# STRATIFIED SAMPLING MEETS MACHINE LEARNING

**KEVIN LANG
KONSTANTIN SHMAKOV
EDO LIBERTY**

Search YDN

# Yahoo Mobile Developer Suite for your apps

Measure, monetize, advertise and improve your apps with Yahoo tools.
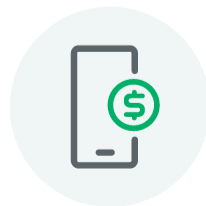
**Sign Up**          **Log in to Flurry**

## Flurry Analytics

Get free insights from the industry's leading mobile app analytics tool.
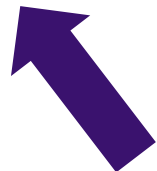
## Yahoo App Publishing

Monetize your app with native and video ads from Yahoo, Flurry, and BrightRoll advertisers.

## Yahoo App Marketing

Reach your target audience with the Gemini native and video marketplace.
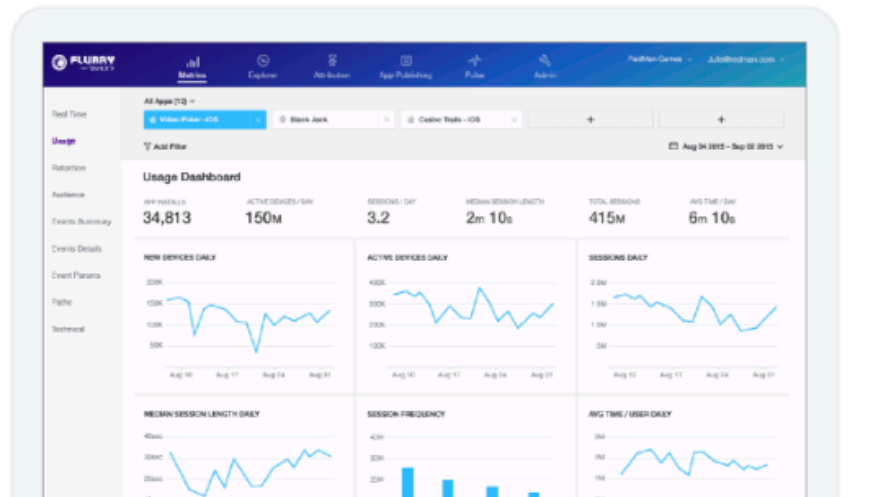
YAHOO!

# Introducing the all new Flurry Analytics

Measure and analyze activity across your app portfolio to answer your hardest questions and optimize your app experience.

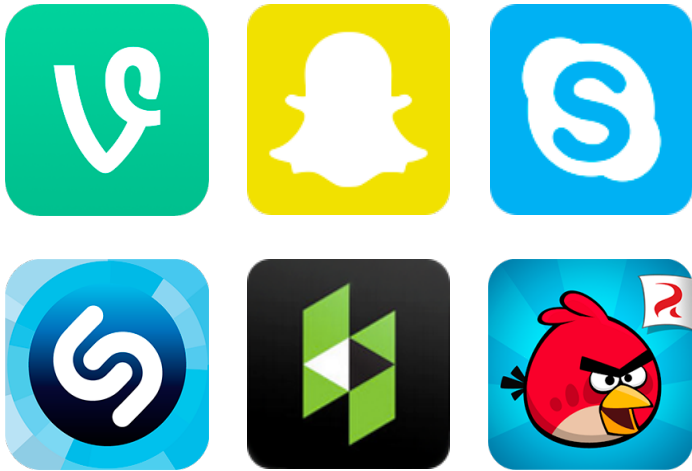| Sign Up | Log in to Flurry | Documentation |
|---|---|---|

## Flurry Grows with You

As your business grows, we are committed to supporting you at scale for free.

Empower all your employees to leverage analytics for data driven decision making.
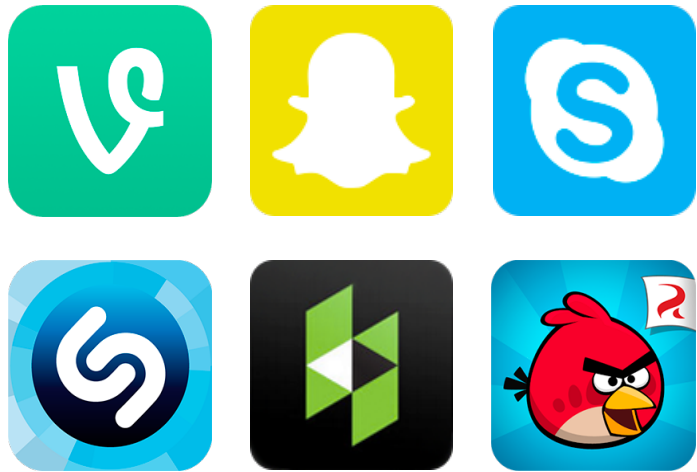
Product Features

YAHOO!

