

Mining Frequent Itemsets through Sampling with Rademacher Averages

Matteo Riondato (Two Sigma Investments, matteo@twosigma.com) and Eli Upfal (Department of Computer Science, Brown University, eli@cs.brown.edu)

1. What problem do we study?

Frequent Itemsets Mining: classic problem with many applications

Dataset \mathcal{D} Each line is a **transaction** of **items** from a set \mathcal{I}
 bread, milk
 bread, beer, eggs
 milk, beer, coke
 bread, milk, beer
 bread, milk, coke
 An **itemset** is a subset of \mathcal{I}
 The **frequency** $f_{\mathcal{D}}(\mathbf{A})$ of itemset \mathbf{A} is the **fraction** of transactions containing \mathbf{A}
 E.g.: $f_{\mathcal{D}}(\{\text{bread, milk}\}) = 3/5 = 0.6$

Problem (Frequent Itemsets (FI) Mining):

Given $\theta \in [0, 1]$ find (i.e., mine) all itemsets $\mathbf{A} \subseteq \mathcal{I}$ with $f_{\mathcal{D}}(\mathbf{A}) \geq \theta$
 i.e., compute the set $\mathbf{FI}(\mathcal{D}, \theta) = \{\mathbf{A} \subseteq \mathcal{I} : f_{\mathcal{D}}(\mathbf{A}) \geq \theta\}$

There are **exact algorithms** for FI mining (Apriori, FP-Growth, ...)

2. How to make FI Mining faster?

Exact algorithms do not scale with $|\mathcal{D}|$:

They scan \mathcal{D} multiple times: slow on large \mathcal{D}

How to get faster? **Trade-off accuracy for speed:**

Only mine **random samples** of \mathcal{D} that fit in main memory

Results are **approximate** but fast to compute

Approximation is OK: FI mining is a **exploratory task**

Key question:

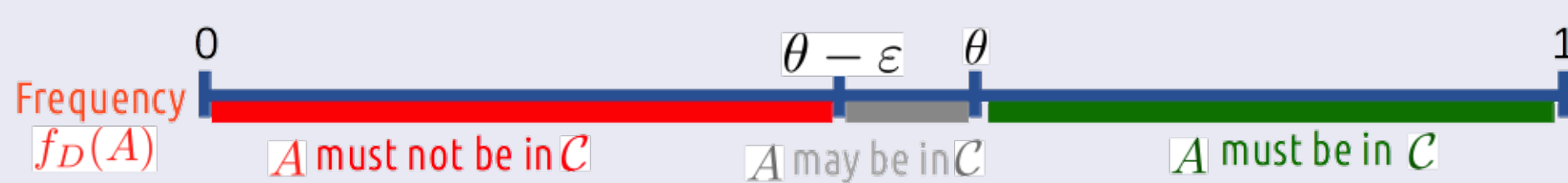
How much to sample to get an approximation of the desired quality?

Our contribution:

A fast approximation algorithm for FI mining that uses **Progressive Random Sampling** and **Rademacher Averages** and has **probabilistic guarantees** on the accuracy of the output

3. How to define an approximation of the FIs?

For $\epsilon, \delta \in (0, 1)$, a (ϵ, δ) -approximation to $\mathbf{FI}(\mathcal{D}, \theta)$ is a collection \mathcal{C} of itemsets s.t., with prob. $\geq 1 - \delta$:



"Close" False Positives are allowed, but **no False Negatives**

\mathcal{C} can act as **set of candidate FIs** to prune with fast scan of \mathcal{D}

Our approximation algorithm for FI mining returns an (ϵ, δ) -approximation to $\mathbf{FI}(\mathcal{D}, \theta)$ by using **Progressive Random Sampling**

4. What is Progressive Random Sampling (PRS)?

How much to sample from \mathcal{D} to obtain an (ϵ, δ) -approximation?

