Visualization-Aware Sampling for Very Large Databases

Yongjoo Park Michael Cafarella Barzan Mozafari

University of Michigan, Ann Arbor

Prevalence of Viz-centric Data Analysis



Brain

This graphic depicts countries and territories with 2550 urban populations exceeding 500,000. Circles are scaled in propertion to urban population size.



Population



Genome





A part of 2 billion points



A part of 2 billion points 71 mins !! (matplotlib)



A part of 2 billion points 71 mins !! (matplotlib)

MathGL took 2+ hours for 100M points.



A part of 2 billion points 71 mins !! (matplotlib)

MathGL took 2+ hours for 100M points.

Tableau crashed for 100 million points on a machine with 122GB memory (r3.4xlarge).

Fast Response is Important

Five times more interactions as the response time was reduced [3].

The most skilled user went from 800 interactions per hour with a 1.5-second response time up to 4,300 interactions per hour with a 0.4-second response time. Five times more interactions as the response time was reduced [3].

The most skilled user went from 800 interactions per hour with a 1.5-second response time up to 4,300 interactions per hour with a 0.4-second response time.

+500ms latency \rightarrow 50% less data exploration subconsciously [1]

6 out of 16 subjects did not report a noticeable difference 7 in terms of system responsiveness. Five times more interactions as the response time was reduced [3].

The most skilled user went from 800 interactions per hour with a 1.5-second response time up to 4,300 interactions per hour with a 0.4-second response time.

+500ms latency \rightarrow 50% less data exploration subconsciously [1]

6 out of 16 subjects did not report a noticeable difference 7 in terms of system responsiveness.

A good reason for **my coffee break**.



Visualization Application

Database















We want to

Reduce computational efforts



We want to

Reduce computational efforts

Without affecting visual perception

Two Approaches to Fast Viz



Two Approaches to Fast Viz



Two Approaches to Fast Viz



















Once you generate a sample offline \rightarrow Enjoy fast viz at query time



Once you generate a sample offline → Enjoy fast viz at query time Just like B-tree used to speed up queries





Question 1: What is a good Sample?



Once you generate a sample offline → Enjoy fast viz at query time Just like B-tree used to speed up queries

Question 1: What is a good Sample?

Question 2: How to obtain such a good Sample?

Existing Samplings Fail



Original (2 billion points, 71 mins)

Existing Samplings Fail



Original (2 billion points, 71 mins)


Existing Samplings Fail



Existing Samplings Fail



Original (2 billion points, 71 mins)









Three common visualization-driven tasks [4]:



Other tasks: (1) shape viz, (2) classification, (3) hierarchy understanding, (4) community detection, etc.

Three common visualization-driven tasks [4]:



Other tasks: (1) shape viz, (2) classification, (3) hierarchy understanding, (4) community detection, etc.

1. Motivation

2. Formal Definition of Good Sample

3. Approximation Algorithm to VAS Problem

4. Large-Scale User Study

5. Offline Runtime Analysis











What is a **good** sample (S) of the original dataset (\mathcal{D}) ?



If the visual difference for every circle is small, \rightarrow two viz will look similar.

What is a **good** sample (S) of the original dataset (\mathcal{D}) ?



If the visual difference for every circle is small, \rightarrow two viz will look similar.

Our Goal

To minimize:

$$Loss(S) = \int subloss(x) \ dx$$

where subloss(x) is the viz distance for the circle centered at x.

What is a **good** sample (S) of the original dataset (D)?



If the visual difference for every circle is small, \rightarrow two viz will look similar.

Our Goal

To minimize:

$$Loss(S) = \int subloss(x) \ dx$$

where subloss(x) is the viz distance for the circle centered at x.

What function to use for subloss(x)?









Let's see what happens **around** the location *x*.



Let's see what happens **around** the location *x*.



Let's see what happens **around** the location *x*.



Let's see what happens **around** the location *x*.



Our Choice

$$subloss(x) = \frac{1}{\sum_{s_i \in S} \kappa(x, s_i)}$$

where κ captures the *proximity* between two coordinates.

Good Viz by Solving VAS Problem

Given a budget |S| = K, we want $\underset{s.t. S \subseteq \tilde{D} \land |S| = K}{\operatorname{arg\,min}} \int \frac{1}{\sum_{s_i \in S} \kappa(x, s_i)} dx$

Good Viz by Solving VAS Problem

Given a budget |S| = K, we want $\underset{s.t. S \subseteq \tilde{\mathcal{D}} \land |S| = K}{\operatorname{arg\,min}} \int \frac{1}{\sum_{s_i \in S} \kappa(x, s_i)} dx$

After some approximations and calculations,

Problem VAS

$$\arg \min_{\substack{s.t. \ S \subseteq \overset{S}{\mathcal{D}} \land |S| = K}} \sum_{\substack{s_i, s_j \in S; \ i < j}} \tilde{\kappa}(s_i, s_j) \longrightarrow Loss(S)$$
where

$$\tilde{\kappa}(s_i, s_j) = \int \kappa(x, s_i) \kappa(x, s_j) \, dx$$

Good Viz by Solving VAS Problem

Given a budget |S| = K, we want $\underset{s.t. S \subseteq \tilde{\mathcal{D}} \land |S| = K}{\operatorname{arg\,min}} \int \frac{1}{\sum_{s_i \in S} \kappa(x, s_i)} dx$

After some approximations and calculations,

Problem VAS

$$\arg \min_{\substack{s.t. \ S \subseteq \mathcal{D} \\ k}} \left| \sum_{\substack{s_i, s_j \in S; \ i < j}} \tilde{\kappa}(s_i, s_j) \right| \longrightarrow Loss(S)$$
where

$$\tilde{\kappa}(s_i, s_j) = \int \kappa(x, s_i) \kappa(x, s_j) \, dx$$

NP-hard (evaluated later); How can we solve?

