cse 442 - Data Visualization Scalable Visualization



Leilani Battle University of Washington

How can we visualize and interact with **billion+ record** databases in real-time?

Topics

Varieties of "big data" Visualizing Large Datasets **Design Subtleties** Examples **Enabling Real-Time Interaction** imMens ForeCache Trust, but Verify: Optimistic Vis

Varieties of "big data"...





Large DBs have petabytes or more (but median DB still fits in RAM!)

Affects system *and* perceptual scalability

How to manage? Parallel data processing Reduction: Filter, aggregate, sample



Many Columns



Lots of variables (100s-1000s...) Select relevant subset **Dimensionality reduction** Statistical methods can suggest and order related variables

Requires human judgment



Many Columns

Many Sources & Structures













Many Columns

Many Sources & Structures





Many Updates



Many Columns

Many Sources & Structures





Many Updates

How can we visualize and interact with **billion+ record** databases in real-time?

Two Challenges: 1. Effective visual encoding 2. Real-time interaction

Perceptual and interactive scalability should be limited by the chosen resolution of the visualized data, not the number of records.

1. Visualizing Large Datasets

Challenge of Perceptual Scalability

We cannot fit millions+ of data records into a typical image resolution (for example, 500x500 image \rightarrow only 250,000 pixels).

Even if we could, the resulting image would be a wall of ink!

How can we visualize the key properties of a large dataset without rendering every data record?



















2.0 3.0

How to Visualize a Billion+ Records



Decouple the visual complexity from the raw data through aggregation.



1. Bin Divide data domain into discrete "buckets" *Categories*: Already discrete (but watch out for high cardinality) *Numbers*: Choose bin intervals (uniform, quantile, ...) *Time*: Choose time unit: Hour, Day, Month, etc. *Geo*: Bin x, y coordinates *after* cartographic projection

1. Bin Divide data domain into discrete "buckets" *Categories*: Already discrete (but watch out for high cardinality) *Numbers*: Choose bin intervals (uniform, quantile, ...) *Time*: Choose time unit: Hour, Day, Month, etc. *Geo*: Bin x, y coordinates *after* cartographic projection **2. Aggregate** Count, Sum, Average, Min, Max, ...

1. Bin Divide data domain into discrete "buckets" *Categories*: Already discrete (but watch out for high cardinality) Numbers: Choose bin intervals (uniform, quantile, ...) *Time*: Choose time unit: Hour, Day, Month, etc. *Geo*: Bin x, y coordinates *after* cartographic projection **2. Aggregate** Count, Sum, Average, Min, Max, ...

3. Smooth Optional: smooth aggregates [Wickham '13]

1. Bin Divide data domain into discrete "buckets" *Categories*: Already discrete (but watch out for high cardinality) Numbers: Choose bin intervals (uniform, quantile, ...) *Time*: Choose time unit: Hour, Day, Month, etc. *Geo*: Bin x, y coordinates *after* cartographic projection **2. Aggregate** Count, Sum, Average, Min, Max, ...

3. Smooth Optional: smooth aggregates **Indickham '13 4. Plot** Visualize the aggregate values

Binned Plots by Data Type



Design Subtleties...

Hexagonal or Rectangular Bins?



Hex bins better estimate density for 2D plots, but the *improvement is marginal* [Scott 92]. Rectangles support *reuse* and *visual queries*.

Color Scale: Discontinuity after Zero



Standard Color Ramp Counts near zero are white.

Add Discontinuity after Zero Counts near zero remain visible.