Prev. works: loose sample size for worst-case dataset

Instead, let's start sampling: the data will tell us when to stop

Progressive Random Sampling is an **iterative sampling scheme**

Outline of PRS algorithm for approximating $\mathbf{FI}(\mathcal{D}, \theta)$

At each iteration,

- create sample \mathcal{S} by drawing transactions from \mathcal{D} uniformly and independently at random
- Check a **stopping condition** on \mathcal{S} , to see if can get (ϵ, δ) -approximation from it
- If stopping condition is satisfied, **mine $\mathbf{FI}(\mathcal{S}, \gamma)$** for some $\gamma < \theta$ and **output it**
- Else, iterate with a **larger sample**

5. What are the challenges? What is our contribution?

The challenges are:

- Developing a **stopping condition** that
 - can be checked without expensive mining of each sample
 - guarantees that the output is a (ϵ, δ) -approximation
 - can be satisfied at small sample sizes
- Devising a method to choose the **next sample size**

Our contribution: We present the first algorithm that

- uses a **stopping condition** that does not mine each sample
- uses PRS to obtain an (ϵ, δ) -approximation of $\mathbf{FI}(\mathcal{D}, \theta)$
- computes the **optimal next sample size on the fly**

Previous contributions gave **no guarantees** and/or required **mining FIs from each sample** (too expensive). They used **predefined sample sizes**

6. What do we really need?

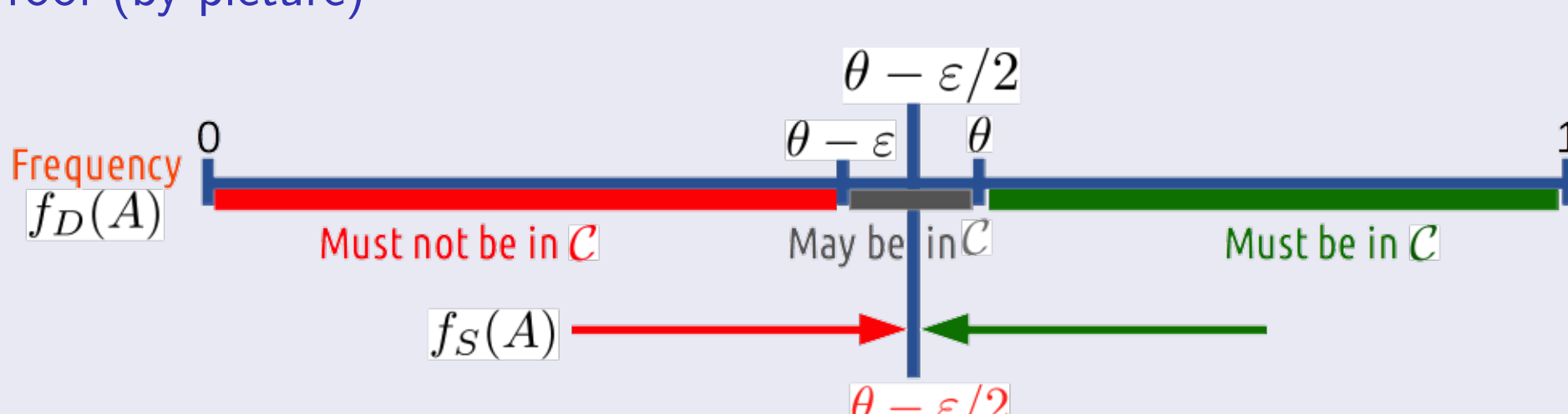
We need an **efficient procedure** that, given a sample \mathcal{S} of \mathcal{D} , computes a value η s.t.

$$\Pr\left(\sup_{\mathbf{A} \subseteq \mathcal{I}} |f_{\mathcal{D}}(\mathbf{A}) - f_{\mathcal{S}}(\mathbf{A})| \leq \eta\right) \geq 1 - \delta$$

Then the **stopping condition** just tests if $\eta \leq \epsilon/2$

Theorem: If $\eta \leq \epsilon/2$, then $\mathbf{FI}(\mathcal{S}, \theta - \epsilon/2)$ is an (ϵ, δ) -approximation to $\mathbf{FI}(\mathcal{D}, \mathcal{I}, \theta)$

Proof (by picture)



How to compute η ? Study statistics, learn about **Rademacher Averages!**

6. What are Rademacher Averages?

The behavior of $|\sup_{\mathbf{A} \subseteq \mathcal{I}} |f_{\mathcal{D}}(\mathbf{A}) - f_{\mathcal{S}}(\mathbf{A})||$ has been extensively studied using VC-dimension, covering numbers, and **Rademacher averages**

For any itemset \mathbf{A} and any transaction $\tau \in \mathcal{D}$, let

$$\phi_{\mathbf{A}}(\tau) = \begin{cases} 1 & \text{if } \mathbf{A} \subseteq \tau \\ 0 & \text{otherwise} \end{cases}$$

Let $\mathcal{S} = \{\tau_1, \dots, \tau_n\}$ be a random sample of \mathcal{D} , and let $\sigma_1, \dots, \sigma_n$ be independent Rademacher r.v. (1 with prob. 1/2, -1 othw.)