1. Motivation

2. Formal Definition of Good Sample

3. Approximation Algorithm to VAS Problem

4. Large-Scale User Study

5. Offline Runtime Analysis

Nemhauser, et al. [2] algorithm (comes with an error guarantee)

Nemhauser, et al. [2] algorithm (comes with an error guarantee)

As scanning over a dataset,

Nemhauser, et al. [2] algorithm (comes with an error guarantee)

As scanning over a dataset,



When |S| < K

Nemhauser, et al. [2] algorithm (comes with an error guarantee)

As scanning over a dataset,



Nemhauser, et al. [2] algorithm (comes with an error guarantee)

As scanning over a dataset,



Valid replacement: Loss(S') < Loss(S)
Basic (Slow) Approach

Nemhauser, et al. [2] algorithm (comes with an error guarantee)

As scanning over a dataset,



Valid replacement: Loss(S') < Loss(S)

Testing for valid replacements is too Slow: $O(K^3)$ for every point

Our algorithm: Expand/Shrink operation

Our algorithm: Expand/Shrink operation



When |S| < K

Our algorithm: Expand/Shrink operation



Our algorithm: Expand/Shrink operation



Expand/Shrink is Fast: O(K) for every point $\rightarrow O(K^2)$ times faster!!

Our algorithm: Expand/Shrink operation



Expand/Shrink is Fast: O(K) for every point $\rightarrow O(K^2)$ times faster!! The result is exactly same as the Nemhauser's algorithm.

1. Motivation

2. Formal Definition of Good Sample

3. Approximation Algorithm to VAS Problem

4. Large-Scale User Study

5. Offline Runtime Analysis

Questions: User performance, Validity of VAS Problem

Questions: User performance, Validity of VAS Problem

Tasks	Questions	Datasets
trend analysis,	#: 72 - 80	· 22 million
density est,	with different	GPS logs [5]
clustering	sample sizes	· Synthetic

Questions: User performance, Validity of VAS Problem

Tasks	Questions	Datasets	
trend analysis,	#: 72 - 80	· 22 million	
density est,	with different	GPS logs [5]	
clustering	sample sizes	· Synthetic	
<section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header></section-header>	echanical turk	× 40 uni	que

Large-scale User Study: Example



Large-scale User Study: Example



Large-scale User Study: Example



Average user performance

- 1. Uniform Random Sampling
- 2. Stratified Sampling
- 3. Visualization-Aware Sampling (VAS)



Average user performance

- 1. Uniform Random Sampling
- 2. Stratified Sampling
- 3. Visualization-Aware Sampling (VAS)



Trend Analysis

success: if a user answers a question correctly.



Average user performance

- 1. Uniform Random Sampling
- 2. Stratified Sampling
- 3. Visualization-Aware Sampling (VAS)



success: if a user answers a question correctly.

Average user performance

- 1. Uniform Random Sampling
- 2. Stratified Sampling
- 3. Visualization-Aware Sampling (VAS)



success: if a user answers a question correctly.

Strong Correlation (between User Performance and Loss)



Strong Correlation (between User Performance and Loss)



Smaller Loss(S) \rightarrow More successes

Strong Correlation (between User Performance and Loss)



Smaller Loss(S) \rightarrow More successes

(Spearman's rank) correlation coefficient = -0.85

1. Motivation

2. Formal Definition of Good Sample

3. Approximation Algorithm to VAS Problem

4. Large-Scale User Study

5. Offline Runtime Analysis

We tested Mixed Integer Programming (MIP) for exact solutions.



Exact approach (MIP) is prohibitively slow even for small data.

Ran our VAS algorithm over 2 billion points.



Ran our VAS algorithm over 2 billion points.



Minimizes Loss fast.

Ran our VAS algorithm over 2 billion points.



Minimizes Loss fast.

Ran our VAS algorithm over 2 billion points.



Minimizes Loss fast.

Ran our VAS algorithm over 2 billion points.



Minimizes Loss fast.

High-quality viz even at early steps

Ran our VAS algorithm over 2 billion points.



Minimizes Loss fast.

High-quality viz even at early steps

Gradually improves if more time is allowed

Also available at yongjoopark.com/vas

Conclusion

We formulated an important problem: Sampling for Visualization.

We formulated an important problem: Sampling for Visualization.

We proposed an efficient algorithm to VAS problem.

We formulated an important problem: Sampling for Visualization.

We proposed an efficient algorithm to VAS problem.

We demonstrated users achieve superior performance with VAS.

Thank You!

References I

🔋 Z. Liu and I. Heer. The effects of interactive latency on exploratory visual analysis.

TVCG, 2014.

G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing submodular set functions-i.

Mathematical Programming, 1978.



B. Shneiderman.

Response time and display rate in human performance with computers.

CSUR. 1984.
B. Shneiderman.

The eyes have it: A task by data type taxonomy for information visualizations.

In Symposium on Visual Languages, 1996.

Y. Zheng, Q. Li, Y. Chen, X. Xie, and W.-Y. Ma. **Understanding mobility based on gps data.** In *UbiComp*, 2008.