Examples

Example: Binned Scatter Plots



Scatterplot Matrix Techniques for Large N [Carr et al. '87]

Example: Basketball Shot Chart



Analytics / Design: Kirk Goldsberry Data Assist Junpin' Matt Adams NBA Shooting 2011-12 [Goldsberry]

[Moritz & Fisher]



Line Chart

Non-Normalized Heatmap

Normalized "DenseLines"

[Moritz & Fisher]

B.2



Repeat for each series





[Moritz & Fisher]



Repeat for each series



[Moritz & Fisher]



Repeat for each series



0	0	0	0	0	0.3	0	0	0	0
0	0	0	0	0	0.3	0.5	0	0	0.5
0	0	0.5	0.5	0.5	0.3	0.5	0.5	0.5	0.5
0	0.5	0.5	0.5	0	0	0	0.5	0.5	0
0.	5 0.5	0	0	0	0	0	0	0	0
0.	5 0	0	0	0	0	0	0	0	0

Approx. Arc-Length Normalized



[Moritz & Fisher]



Approx. Arc-Length Normalized

0

0

0

0

0

0.5 0

0

0

Aggregate

4.5 4

4

4

4 3.5 3.5 3.5 3.5 2





Color

Summary: Perceptual Scalability

Encoding millions+ of data points into one image is a challenge.

Our goal is to reduce the data size to match our target pixel resolution.

Data reduction strategies such as sampling, filtering, modeling, and (binned) aggregation each have their own pros and cons.

Binned aggregation is the most common technique and can be applied to many visualization types.
2. Enabling Real-Time Interaction

Challenge of Interactive Scalability

Each **interaction** triggers a **query** over the original data, which will alter the data we need to render in the visualization.

We may have to retrieve the necessary query results from a **back-end database system** before we can render them.

- 1. Query Database
- 2. Client-Side Indexing / Data Cubes
- 3. Prefetching
- 4. Approximation

1. Query Database Offload to a scalable backend Tableau, for example, issues aggregation queries. Analytical databases are designed for fast, parallel execution. But round-trip queries to the DB may still be too slow... 2. Client-Side Indexing / Data Cubes 3. Prefetching

4. Approximation

- 1. Query Database
- 2. Client-Side Indexing / Data Cubes Query data summaries

Build sorted indices or data cubes to quickly re-calculate aggregations as needed on the client.

- 3. Prefetching
- 4. Approximation

- 1. Query Database
- 2. Client-Side Indexing / Data Cubes
- 3. Prefetching Request data before it is needed

Reduce latency by speculatively querying for data before it is needed. Requires prediction models to guess what is needed.

4. Approximation

- 1. Query Database
- 2. Client-Side Indexing / Data Cubes
- 3. Prefetching
- 4. Approximation Give fast, approximate answers

Reduce latency by computing aggregates on a sample, ideally with approximation bounds characterizing the error.

- 1. Query Database
- 2. Client-Side Indexing / Data Cubes
- 3. Prefetching
- 4. Approximation

These strategies are **not** mutually exclusive! Systems can apply them in tandem.



imMens [Liu, Jiang & Heer '13]

Strategies: Client-Side Data Cubes



Binned Aggregation imMens



Binned Aggregation imMens





















Multivariate Data Tiles1. Send data, not pixels2. Embed multi-dim data

































Full 5-D Cube



For any pair of 1D or 2D binned plots, the maximum number of dimensions needed to support brushing & linking is **four**.

3-D cubes


13 3-D Data Tiles

3-D cubes

3-D data tiles







3-D cubes

3-D data tiles

~17.6M bins (in 352KB!)

Performance Benchmarks



Simulate interaction: brushing & linking across binned plots.

- 4x4 and 5x5 plots - 10 to 50 bins

Measure time from selection to render.

2.3 GHz MacBook Pro NVIDIA GeForce GT 650M Google Chrome v.23.0



Limitations and Questions

But where do the multivariate data tiles come from? They must be provided by a backend server. This can be time-consuming, particularly if supporting deep levels of zooming. imMens assumes that tiles have either been pre-computed or that a backing database can suitably generate them on demand.

Does super-low-latency interaction really matter? Is it worth it to go to all of this trouble? (Short answer: yes!) High latency leads to reduced analytic output [Liu & Heer, InfoVis 2014]

How does interactive latency affect exploratory analysis with visualizations?

