

DISTINCT VALUE ESTIMATORS FOR ZIPFIAN DISTRIBUTIONS

SERGEI VASSILVITSKII
RAJEEV MOTWANI

STANFORD UNIVERSITY

PROBLEM STATEMENT

Given a large multiset X with n elements, count the number of distinct elements in X .

$$X = \{a, b, c, a, a, c, b, a\} \implies \text{Distinct}(X) = 3$$

Alternatively, given samples from a distribution \mathcal{P} , estimate the 0-th frequency moment.

WHY DO WE CARE?

Good Planning for SQL Queries. Consider:

`select * from R, S where R.A = S.B and $f(S.C) > k$`

where f is expensive to compute.

If $S.C$ has few distinct elements, compute f first, cache results, then join.

If $S.C$ has many elements, compute the join first, then check the f condition.

Orders of Magnitude Improvements

CLASSICAL PROBLEM

Different approaches:

Streaming Input - Minimize space used.

Sample from Input - Guarantee on approximations?

Given a sample of size r from X , find \hat{D} an approximation to $Distinct(X)$.

PREVIOUS WORK

Good-Turing Estimator: “The Population Frequencies of Species, and the estimation of Population Parameters,” 1953.

Other Heuristic Estimators:

Smoothed Jackknife Estimator (Haas et. al)

Adaptive Estimator (Charikar et. al)

Many Others

PREVIOUS WORK - THEORY

Given r samples from a set of size n

Guaranteed Error Estimator(GEE) [CCMN]

Approximation Ratio: $O(\sqrt{n/r})$

Lower Bound:

There exist inputs such that with constant probability any estimator will have approximation ratio at least:

$$\sqrt{\frac{n-r}{2r}}$$

LOWER BOUND DETAIL

Scenario 1: $S = \{x, x, \dots, x\}, |S| = n$

Scenario 2: $S' = \{x, x, \dots, x, y_1, \dots, y_k\}, |S'| = n$

With $k = \frac{n-r}{2r} \ln \frac{1}{\delta}$, after r samples cannot distinguish between the two scenarios with probability at least δ .

SO WHY ARE WE HERE?

Many large datasets are not worst-case. In fact, many follow Zipfian Distributions.

$$\text{Zipf}_\theta(i) \propto \frac{1}{i^\theta}$$

Examples:

- In/Out-Degrees of the Web Graph
- Word frequencies in many languages
- many, many more.

PROBLEM DEFINITION

Suppose $X \sim \text{Zipf}_\theta$ on D elements.

θ is known, D is unknown

Estimate D by sampling from X .

Two Kinds of Results:

- Adaptive Sampling: Will sample from X until a stopping condition is met.
- Best-you-can Estimation: Given a sample from X , return best estimate of D .

RESULTS

Let p^* be the probability of the least likely element.

Adaptive sampling will return D after at most $O\left(\frac{\log D}{p^*}\right)$ samples with constant probability.

Given $r = \frac{(1 + 2\epsilon)^{1+\theta}}{p^*}$ samples, can return an $1 + \epsilon$ estimate to D with probability at least $1 - \exp(-\Omega(D\epsilon^2))$

OUTLINE

Introduction

Techniques

Experimental Results

Conclusion

APPROXIMATION TECHNIQUES

For a sample of size r let f_r be the number of distinct values in the sample.

Suppose D and θ are known, then we can compute $E_{D,\theta}[f_r]$ the expected number of distinct values in the sample.

If f_r^* is the number of distinct values observed, the estimator returns \hat{D} such that $E_{\hat{D},\theta}[f_r] = f_r^*$.

ANALYSIS

LEMMA: Tight Distribution of f_r .

For large enough r ,

$$\Pr \left[|E[f_r] - f_r| \geq \epsilon E[f_r] \right] \leq \exp(-\epsilon^2 \Omega(D))$$

Proof: Parallels the sharp threshold coupon collector arguments for uniform distributions.

ANALYSIS (2)

LEMMA: MLE preserves approximation

Given: $f_r \leq (1 + \epsilon)E_{D,\theta}[f_r]$, observed f_r^* elements

Let \hat{D} such that $f_r^* = E_{\hat{D},\theta}[f_r]$, and $r \geq 1/p^*$.

Then: $(1 - 2\epsilon)\hat{D} \leq D \leq (1 + 2\epsilon)\hat{D}$

OUTLINE

Introduction

Techniques

Experimental Results

Conclusion

THE COMPETITION

Zipfian Estimator (ZE): Performance guarantees only for Zipfian Distributions.