Flurry Analytics
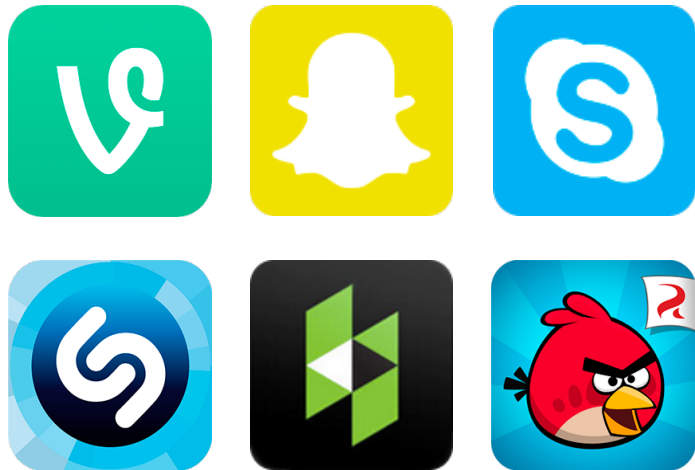
Apps with Flurry SDK

YAHOO!

# Flurry Analytics

Events

YAHOO!

# Flurry Analytics

Query $q(\cdot)$

Answer $\sum_i q(u_i)$

Examples:
- Number of event of a certain type
- Number of unique user
- Number of unique users in a specific day
- Total time spent in certain geo
- Average $ spent by age

YAHOO!

# SAMPLING

Challenges:
1. The data is very large. Computing $\sum_i q(u_i)$ exactly is too costly.
2. The function $q(\cdot)$ is user specified and completely unconstrained.

Good News:
And approximate answer is acceptable (if the error is small)

Solution:
Estimate the answer on a random subset of the records

YAHOO!

# NOTATIONS

- $q_i := q(u_i)$ for brevity

- $y := \sum_i q_i$ the exact answer for the query $q$

- $p_i$ the probability of choosing record $i$

- $S$ the set of sampled records, each chosen with probability $p_i$

- $\tilde{y} = \sum_{i \in S} q_i/p_i$ the Horvitz-Thompson estimator for $y$

YAHOO!

# PROPERTIES

- $\mathbb{E}[\tilde{y} - y] = 0$  Horvitz-Thompson **estimator is unbiased**

- $\sigma[\tilde{y} - y] \leq y\sqrt{1/(\zeta \cdot \mathrm{card}(q))}$  its standard deviation isn't large

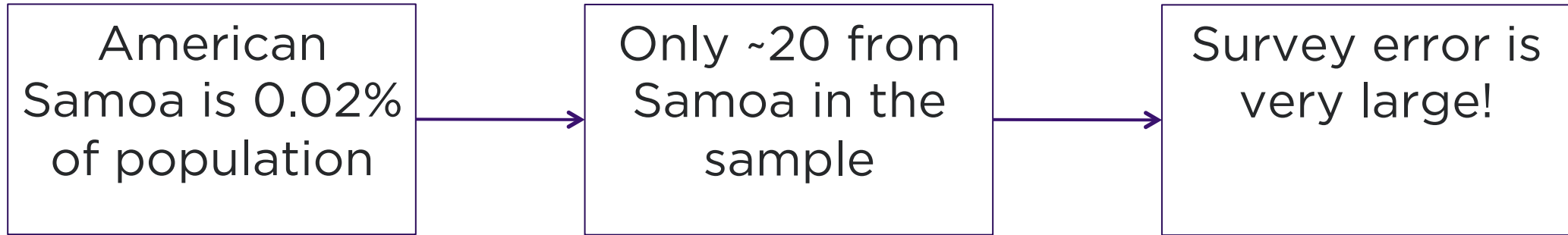$$\zeta = \min_i p_i \qquad \mathrm{card}(q) := \sum |q_i| / \max |q_i|$$

- $\Pr[|\tilde{y} - y| \geq \varepsilon y] \leq e^{-O(\varepsilon^2 \zeta \cdot \mathrm{card}(q))}$  probability for large error is small

$$\boxed{\mathrm{card}(q) \sim \Omega(n) \quad \rightarrow \quad |S| \sim 1/\varepsilon^2}$$

(Olken, Rotem and Hellerstein 1986, and 1990) application to databases
(Acharya, Gibbons, Poosala 2000) uniform sampling is best in the worst case

YAHOO!

# STRATIFIED SAMPLING

- Sample = 100,000 US individuals.
- Query = Republicans vs. Democrats in American Samoa?

| American Samoa is 0.02% of population | → | Only ~20 from Samoa in the sample | → | Survey error is very large! |

$$\text{If } \mathrm{card}(q) \text{ is small } |S| \text{ must be large}$$

- Sample <u>different *strata*</u> (e.g. US territories) with <u>different probabilities</u>.