The **Rademacher Average** on \mathcal{S} is:

$$\mathbf{R}(\mathcal{S}) = \mathbb{E}_{\sigma} \left[\sup_{\mathbf{A} \subseteq \mathcal{I}} \frac{1}{n} \sum_{i=1}^n \sigma_i \phi_{\mathbf{A}}(\tau_i) \right]$$

The expectation is taken wrt the σ_i , i.e., conditionally on \mathcal{S}

7. How to use the Rademacher average?

Theorem (Key result from Statistical Learning Theory): Let

$$\eta = 2\mathbf{R}(\mathcal{S}) + \sqrt{\frac{2 \ln(2/\delta)}{n}}$$

Then

$$\Pr(\sup_{\mathbf{A} \subseteq \mathcal{I}} |f_{\mathcal{D}}(\mathbf{A}) - f_{\mathcal{S}}(\mathbf{A})| \leq \eta) \geq 1 - \delta$$

η is a **sample-dependent upper bound** to $\sup_{\mathbf{A} \subseteq \mathcal{I}} |f_{\mathcal{D}}(\mathbf{A}) - f_{\mathcal{S}}(\mathbf{A})|$

The key question is: **how do we compute or bound $\mathbf{R}(\mathcal{S})$?**

Computing $\mathbf{R}(\mathcal{S})$ efficiently allows us to compute η and check our stopping condition efficiently

8a. How can we bound the Rademacher average? (High-level)

We compute an **upper bound** to the distribution of the frequencies in \mathcal{S} of the **Closed Itemsets (CIs)** in \mathcal{S} (An itemset is **closed** iff none of its supersets has the same frequency)

Connection with the CIs:

$$\sup_{\mathbf{A} \subseteq \mathcal{I}} |f_{\mathcal{D}}(\mathbf{A}) - f_{\mathcal{S}}(\mathbf{A})| = \sup_{\mathbf{A} \in \text{CIs}} |f_{\mathcal{D}}(\mathbf{A}) - f_{\mathcal{S}}(\mathbf{A})|$$

Efficiency Constraint: use only information that can be obtained with a single scan of \mathcal{S}

How:

- We use the **frequency of the single items** and the **lengths of the transactions** to define a (conceptual) partitioning of the CIs into classes, and to **compute upper bounds to the size of each class and to the frequencies of the CIs in the class**
- We use these bounds to **compute an upper bound to $\mathbf{R}(\mathcal{S})$** by minimizing a convex function in \mathbb{R}^+ (no constraints)

8b. How to bound the Rademacher Averages? (In-depth)

For any itemset $\mathbf{A} \subseteq \mathcal{I}$, let $\mathbf{v}_{\mathcal{S}}(\mathbf{A})$ be the n -dimensional vector

$$\mathbf{v}_{\mathcal{S}}(\mathbf{A}) = (\phi_{\mathbf{A}}(\tau_1), \dots, \phi_{\mathbf{A}}(\tau_n)),$$

and let $\mathbf{V}_{\mathcal{S}} = \{\mathbf{v}_{\mathcal{S}}(\mathbf{A}), \mathbf{A} \subseteq \mathcal{I}\}$ ($\mathbf{V}_{\mathcal{S}}$ is a set)

Theorem (Variant of Massart's Lemma):

Let $\mathbf{w} : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ be the function

$$\mathbf{w}(\mathbf{s}) = \frac{1}{\mathbf{s}} \ln \sum_{\mathbf{v} \in \mathbf{V}_{\mathcal{S}}} \exp(\mathbf{s}^2 \|\mathbf{v}\|^2 / (2n^2))$$

Then

$$\mathbf{R}(\mathcal{S}) \leq \min_{\mathbf{s} \in \mathbb{R}^+} \mathbf{w}(\mathbf{s})$$

Since $\tilde{\mathbf{w}}$ is convex, its global minimum can be found efficiently

What does the set of vectors $\mathbf{V}_{\mathcal{S}}$ look like?

