cse 512 - Data Visualization Scalable Visualization



Jeffrey Heer University of Washington

Session Outline

The Varieties of "Big Data" Scalable Plotting Techniques Scalable Interaction Why Latency Matters Sampling Methods



The Varieties of "Big Data"

Tall Data

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

How to manage? Parallel data processing Reduction: Filter, aggregate Sample or approximate

Not just about systems. Consider perceptual / cognitive scalability.

Tall Data



Wide data



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

Requires human judgment

Tall DataWide data





Diverse data





Wide data



Diverse data



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.

Scalable Plotting Techniques

















3.0 _T

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 [Wickham '13]

4. Plot Visualize the aggregate values

Binned Plots by Data Type







Binned Aggregation (imMens) [Liu, Jiang, Heer '13]



Binned Aggregation

[Liu, Jiang, Heer '13]

Example: Binned Scatter Plots



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

Example: Basketball Shot Chart



NBA Shooting 2011-12 [Goldsberry]

Time Series











Insight: the resolution is bound by the number of pixels.



Insight: the resolution is bound by the number of pixels.

1. Compute average value per pixel (1 point/pixel) ...this may miss extreme (min, max) values



Insight: the resolution is bound by the number of pixels.

- 1. Compute average value per pixel (1 point/pixel) ...this may miss extreme (min, max) values
- 2. Plot min/max values per pixel (2 points/pixel) ...this does better, but still misrepresents





Insight: the resolution is bound by the number of pixels.

- 1. Compute average value per pixel (1 point/pixel) ...this may miss extreme (min, max) values
- 2. Plot min/max values per pixel (2 points/pixel) ...this does better, but still misrepresents
- 3. M4: min/max values & timestamps (4 points/pixel)

...this provides provable fidelity to the full data!






Data Reduction in the Database

```
SELECT t, v FROM Q JOIN
(SELECT round($w*(t-$t1)/($t2-$t1)) as k, --define key
        min(v) as v_min, max(v) as v_max, --get min, max
        min(t) as t_min, max(t) as t_max --get 1st,last
        FROM O GROUP BY k) as OA
                                           --group by k
ON k = round(\frac{(+++)}{(+++)})
                                           --join on k
       AND (v = v_min \ OR \ v = v_max \ OR
                                           --&(min|max|
            t = t \min OR t = t \max
                                           -- 1st|last)
```

Q: query that returns a time series (t,v)
\$w: chart width in pixels
\$t1, \$t2: global min/max timestamps

Time Series: 1M samples, 1 sample/second



M4: 1M samples -> 2,653 plotted points



But what about multiple time-series?



Perceptual scalability breaks down...



The non-normalized heatmap suffers from artifacts, seen as vertical stripes.

Binned charts convey high points across the top, a collective dip in stocks during the crash of 2008, and two distinct bands of \$25 and \$15 stocks.

















Approx. Arc-Length Normalized



Approx. Arc-Length Normalized

Aggregate

Color



The density of the second group appears to increase to the right! Without normalization, the steep lines are over-represented.

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.

Color / Opacity Ramps



Linear interpolation in RGBA is not perceptually linear.



Perceptual color spaces approximate perceptual linearity.

Scalable Interaction

- 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 ... or alternative data frame implementation

Python: Vaex, Polars, Modin, cuDF

R: <u>dbplyr</u>

All: <u>DuckDB</u>

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.

Client-Side Indexes

Binned Aggregation

[Liu, Jiang, Heer '13]

























Multivariate Data Tiles

- 1. Send data, not pixels
- 2. 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**.



13 3-D Data Tiles



(in 352KB!)

5 dimensions x 50 bins/dim x 25 plots



Limitations and Questions

But where do the multivariate data tiles come from?

They must be provided by a backend server. This can be timeconsuming, particularly if supporting deep levels of zooming.

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]

Why Latency Matters

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

Higher latency entails higher action costs, subjects satisfice by selecting strategies that *reduce shortterm 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]

Higher latency entails higher action costs, subjects satisfice by selecting strategies that *reduce shortterm 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]

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

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]

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

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?

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

Higher latency leads to...

Higher latency leads to...

Reduced user activity and data set coverage

Higher latency leads to...