(Neyman, Jerzy 1934)

YAHOO!

# DBLP EXAMPLE

Choosing the right strata is hard!

- 2,101,151 papers
- 1000 most populous venues
- Query example
  - title contains "learning" and # authors <= 3
  - title contains "mechanism" and year > 2004

What is the right stratification here?

- Stratifying by venue made things worse!
- Stratifying by year was better but still worse than uniform sampling.
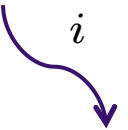
YAHOO!

# SAMPLING, STRATIFICATION, AND DATABASES

- Design strata that minimize worst case variance on possible queries
- Linearly combine strata based on record features
- Combine stratifies and uniform sampling: Congressional Sampling
  - Acharya, Gibbons, Poosala 2000:

Important idea: consider past queries to the database!

- Each stratum is a set of records that agree on all queries
  - Chaudhuri, Das and Narasayya 2007: optimize for the query log

- Split to two strata, per each query. Take linear combinations
  - Joshi, Jermaine, 2008: linear combinations of stratified probabilities

YAHOO!

# OUR APPROACH

- Assume queries are drawn from a distribution $\mathbb{Q}$

- Use the query log $Q$ as a "training set" (assumed w.r.t. $\mathbb{Q}$)

- Allow each record to be sampled with a different probability $p_i$

- Minimize the Risk $\mathbb{E}[(\tilde{y} - y)^2]$

- This translates to $\mathbb{E}_{q \sim \mathbb{Q}} \sum_i q_i^2 (1/p_i - 1)$

unknown

YAHOO!

# OUR APPROACH

- ERM: Minimize $\quad \displaystyle\sum_{q \in Q} \sum_i q_i^2 (1/p_i - 1)$

  Query log

- Sampling budget $\quad \displaystyle\sum_i p_i c_i \leq B \qquad ( \sum_i c_i \ll B )$

- Regularization $\quad \forall i \; p_i \in [\zeta, 1] \qquad ( \zeta \leq B/\sum_i c_i )$

YAHOO!

# OUR APPROACH

- Solve with Lagrange multipliers

$$
\max_{\alpha,\beta,\gamma}[\frac{1}{|Q|}\sum_{q\in Q}\sum_i q_i^2(1/p_i-1)-\sum_i \alpha_i(p_i-\zeta)
$$
$$
-\sum_i \beta_i(1-p_i)-\gamma(B-\sum_i p_i c_i)]
$$

- By KKT conditions

$$
p_i = \zeta \qquad \text{or} \qquad p_i \propto \sqrt{\frac{1}{c_i}\frac{1}{|Q|}\sum_{q\in Q}q_i^2} \qquad \text{or} \qquad p_i = 1
$$

YAHOO!

# OUR APPROACH

1: **input:** training queries $Q$,
2:          budget $B$, record costs $c$,
3:          regularization factor $\eta \in [0, 1]$
4: $\zeta = \eta \cdot (B / \sum_i c_i)$
5: $\forall i \;\; z_i = \sqrt{\frac{1}{c_i} \frac{1}{|Q|} \sum_{q \in Q} q_i^2}$
6: Binary search for $\lambda$ satisfying $\sum_i c_i \, \mathrm{CLIP}_\zeta^1(\lambda z_i) = B$
7: **output:** $\forall i \;\; p_i = \mathrm{CLIP}_\zeta^1(\lambda z_i)$