Let $\mathbf{CI}(\mathcal{S})$ be the set of all **Closed Itemsets** in \mathcal{S}

Lemma: $\mathbf{V}_{\mathcal{S}}$ contains all and only the vectors $\mathbf{v}_{\mathcal{S}}(\mathbf{A})$ for all $\mathbf{A} \in \mathbf{CI}(\mathcal{S})$.

Issue: Can not mine $\mathbf{CI}(\mathcal{S})$ to compute $\mathbf{w}(\mathbf{s})$: it is **too expensive!**

Solution: Define a function $\tilde{\mathbf{w}}(\mathbf{s})$ efficient to compute and minimize and s.t. $\tilde{\mathbf{w}}(\mathbf{s}) \geq \mathbf{w}(\mathbf{s})$ for all \mathbf{s} . Then use $\tilde{\mathbf{w}}(\mathbf{s})$ to compute $\eta_{\mathcal{S}}$

We define a **partitioning \mathcal{P}** of $\mathbf{CI}(\mathcal{S})$

- Assume an ordering $<_{\mathcal{I}}$ of \mathcal{I} . For any $\mathbf{a} \in \mathcal{I}$, assume an ordering $<_{\mathbf{a}}$ of the transactions of \mathcal{S} that contain \mathbf{a}
- For any $\mathbf{A} \in \mathbf{CI}(\mathcal{S})$, let $\mathbf{a} \in \mathbf{A}$ be the item in \mathbf{A} that comes first wrt $<_{\mathcal{I}}$, and let τ be the transaction containing \mathbf{A} that comes first wrt $<_{\mathbf{a}}$. We assign \mathbf{A} to the class $\mathcal{P}_{\mathbf{a}, \tau}$

For each class $\mathcal{P}_{\mathbf{a}, \tau}$ we

- compute an upper bound to $|\mathcal{P}_{\mathbf{a}, \tau}|$ using $|\tau|$ and $<_{\mathbf{a}}$
- use $f_{\mathcal{S}}(\mathbf{a})$ as upper bound to $f_{\mathcal{S}}(\mathbf{A})$, for $\mathbf{A} \in \mathcal{P}_{\mathbf{a}, \tau}$

Very efficient to compute $f_{\mathcal{S}}(\mathbf{a})$ while creating the sample

The new function $\tilde{\mathbf{w}}$ used to compute $\mathbf{R}(\mathcal{S})$ is:

$$\tilde{\mathbf{w}}(\mathbf{s}) = \frac{1}{\mathbf{s}} \ln \sum_{\mathbf{a} \in \mathcal{I}_{\mathcal{S}}} \left(\left(1 + \sum_{r=1}^{\chi_{\mathbf{a}}} \sum_{j=1}^{g_{\mathbf{a}, r}} 2^{\min\{r, h_{\mathbf{a}, r} - j\}} \right) e^{-\frac{\mathbf{s}^2 f_{\mathcal{S}}(\mathbf{a})}{2n}} \right)$$

Then

$$\eta_{\mathcal{S}} = \min_{\mathbf{s} \in \mathbb{R}^+} \tilde{\mathbf{w}}(\mathbf{s}) + \sqrt{\frac{2 \ln(2/\delta)}{n}}$$

9. How do we choose the next sample size?

We can compute the next sample size on the fly

First iteration: Use a sample of size at least $8 \ln(2/\delta) \epsilon^{-2}$

Why? It is impossible that $\eta \leq \epsilon/2$ at smaller sample sizes

Successive iterations:

multiply the sample size from the previous iteration by $(2\eta/\epsilon)^2$

Intuition: If the frequencies of the items in the current iteration and the distribution of the transaction lengths are the same as in the previous iteration, then the stopping condition will be satisfied at this iteration

12. Experimental Evaluation

Implementation and Datasets: Sampling and stopping condition implemented in C++11, FIs mining using existing C implementation, NLOpt for minimization step. Implementation available from <http://cs.brown.edu/~matteo/radeprogrfi.tar.bz2>. Datasets from the FIMI'03 repository. Experiments run on a machine with a quad-core AMD PhenomTM II X4 955 processor and 16GB of RAM, running GNU/Linux 3.2.0

