemsets** $I(\Psi)$ – the sets of **pairwise incomparable** items.
(e.g. {coffee, tennis} but not {coffee, drink})
- If an itemset is frequent then its **subsets** are also frequent.
- If an itemset is frequent then itemsets with **more general items** are also frequent.
- We define an **order relation** \leq on itemsets: $A \leq B$ for “A is more general than B”.
- **Formally**, $\forall i \in A, \exists j \in B$ s.t. i is more general than j .
- freq is **monotone**: if $A \leq B$ and B is frequent then A also is.

Itemset taxonomy example

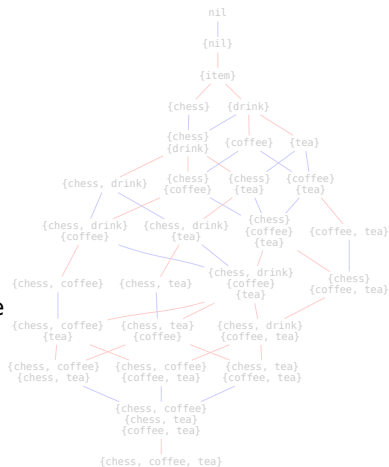
Taxonomy Ψ



Itemset taxonomy $I(\Psi)$

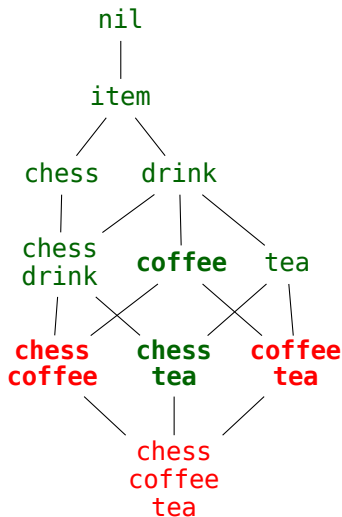


Solution taxonomy $S(\Psi)$



Maximal frequent itemsets

- **Maximal frequent itemset (MFI)**: a frequent itemset with no frequent descendants.
 - **Minimal infrequent itemset (MII)**.
 - The MFIs (or MIIs) **concisely** represent freq.
- ⇒ We can study complexity as a function of the size of the **output**.



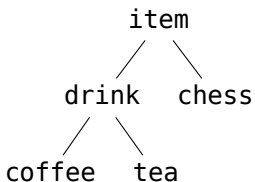
Solution taxonomy

- Conversely, (we can show) **any** set of pairwise incomparable itemsets is a possible MFI representation.
 - Hence, the set of all possible solutions has a **similar structure** to the “itemsets” over the itemset taxonomy $I(\Psi)$.
- ⇒ We call this the **solution taxonomy** $S(\Psi) = I(I(\Psi))$.

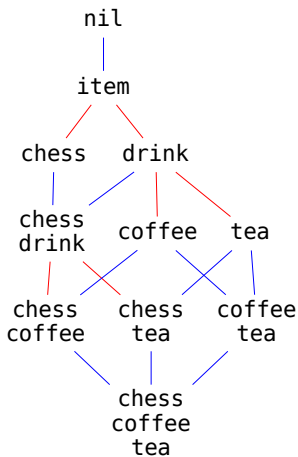
Identifying the freq predicate amounts to **finding the correct node** in $S(\Psi)$ through itemset frequency queries.

Solution taxonomy example

Taxonomy Ψ



Itemset taxonomy $I(\Psi)$



Solution taxonomy $S(\Psi)$

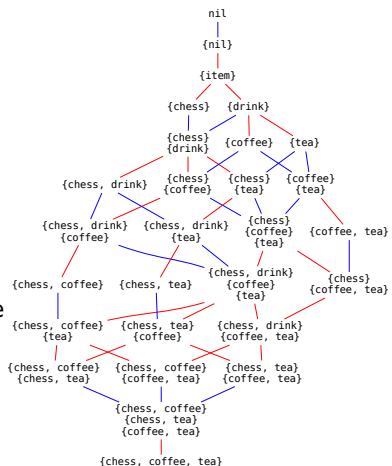


Table of contents

- 1 Background
- 2 Preliminaries
- 3 Crowd complexity**
- 4 Computational complexity
- 5 Conclusion

Lower bound

- Each query yields **one bit** of information.
- **Information-theoretic lower bound**: we need at least $\Omega(\log |S(\Psi)|)$ queries.
- This is bad in general, because $|S(\Psi)|$ can be **doubly exponential** in Ψ .
- As a function of the **original taxonomy** Ψ , we can write:
$$\Omega\left(2^{\text{width}[\Psi]} / \sqrt{\text{width}[\Psi]}\right).$$

Upper bound

- We can **achieve** the information-theoretic bound if there always is an unknown itemset that is frequent in about half of the possible solutions.
- A **result from order theory** shows that there is a constant $\delta_0 \approx 1/5$ such that some element always achieves a split of at least δ_0 .
- Hence, the previous bound is **tight**: we need $\Theta(\log |S(\Psi)|)$ queries.

nil	6/7
a1	5/7
a2	4/7
a3	3/7
a4	2/7
a5	1/7

Lower bound, MFI/MII

- To describe the solution, we need the MFIs **or** the MIIs.
- However, we need to query **both** the MFIs **and** the MIIs to identify the result uniquely: $\Omega(|\text{MFI}| + |\text{MII}|)$ queries.
- We can have $|\text{MFI}| = \Omega(2^{|\text{MII}|})$ and vice-versa.
- This bound is **not tight** (e.g., chain).

nil
|
a1
|
a2
|
a3
|
a4
|
a5

Upper bound, MFI/MII

- There is an **explicit algorithm** to find a new MFI or MII in $\leq |\mathcal{I}|$ queries.
- **Intuition**: starting with any frequent itemset, add items until you cannot add any more without becoming infrequent.
- The number of queries is thus $O(|\mathcal{I}| \cdot (|\text{MFI}| + |\text{MII}|))$.