$$\mathrm{Risk}(p) \leq \mathrm{Risk}(p^*) \left( 1 + O \left( \mathrm{skew} \sqrt{\frac{\log(n/\delta)}{|Q|}} \right) \right)$$
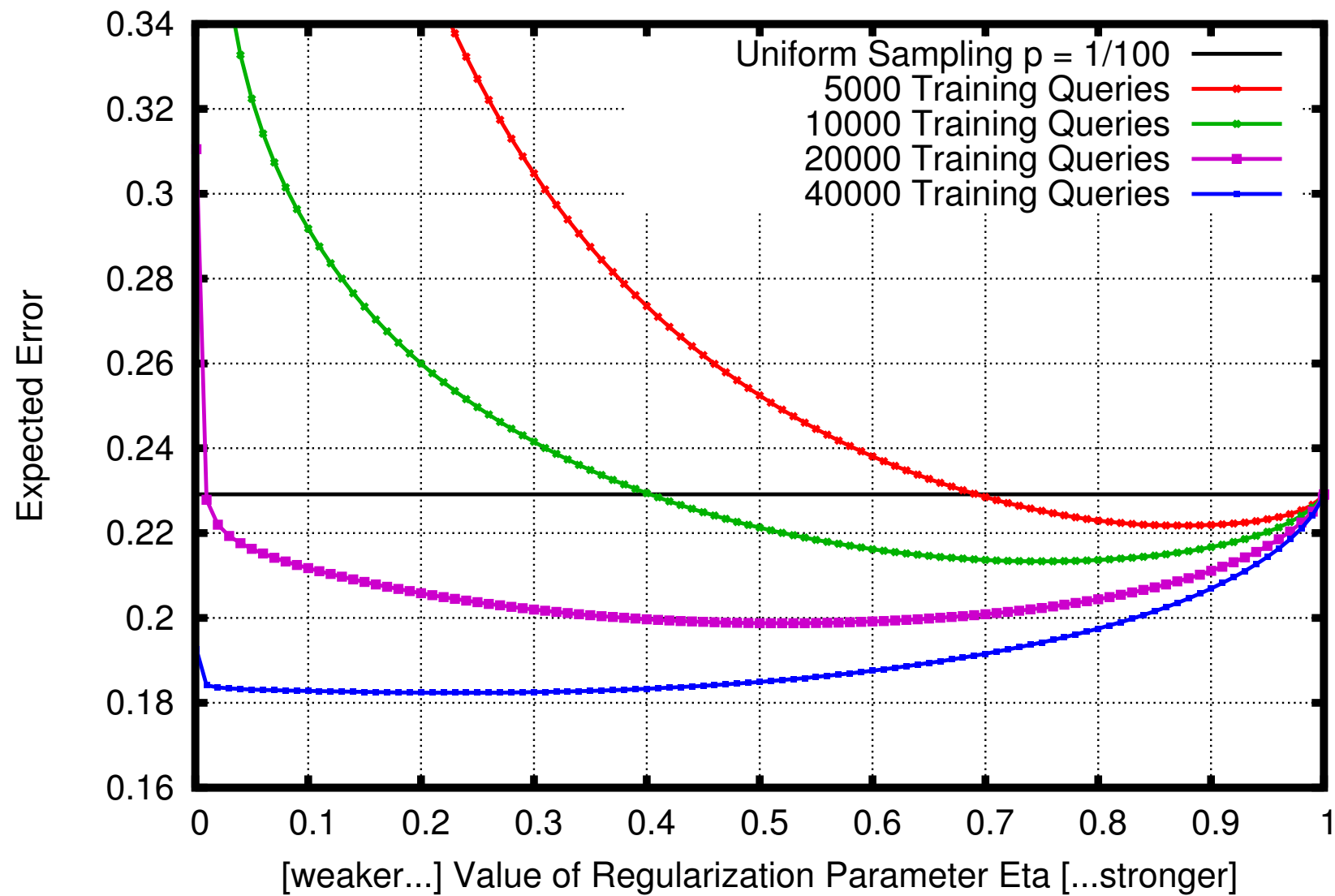
Alg' Risk      Best Risk      Database "badness"      Training Set size

YAHOO!