Recall: in 10K+ runs, **always** returned an ϵ -approx., not just with prob. $1 - \delta$. All itemsets from $\mathbf{FI}(\mathcal{D}, \theta)$ are in output: the **recall is always 100%**

Precision: Algorithm gives no guarantee. Varied between 15% and 92%, depending on the parameters and on the dataset. Price to pay when mining a subset of transactions. **The output can be used as a set of candidates** from which to compute efficiently the exact collection of FIs with a single linear scan of the dataset (negligible cost)

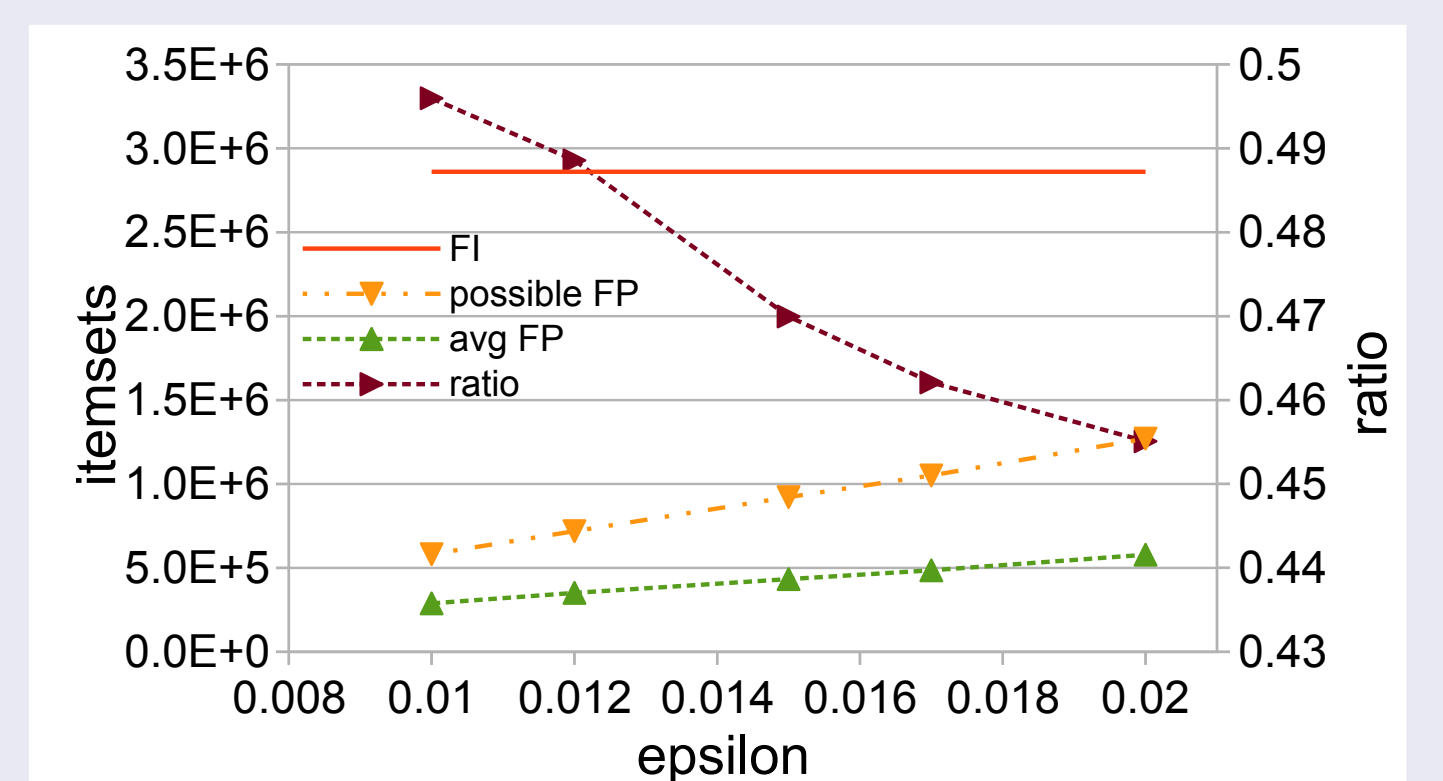


Figure: Precision for connect, $\theta = 0.72$

Frequency Estimation: For each itemset \mathbf{A} in output, we measure the absolute frequency error. **Average and maximum abs. error are almost 10x less than $\epsilon/2$** and very concentrated. We also measure the **relative frequency error** (right axis): it is **always less than 1.4%**