Reduced user activity and data set coverage Significantly fewer brushing actions

Higher latency leads to...

Reduced user activity and data set coverage Significantly fewer brushing actions Less observation, generalization & hypothesis

Verbal Category	likelihood-ratio test:	n 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		-0.0 <mark>14</mark>								
Recall	0.0202	0.8869			0.003							
Simulation	0.6983	0.4033			0.016							
				0.	00	0.05	0.1	0	0.15	0.20	0.25	

Latency Coefficient

Higher latency leads to...

Reduced user activity and data set coverage Significantly fewer brushing actions Less observation, generalization & hypothesis **Interaction effect**: Exposure to delay reduces subsequent performance in low-latency interface.

Higher latency leads to...

Reduced user activity and data set coverage Significantly fewer brushing actions Less observation, generalization & hypothesis

Interaction effect: Exposure to delay reduces subsequent performance in low-latency interface.

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

Higher latency leads to...

Reduced user activity and data set coverage Significantly fewer brushing actions Less observation, generalization & hypothesis

- **Interaction effect**: Exposure to delay reduces subsequent performance in low-latency interface.
- **Different interactions exhibit varied sensitivity** to latency. Brushing is highly sensitive!
- In short: milliseconds matter! And optimizing for latency was not a waste of time... 😅

Sampling Methods

Visual Data Exploration



Common Sampling Methods

First-N: Useful for transformation, but not inference.

Random: Good default, but may miss features of interest. Possible in one pass via reservoir sampling, or faster if stored in randomized order.

Stratified: Sample within groups, ensure coverage and balance across those categories.



Binned Aggregation

[Liu, Jiang, Heer '13]


Online Aggregation [Hellerstein, Haas, Wang '97]

Provide dynamic, *progressive* results as queries run: see results over growing samples. Visualize current results with confidence intervals to convey uncertainty of estimate.

Challenge: difficult to ensure truly random sampling.



What if 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

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



Optimistic Visualization

Trust but Verify

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

Optimistic Visualization

Trust but Verify. Moritz et al. CHI 2017.



- 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.

Optimistic Visualization



Visualize Uncertainty



Show a History of Previous Charts



Help Analysts Confirm Results

Data: FAAData 🔻	Heatmap	What have you learned?	°2_× ^
Type to filter schema	X-Axis	The visualization is read only because you're looking at the history. <u>Return to the working vis</u>	
# Year	Field: DepDelay	Fract Data	
# Quarter	Rinning: 54	400-	The second
# Month	Dimining. Pr	350 —	Exact data loaded (51s)
# DayofMonth	Sort by key: 🕑	300 — 20M	2 ×
# DayOfWeek	Y-Axis	250 — 15M	1 E
FlightDate	Field: ArrDelay	200 — 10M	
A UniqueCarrier	Binning: 64	5M	Contraction of the second s
# AirlineID	Sort by key:	100 — 0	Exact data loaded (94s)
A Carrier	Mahar	50-	
A TailNum	value	-50	
# FlightNum	Function: Count	-100 —	
# OriginAirportID	Paraistant Filtara	-150	
# OriginAirportSeqID	reisistent riiters	-20 0 20 40 60 80 100 120 140 160 180	·
# OriginCityMarketID	<pre>e.g. AND(Carrier \$IN\$[ha, d1])(DepDelay>=0)</pre>	DepDelay	Exact data loaded (48s)
A Origin		Difference to Approximate Data	
A OriginCityName		Relative	
A OriginState			
A OriginStateFips		300 40k	
A OriginStateName		20k	•
# OriginWac		200 —///.	
# DestAirportID	Zoom	-40k	You are looking at the history and
# DestAirportSeqID	(ArrDelay \$RNG\$		cannot make any changes.
# DestCityMarketID	[[-148.80619517543857,390.49205043859655]])	50	
A Dest	(DepDelay \$RNG\$	0	i and
A DestCityName	[[[-14.619056216570562,167.25649057554257]])	-50	
A DestState		150	
A DestStateFips		-20 0 20 40 60 80 100 120 140 160 180	
A DestStateName		DepDelay	Return to editing
# DestWac			
A CRSDepTime			Clear History Reset App

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 [Wickham '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.