# RESULTS

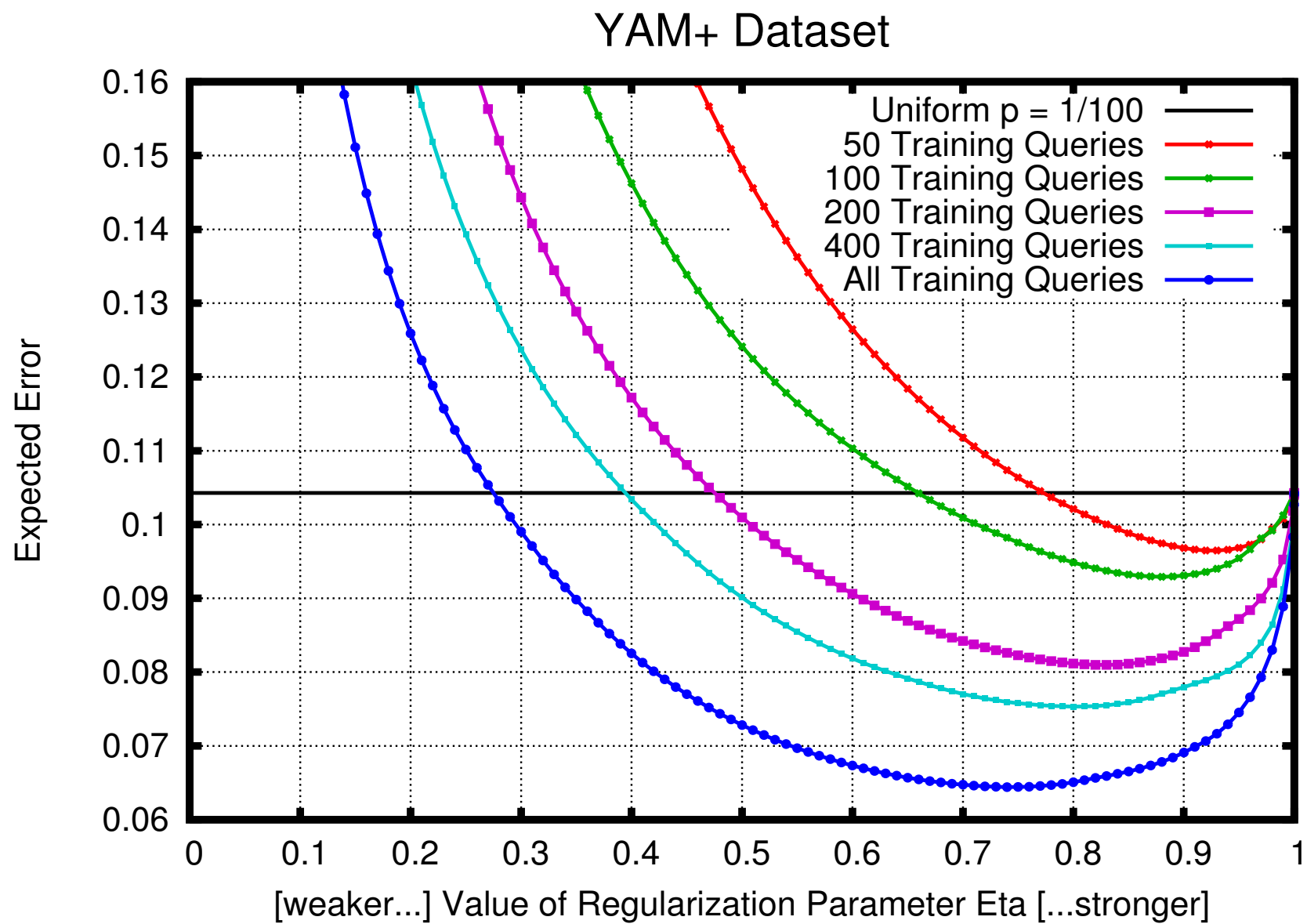| Dataset | Cube | DBLP | YAM+ |
| --- | --- | --- | --- |
| Sampling Rate | 0.1 | 0.01 | 0.01 |
| Uniform Sampling | 0.664 | 0.229 | 0.104 |
| Neyman Allocation | 0.643 | 0.640 | 0.286 |
| Regularized Neyman | 0.582 | 0.228 | 0.102 |
| ERM-$\eta$, small training set | 0.637 | 0.222 | 0.096 |
| ERM-$\rho$, small training set | 0.623 | 0.213 | 0.092 |
| ERM-$\eta$, large training set | **0.233** | 0.182 | 0.064 |
| ERM-$\rho$, large training set | **0.233** | **0.179** | **0.059** |

YAHOO!