Prior Work – Negatives to Latency

Higher latency entails higher action costs, subjects satisfice by selecting strategies that *reduce short-term effort* with no guarantee that the final outcome is optimized. [Gray & Boehm-Davis]

300ms latency reduces the number of Google searches; effect persists for days. [Brutlag et al]

When the cost of acquiring information is increased, subjects change strategy and rely more on working memory. [Ballard et al]

Prior Work – Positives to Latency

When confronted with increased latencies, users resort to more mental planning, at times making fewer errors and performing better on tasks with *verifiable outcomes*. [O'Hara & Payne]

Prior Work – Positives to Latency

When confronted with increased latencies, users resort to more mental planning, at times making fewer errors and performing better on tasks with *verifiable outcomes*. [O'Hara & Payne]

But what about open, exploratory analysis tasks? Addressed by Liu & Heer.

Experiment Design

2 (Latency) x 2 (Scenario) Design *Latency*: +0ms / +500ms *Scenario*: Mobile Check-ins / FAA Flight Delays

Exploratory Analysis Tasks (2 per session) imMens with brush, pan, zoom, adjust scales Users asked to explore data and share findings Log events, record audio and screen capture

16 subjects, all familiar with data analysis + vis









4.5m Mobile Check-Ins



140m FAA Flight Delay Records

Data Collection & Analysis

Event Log Analysis Analyze triggered & processed user input events Assess data set coverage (# unique tiles) Verbal Protocol Analysis Think-aloud protocol: verbalize thought process Transcribe sessions; Code actions and insights Analyze number and type of coded events

Latency Study Results

Higher latency leads to... Reduced user activity and data set coverage Less observation, generalization & hypothesis

Verbal Category	likelihood-ratio test: Chisq(1, N=32)	p value	significance						
Observation	5.4812	0.01922	*		0.283				
Observation (Single View)	1.5706	0.2101			0.070				
Observation (Multiple Views)	3.3119	0.06878			0.215				
Generalization	8.9763	0.002735	**		0.103				
Generalization (Single View)	0.2641	0.6073			0.002				
Generalization (Multiple Views)	8.5054	0.003541	**		0.100				
Hypothesis	8.3999	0.003752	**		0.169				
Question	0.7416	0.3891			0.043				
Interface	0.4651	0.4953		-0.014					
Recall	0.0202	0.8869			0.003				
Simulation	0.6983	0.4033			0.016				
				0.	00	0.05	0.1 L	0 atency	0.15 Coefficient

0.20	0.25	

Latency Study Results

Higher latency leads to... Reduced user activity and data set coverage Less observation, generalization & hypothesis

Different interactions exhibit varied sensitivity to latency. Brushing is highly sensitive!

Latency Study Results

Higher latency leads to... Reduced user activity and data set coverage Less observation, generalization & hypothesis

Different interactions exhibit varied sensitivity to latency. Brushing is highly sensitive! In short: milliseconds matter! And imMens was not a waste of time... 😅

Break Time!

Administrivia

Final Project Schedule

Wed Feb 19 **Proposal** Tues Mar 4 Prototype *Demo Video* Tue Mar 11 *Video Showcase* Thu Mar 13 (in class)



Deliverables Tue Mar 18

Milestone Prototype

Publish work to Gitlab pages for us to examine and share feedback. You **are not** expected to have complete, polished content at this point.

You **are** expected to provide prototype work that communicates your design goals. For example: initial visualizations, sketches, storyboards, and text annotations / idea descriptions.

We should get a sense of what you intend to ultimately submit! Also feel free to ask **us** questions.

ForeCache [Battle, Chang, & Stonebraker '16]

Strategies: Query Database, Data Cubes, Prefetching

ForeCache is also a Data Tile-Based System



Manage a Cache of Tiles from DBMS

Example Tile-Based Views

Key Idea: Model & Predict User Behavior

Classify the User's Analysis Phase
 Foraging: Searching for patterns of interest
 Sensemaking: Closely examine a region-of-interest (ROI)
 Navigation: Transition between levels of detail

2. Predict Which Data Tiles Will be Requested
Train a machine learning classifier (SVM) to predict phase.
The input data is the activity trace of user interactions.