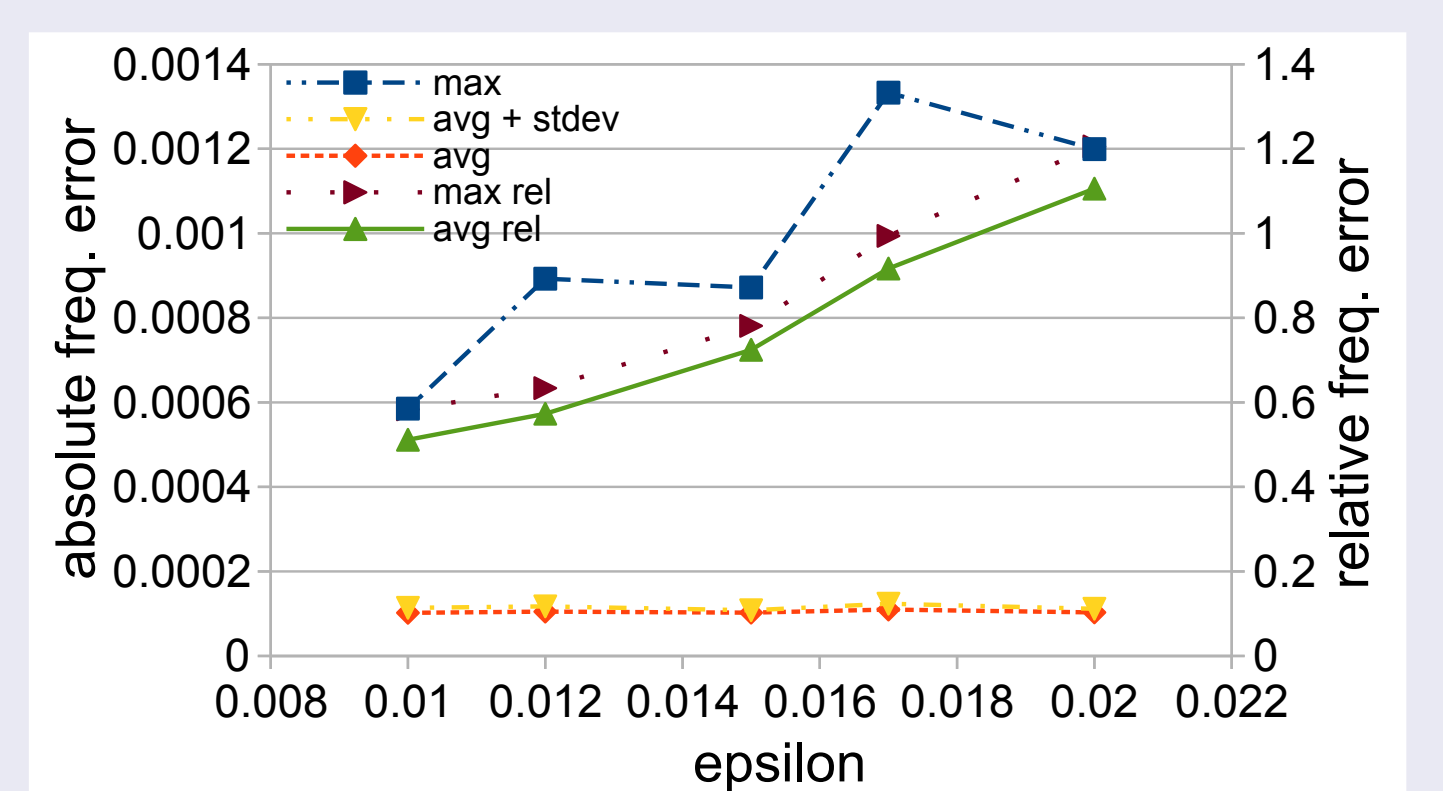


Figure: Frequency error for retail, $\theta = 0.015$

Runtime: We compare the running time of our algorithm to that of a one-shot sampling approach (VC), to that of FP-Growth, and to the running time of our algorithm using a geometric sampling schedule $|\mathcal{S}_i| = \alpha^i |\mathcal{S}_1|$ for $\alpha \in \{2, 2.5, 3\}$ (to evaluate the automatic sample schedule). Our algorithm (avg line) vastly outperforms the exact algorithm and VC. The automatic sampling schedule is more efficient than the geometric sample