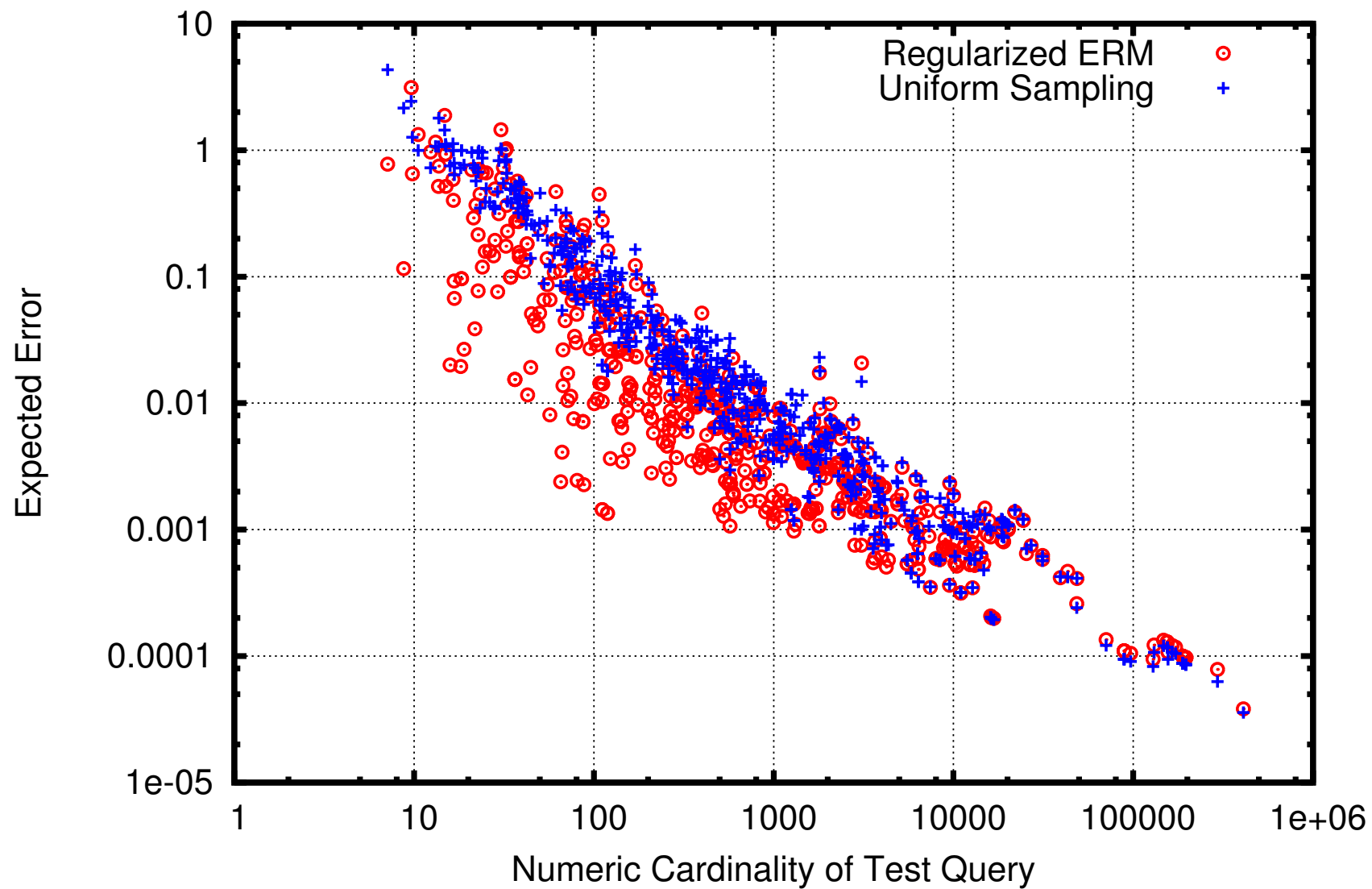# RESULTS

## DBLP Dataset

# RESULTS



Cube Dataset
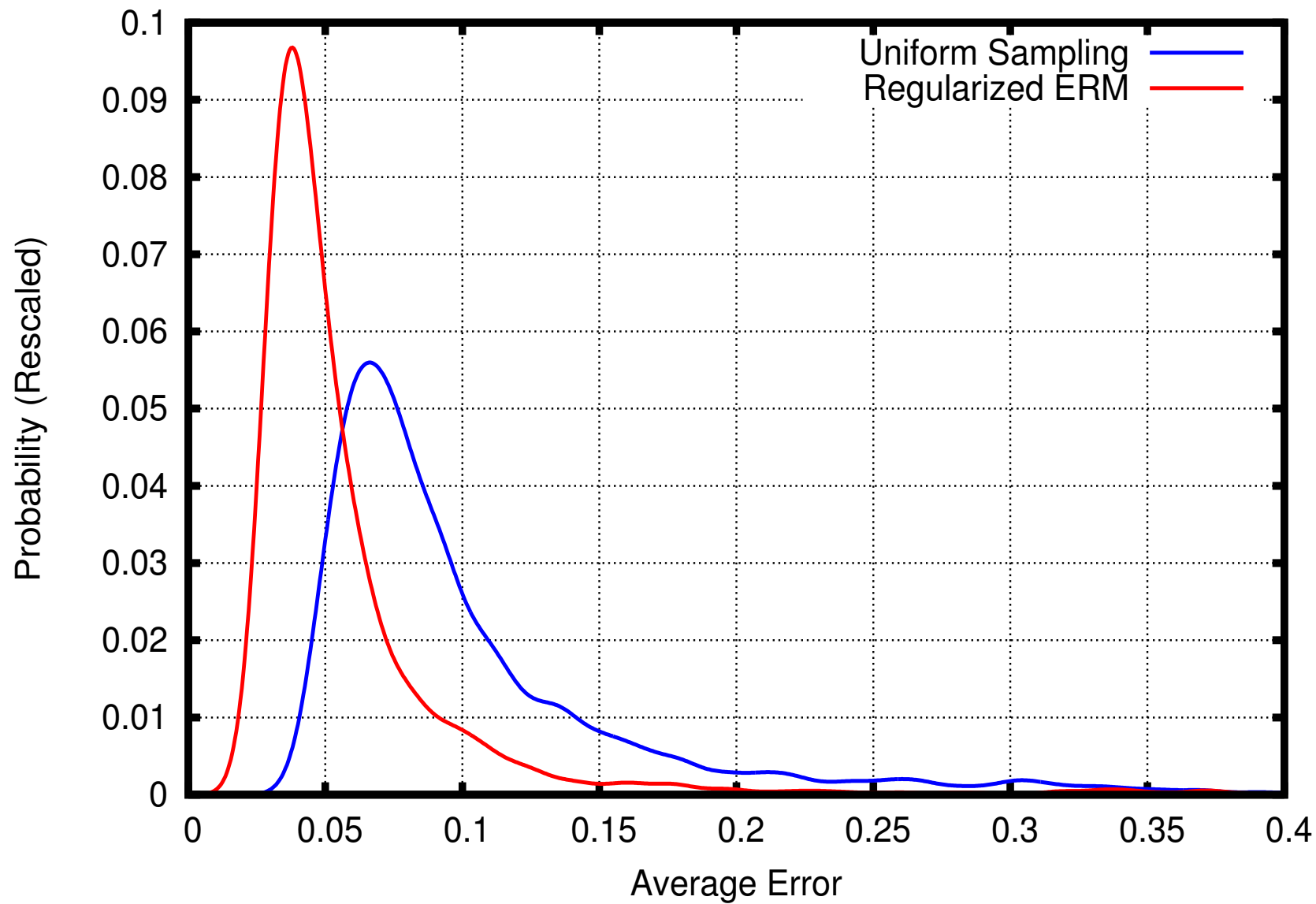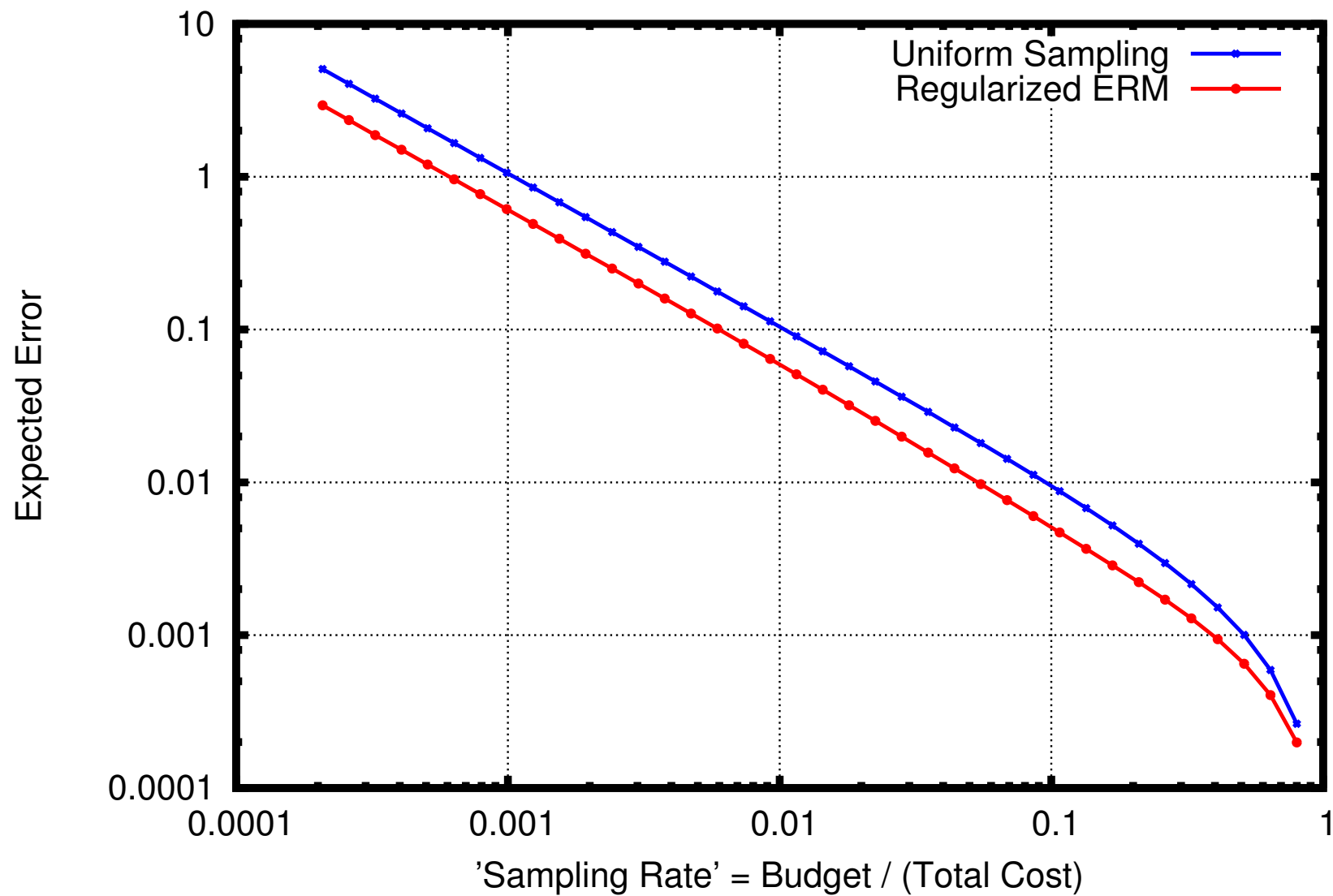