Foraging





Navigation

User zooms in



Navigation

User zooms out

Foraging



































Model 2 T_C, T_D, T_E, T_F







Model 2





Action-Based Tile Recommendations

Idea: user consistently moves in predictable directions










Evaluating ForeCache: A User Study

Participants: 18 earth science researchers Explored NASA MODIS snow cover queries





Retrospective Performance Experiments

Compared response times and prediction accuracy to a non-prefetching baseline and two existing pre-fetching methods:

Momentum







Hotspot

[Doshi et al. 2003]

Results: ForeCache was 20% More Accurate and 88% Faster than **Existing Pre-fetching Methods**



Model ForeCache Hotspot Momentum

What if the data is too large to query in a reasonable time?

Trust, but Verify: Optimistic Vis [Moritz, Fisher, Ding & Wang '17]

Strategies: Query Database, Approximation





Latencies reduce engagement and lead to fewer observations.

The Effect of Interactive Latency. Liu, Heer. IEEE InfoVis 2014.





Approximation: Trade Accuracy for Speed

Approximate query processing (AQP) Uncertainty estimation in statistics Uncertainty visualization Probabilistic programming Approximate hardware

Small chance of error



Pick your poison: 1. Trust the approximation, or 2. Wait for everything to complete.





Optimistic Visualization

Trust but Verify

This glass is half full



What if we think of the issues with approximation as user experience problems?

Optimistic Visualization



- 1. Analysts uses initial estimates.
- 2. Precise queries run in the background.
- 3. System confirms results. Analyst detects errors. Analysts can use approximations and also trust them.

Trust but Verify. Moritz et al. CHI 2017.

Pangloss Implements Optimistic Visualization

Data: FAAData	Heatmap •	\checkmark	Load more data	xpect some er	rors: 2.3%	
Type to filter schema	Y-Avie	What h	ave you learned?			
# Year	Field: DepDelay	Appro	oximate Values			
# Quarter						
# Month	Binning: 64 don't bin		300 —			
# DayofMonth	Sort by key: 🗹		250 —			
# DayOfWeek	Y-Axis		200 —			
i FlightDate	Field: ArrDelay	ay	150			
A UniqueCarrier	Binning: 40 don't bin	rDel	150 —			
# AirlineID	Sort by key:	Ar	100 —			
A Carrier			50 — —	[10,20], [30,40]		
A TailNum	value		0	534k±72k		
# FlightNum	Function: Count		-50 —			
# OriginAirportID	Dereistent Eiltere	_				
# OriginAirportSeqID	Persistent Filters		-20 0 20	40 60 80	100 120 14	40 160 180
# OriginCityMarketID	<pre>e.g.AND(Carrier \$IN\$[ha, dl])(DepDelay>=0)</pre>	_		DepDe	лау	
A Origin		Expect Relativ	cted Error			
A OriginCityName		rtelativ				
A OriginState			300 —			
A OriginStateFips			250 —			
A OriginStateName	Filter set clear		200 —			
# OriginWac		۲. م	200 —			
# DestAirportID		Dels	150 —			
# DestAirportSeqID	Zoom clear Capture as Filter	Arr	100 —			
# DestCityMarketID	(ArrDelay \$RNG\$		50 —			
A Dest	[[-148.80619517543857,390.49205043859655]])		0 —			
A DestCityName	(DepDelay \$RNG\$		-50 —			
A DestState						
A DestStateFips			-20 0 20	40 60 80	100 120 1	40 160 180
A DestStateName				DepDe	lay	
# DestWac	•					