schedule by avoiding the creation and analysis of samples whose size is probably not sufficient for the stopping condition to be satisfied, based on information obtained from the current sample

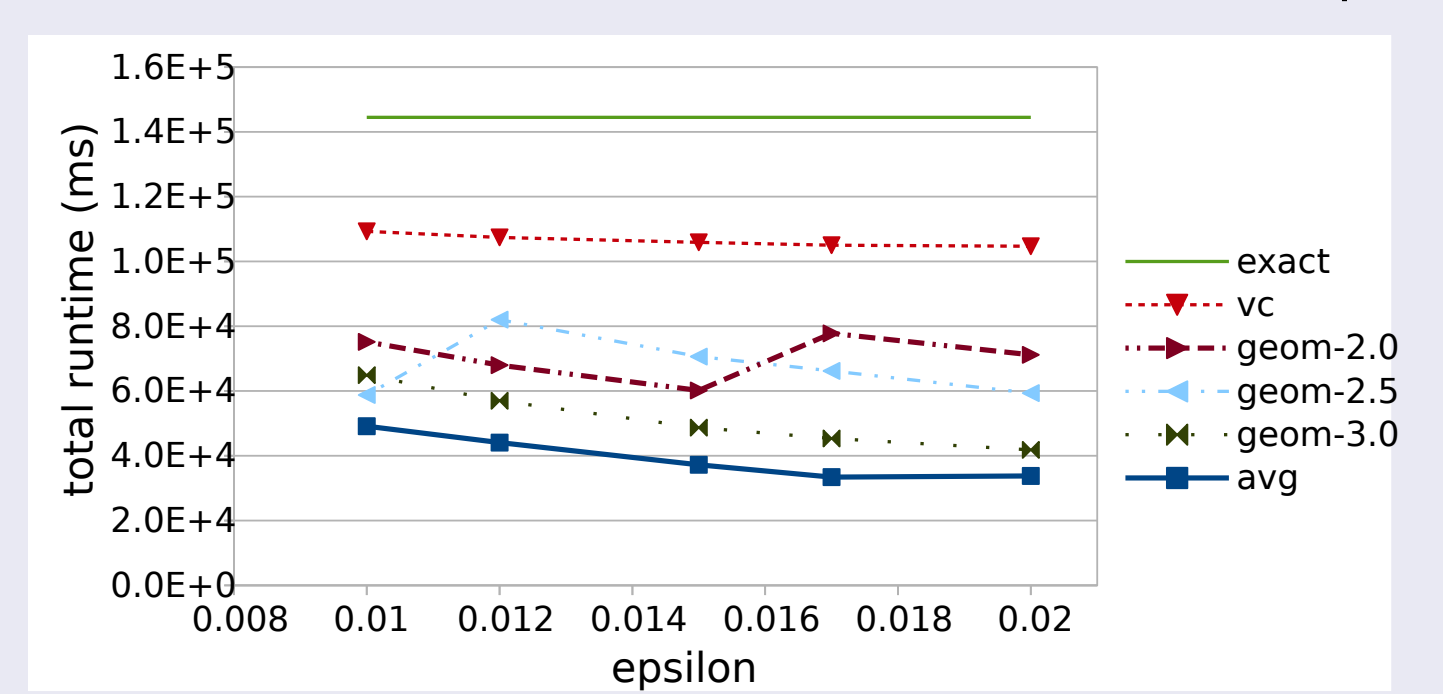


Figure: Running time for BMS-POS, $\theta = 0.015$

We also analyze the **breakdown of the runtime** of our algorithm, splitting it between time needed to sample the transaction, time needed to evaluate the stopping condition, and time needed to perform the mining of the sample after the stopping condition is satisfied