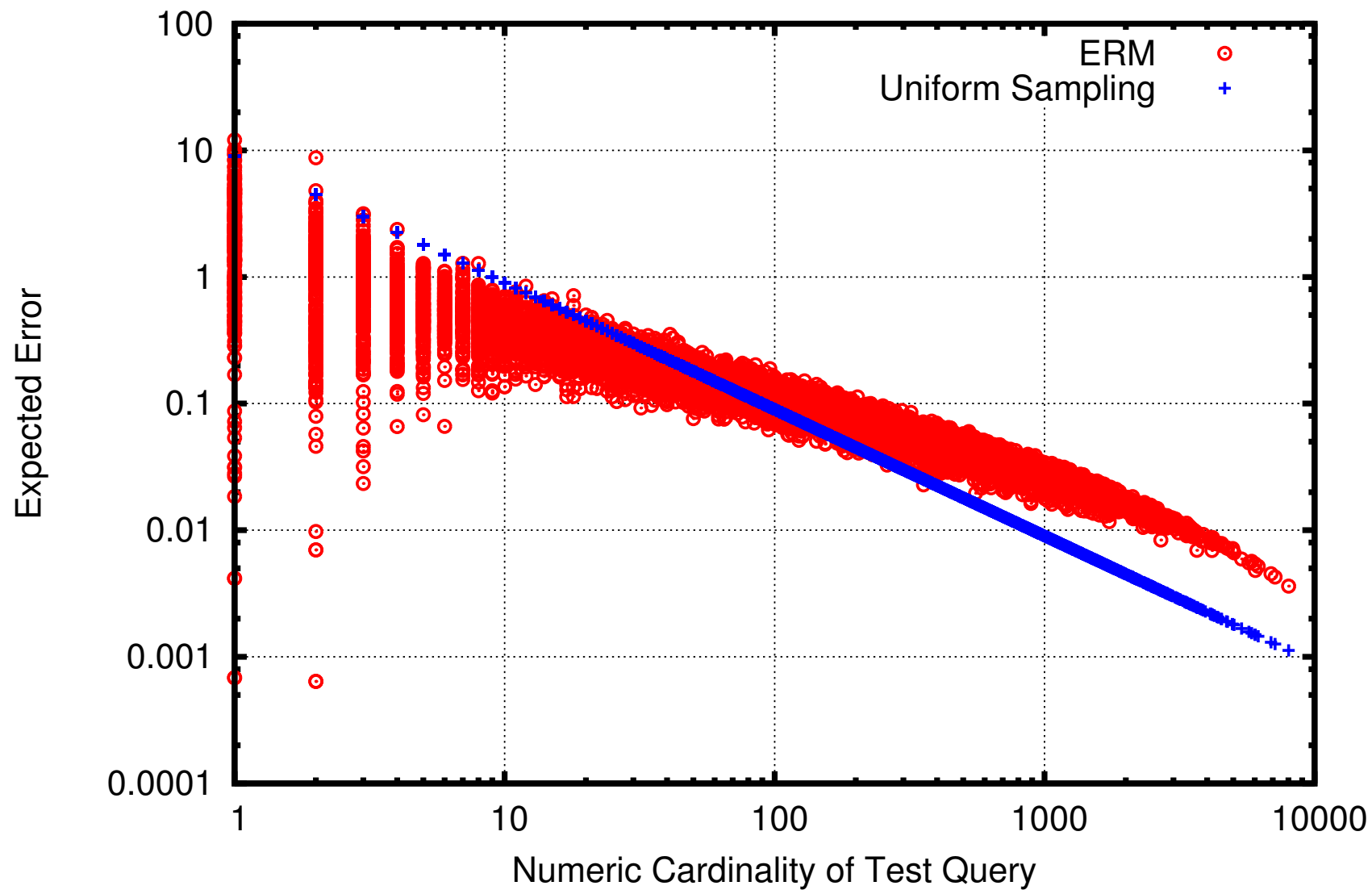
# RESULTS



YAM+ Dataset

# RESULTS



YAM+ Dataset

# RESULTS



YAM+ Dataset

# RESULTS



YAM+ Dataset

# RESULTS

## Cube Dataset



Legend: ERM (red ○), Uniform Sampling (blue +)

Y-axis: Expected Error
X-axis: Numeric Cardinality of Test Query

YAHOO!

YAHOO!