Guaranteed Error Estimator (GEE): $O(\sqrt{n/r})$ error guarantee. (Works for all distributions)

Analytic Estimator (AE): Best performing heuristic - no theoretical guarantees.

DATASETS

Synthetic Data:

- Vary number of distinct elements $D \in \{10k, 50k, 100k\}$
- Vary the Database size $n \in \{100k, 500k, 1000k\}$
- Vary the skew of the distribution $\theta \in \{0, 0.5, 1\}$

Real Datasets

- “Router” dataset - Packet trace from the Internet Traffic Archive. $\theta \approx 1.6, n \approx 4M, D \approx 250k$

ESTIMATING θ

Recall: $Zipf_{\theta}(i) \propto \frac{1}{i^{\theta}}$

Let f_i be the frequency of the i -th element.

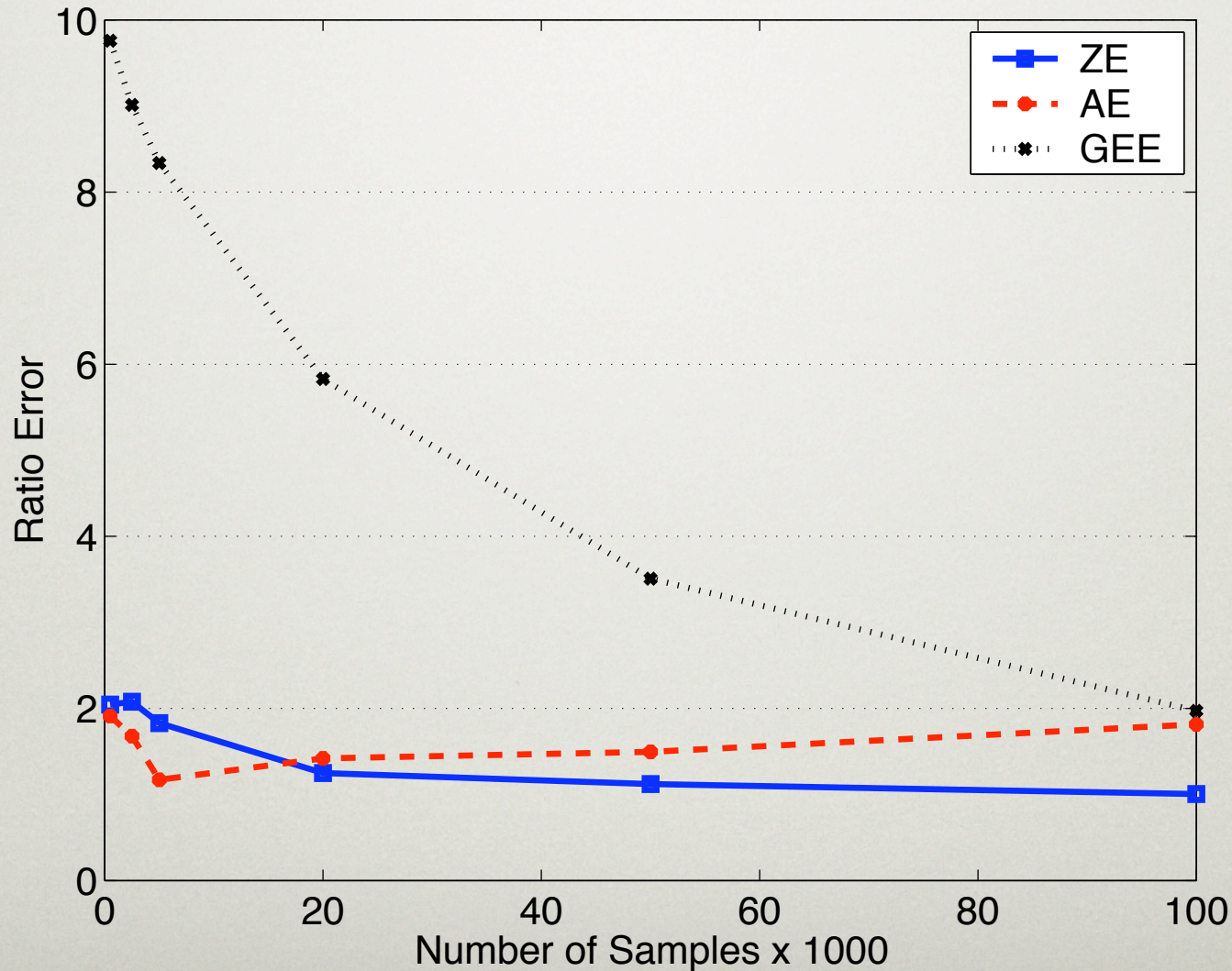
$$E[f_i] = cr i^{-\theta} \implies \log E[f_i] = \log cr - \theta \log i$$

Estimate θ by doing linear regression on

$\log f_i$ vs $\log i$ plot.