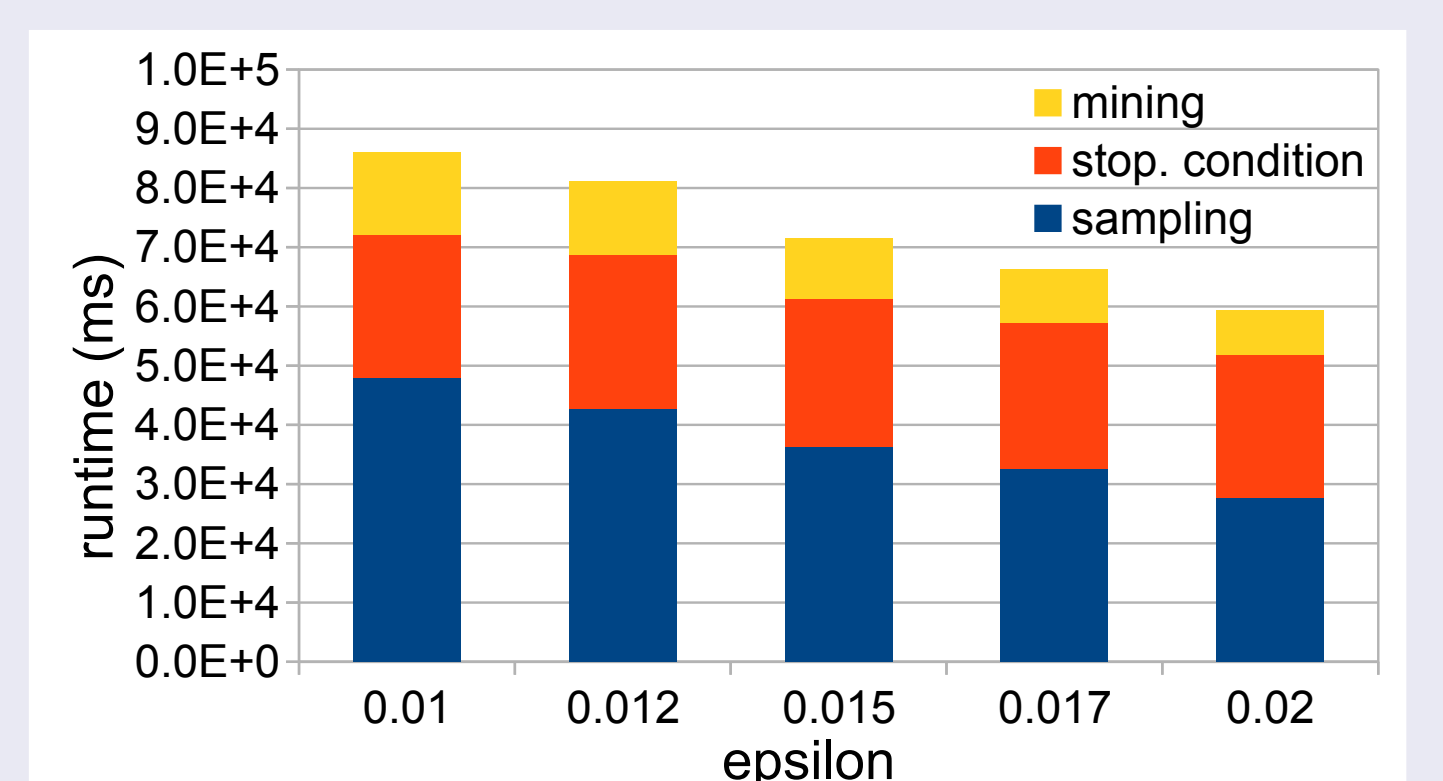


Figure: Breakdown of runtime for pumsb_star, $\theta = 0.32$.

11. Future Work

Improve bound to the Rademacher Average by studying the orderings and the function $\tilde{\mathbf{w}}$

Improve the automatic sample schedule by better taking into account the behavior of $\tilde{\mathbf{w}}$

Provide multiplicative error guarantees

Extend to other measures of interestingness not bounded by anti-monotonic functions

Acknowledgements

This work was supported, in part, by NSF grant IIS-1247581 and NIH grant R01-CA180776

Extended Version

Available from <http://goo.gl/a7dceX>