Pangloss Visualizes Uncertainty

Data: FAAData	Heatmap •	 Load more data Expect some errors: 2.3%
Type to filter schema	Y Avic	What have you learned?
# Year		Approximate Values
# Quarter		
# Month	Binning: 64 don't bin	300 —
# DayofMonth	Sort by key: 🗹	250 —
# DayOfWeek	Y-Axis	200 —
≝ FlightDate	Field: ArrDelay X	
A UniqueCarrier	Binning: 40 don't bin	
# AirlineID	Sort by key:	
A Carrier	Value	50 — [10,20], [30,40]
A TailNum	value	0 — 534k±72k
# FlightNum	Function: Count	-50 —
# OriginAirportID	Porcistont Filtors	
# OriginAirportSeqID		-20 0 20 40 60 80 100 120 140 160
# OriginCityMarketID	<pre>e.g. AND(Carrier \$IN\$[ha, dl])(DepDelay>=0)</pre>	
A Origin		Relative
A OriginCityName		
A OriginState		300 —
A OriginStateFips		250 —
A OriginStateName	Filter set clear	200 —
# OriginWac		
# DestAirportID		
# DestAirportSeqID	Zoom clear Capture as Filter	100 —
# DestCityMarketID	(ArrDelay \$RNG\$	5 0 — 5
A Dest	[[-148.80619517543857,390.49205043859655]])	11) o —
A DestCityName	(DepDelay \$RNG\$ (II-19 819658218570382 187 25649037534237II)	-50 — -50 —
A DestState		
A DestStateFips		-20 0 20 40 60 80 100 120 140 160
		DepDelay
A DestStateName		



Pangloss shows a History of Previous Charts

Data: FAAData	Heatmap •	✓ Load more data Expect some errors: 2.3%
Type to filter schema	Y-Avis	What have you learned?
# Year	Field: DepDelay	Approximate Values
# Quarter		
# Month	Binning: 64 don't bin	300 —
# DayofMonth	Sort by key: 🗹	250 —
# DayOfWeek	Y-Axis	200 —
i FlightDate	Field: ArrDelay	
A UniqueCarrier	Binning: 40 don't bin	
# AirlineID	Sort by key:	¥ 100 —
A Carrier		50 — [10,20], [30,40]
A TailNum	value	0 — 534k±72k
# FlightNum	Function: Count	-50 —
# OriginAirportID	Paraistant Filtara	
# OriginAirportSeqID	Feisistent Filters	-20 0 20 40 60 80 100 120 140 160 180
# OriginCityMarketID	<pre>e.g.AND(Carrier \$IN\$[ha, dl])(DepDelay>=0)</pre>	DepDelay
A Origin		Expected Error
A OriginCityName		
A OriginState		300 —
A OriginStateFips		250 —
A OriginStateName	Filter set clear	200 —
# OriginWac		
# DestAirportID		
# DestAirportSeqID	Zoom clear Capture as Filter	100 —
# DestCityMarketID	(ArrDelay \$RNG\$	50 —
A Dest	[[-148.80619517543857,390.49205043859655]])	0-
A DestCityName	(DepDelay \$RNG\$	-50 —
A DestState		
A DestStateFips		-20 0 20 40 60 80 100 120 140 160 180
A DestStateName		DepDelay
# DestWac	•	



In Pangloss, Analysts can Confirm results





Evaluation

Case studies with teams at Microsoft who brought in *their own data*.

Approximation works

"seeing something right away at first glimpse is really great"

Need for guarantees

"[with a competitor] I was willing to wait 70-80 seconds. It wasn't ideally interactive, but it meant I was looking at all the data."

Optimism works

"I was thinking what to do next— and I saw that it had loaded, so I went back and checked it ... [the passive update is] very nice for not interrupting your workflow."

In Conclusion...

Two Challenges: 1. Effective visual encoding 2. Real-time interaction

Perceptual and interactive scalability should be limited by the chosen resolution of the visualized data, not the number of records.

Bin > Aggregate (> Smooth) > Plot

1. Bin Divide data domain into discrete "buckets"

- 2. Aggregate Count, Sum, Average, Min, Max, ...
- 3. Smooth Optional: smooth aggregates [M/ickham '13]
- 4. Plot Visualize the aggregate values

Interactive Scalability Strategies

- 1. Query Database
- 2. Client-Side Indexing / Data Cubes
- 3. Prefetching
- 4. Approximation

These strategies are not mutually exclusive! Systems can apply them in tandem.