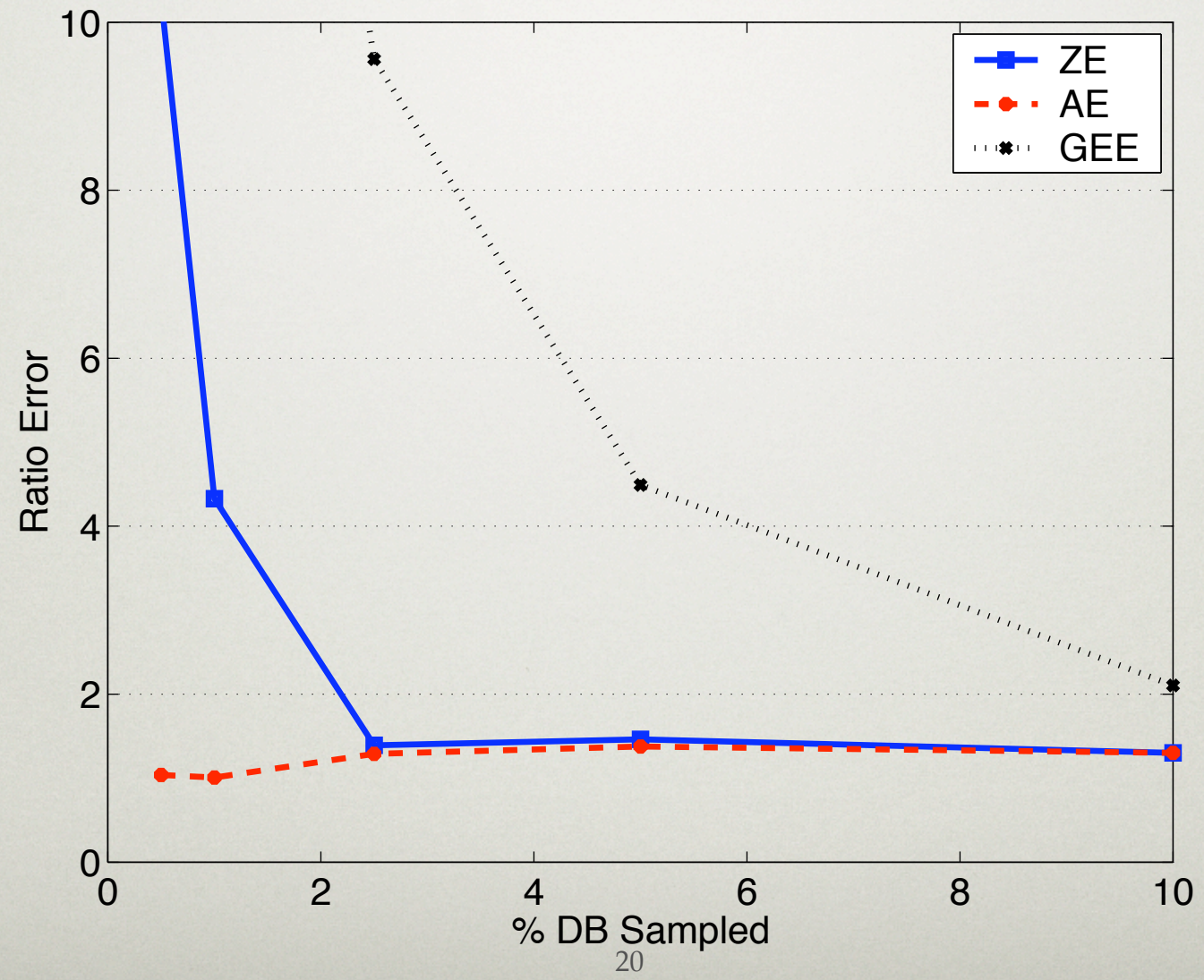
EXPERIMENTAL RESULTS

Theta = 0.5, D = 50000, n = 1M



EXPERIMENTAL RESULTS (2)

Router Dataset



OUTLINE

Introduction

Techniques

Experimental Results

Conclusion

CONCLUSION

Can have error guarantees if the family of distributions is known ahead of time.

How does the approximation of θ affect error guarantees?

Subtle problem: disk reads occur in blocks. Time to sample 10% is equivalent to reading the whole DB.

THANK YOU