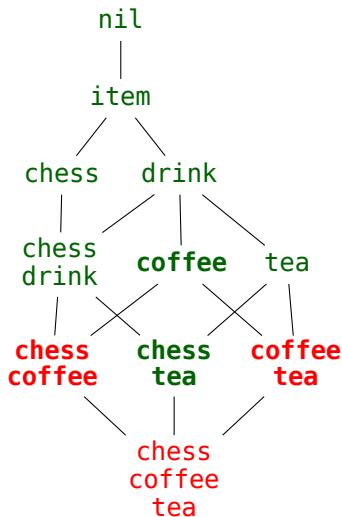


Table of contents

- 1 Background
- 2 Preliminaries
- 3 Crowd complexity
- 4 Computational complexity**
- 5 Conclusion

Hardness for standard (input) complexity

- We want an unknown itemset of $I(\Psi)$ that is frequent for **about half** of the possible solutions of $S(\Psi)$.
- We can **count** over $S(\Psi)$ but it may be exponential in $|I(\Psi)|$.
- Counting the **antichains** of $I(\Psi)$ is $FP^{\#P}$ -complete.
- Finding the best-split element in $I(\Psi)$ is **$FP^{\#P}$ -hard** in $|I(\Psi)|$?
- **Problem:** $I(\Psi)$ is not a general DAG, so we only show hardness in $|\Psi|$ for restricted (fixed-size) itemsets.
- **Intuition:** count antichains by comparing to a known poset; use a best-split oracle to compare; perform a binary search.

Hardness for output complexity

- In the incremental algorithm, **materializing** $I(\Psi)$ is expensive. Do we need to?
- Actually, how to decide if we can **stop** with our MFIs and MIIIs?
- Proved **EQ-hardness** for problem EQ (exact complexity open).

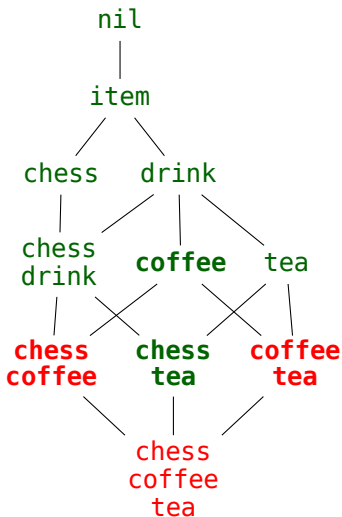


Table of contents

- 1 Background
- 2 Preliminaries
- 3 Crowd complexity
- 4 Computational complexity
- 5 Conclusion**

Summary and further work

We have studied the crowd and computational complexity of **crowd mining under a taxonomy**. What now?

- Improve the **bounds** and close gaps.
- Benchmark **heuristics** (chain partitioning, random, etc.).
- Integrate **prior knowledge**.
- Manage **uncertainty** (black box for now).
- Guide exploration with a **query** (under review).
- Work with **numerical values** for support.
- Mine more expressive **patterns**.
- Focus on **top- k** itemsets (work in progress).

Summary and further work

We have studied the crowd and computational complexity of **crowd mining under a taxonomy**. What now?

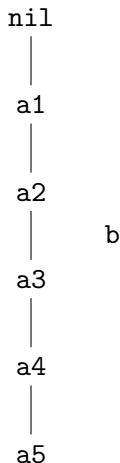
- Improve the **bounds** and close gaps.
- Benchmark **heuristics** (chain partitioning, random, etc.).
- Integrate **prior knowledge**.
- Manage **uncertainty** (black box for now).
- Guide exploration with a **query** (under review).
- Work with **numerical values** for support.
- Mine more expressive **patterns**.
- Focus on **top- k** itemsets (work in progress).

Thanks for your attention!

Greedy algorithms

- Querying an element of the chain may remove $< 1/2$ possible solutions.
 - Querying the isolated element b will remove **exactly** $1/2$ solution.
 - However, querying b classifies **far less** itemsets.
- ⇒ Classifying **many itemsets** isn't the same as eliminating **many solutions**.

Finding the **greedy-best-split** item is $\text{FP}^{\#P}$ -hard.



Restricted itemsets

- Asking about **large** itemsets is irrelevant.

“Do you often go cycling and running while drinking coffee and having lunch with orange juice on alternate Wednesdays?”

- If the itemset size is bounded by a **constant**, $I(\Psi)$ is tractable.
- ⇒ The crowd complexity $\Theta(\log |S(\Psi)|)$ is **tractable** too.

Chain partitioning

- Optimal strategy for **chain taxonomies**: binary search.
- We can determine a **chain decomposition** of the itemset taxonomy and perform binary searches on the chains.
- **Optimal** crowd complexity for a chain, performance in general is unclear.
- Computational complexity is **polynomial** in the size of $I(\Psi)$ (which is still exponential in Ψ).

```
nil
|
a1
|
a2
|
a3
|
a4
|
a5
```