# Interactive A Data Analysis

**Jeffrey Heer** @jeffrey\_heer

# Interactive A Data Analysis

Jeffrey Heer Stanford University

# 

Jeffrey Heer University of Washington

# How much data (bytes) did we produce in 2010?

### **2010:** 1,200 exabytes

Gantz et al, 2008, 2010

## **2010:** 1,200 exabytes 10x increase over 5 years

Gantz et al, 2008, 2010

The ability to take data—to be able to understand it, to process it, to extract value from it, to **visualize** it, to **communicate** it that's going to be a hugely important skill in the next decades, ... because now we really do have essentially free and ubiquitous data. So the complimentary scarce factor is the ability to understand that data and extract value from it.

> Hal Varian, Google's Chief Economist The McKinsey Quαrterly, Jan 2009



🗹 Animate

#### 

Graph	Viewer
-------	--------



00	Graph Viewer
Graph Viewer	
Roll-up by:	ى ئى ئى بىرى بىرى بىرى بىلى بىرىيە بىرى بىرى بىرى بىرى بىرى بىرى بىرى بىر
All	에 가장에 있는 것이 가장에 있는 것이 같이 있는 것이 있 같이 같이 같이 있는 것이 같이 있는 것이 같이 있는 것이 있
Visualization:	
Matrix 🛟	
Sort by:	
None	- A 「計論」は 20世紀 (金融) 2010年 - A 「A」 「A」 「A」 「A」 「A」 「A」 「A」 A」 A
Edge centrality filters:	
•	(1) Sing states (1) Sing st
	(1) 新生产的合同的 网络白刺纲 电子和电力 计可以注意 计可能分析的 计可能分析
	n en de la companya de la substance de pública de la defenda de la defenda de la companya de la substance de l A 1997 - La companya de la companya
	n fernanden er en
	· · · · · · · · · · · · · · · · · · ·
	e se la companya de la companya de La companya de la comp

### Visualization

#### Acquisition

Cleaning

Integration

Visualization

Modeling

Presentation

Dissemination







How might we support expressive and effective visualization designs?

### d3.js Data-Driven Documents



with Mike Bostock & Vadim Ogievetsky

) (	GitHub, Inc. [US] http	os://github.com/popular/w	vatched		5
,£	Coffee-script Unfancy JavaScript	jashkenas	4,733 watching	all commits sy owner	52 week participation
<b>*</b>	Realtime application frame cross-browser fallbacks su	LearnBoost work for Node.JS, with HTML pport.	4,454 watching 5 WebSockets and	all commits by owner	52 week participation
, C	A collection of useful .gitig	github nore templates	4,437 watching	■ all commits ■ commits by owner	52 week participation
~~	underscore JavaScript's utility _ belt	documentcloud	4,434 watching	all commits a commits by owner	52 week participation
F	d3 A JavaScript visualization	mbostock library for HTML and SVG.	4,286 watching	all commits Commits by owner	52 week participation
Q	Linux kernel source tree	torvalds	4,280 watching	🖩 all commits 🔳 commits by owner	52 week participation
Ś	Symfony The Symfony2 PHP frame	symfony work	4,209 watching	all commits Commits by owner	52 week participation
0	paperclip Easy file attachment mana	thoughtbot gement for ActiveRecord	4,031 watching	all commits Commits by owner	52 week participation
	redis Redis is an in-memory data key-value, but many differe	antirez abase that persists on disk. T ent kind of values are support	3,941 watching The data model is ed: Strings, Lists,	I all commits commits by owner	52 week participation



#### 512 Paths to the White House

Select a winner in the most competitive states below to see all the paths to victory available for either candidate.



## stamen design



## stamen design

**d**3

d3

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



**Binned Aggregation** 

### *imMens:* Real-Time Visual Querying of Big Data with Zhicheng (Leo) Liu & Biye Jiang









*imMens:* Real-Time Visual Querying of Big Data with Zhicheng (Leo) Liu & Biye Jiang







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



13 3-D Data Tiles

→ ~17.6M bins (in 352KB!)


































3 5 7

1

9 11 13 15 17 19 21 23

























1 3 5 7 9 11 13 15 17 19 21 23 25 27 29 31

0





9 11 13 15 17 19 21 23

5 7

1 3









## Query & Render on GPU via WebGL



## Query & Render on GPU via WebGL



Bind data tiles as image textures. Compute aggregates in parallel on GPU.



### 60 50 Average Frames per Second 40 30 20 10 **In-Memory Data Cube** 0 10k 100k 1M 10M 100M 1B Number of Data Points











I spend more than half of my time integrating, cleansing and transforming data without doing any actual analysis. Most of the time I'm lucky if I get to do any "analysis" at all.

Anonymous Data Scientist from our interview study, 2012



Bureau http:/,	of Justice Stati /bjs.ojp.usdoj.go	stics – Data Online w/			
Report	ed crime in Alaba	ıma			
Year 2004 2005 2006 2007 2008	Population 4525375 4029.3 4548327 3900 4599030 3937 4627851 3974.9 4661900 4081.9	Property crime rate 987 2732.4 309.9 955.8 2656 289 968.9 2645.1 322.9 980.2 2687 307.7 1080.7 2712.6 288.6	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
Report	ed crime in Alask	a			
Year 2004 2005 2006 2007 2008	Population 657755 3370.9 663253 3615 670053 3582 683478 3373.9 686293 2928.3	Property crime rate 573.6 2456.7 340.6 622.8 2601 391 615.2 2588.5 378.3 538.9 2480 355.1 470.9 2219.9 237.5	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
Report	ed crime in Arizo	ina			
Year 2004 2005 2006 2007 2008	Population 5739879 5073.3 5953007 4827 6166318 4741.6 6338755 4502.6 6500180 4087.3	Property crime rate 991 3118.7 963.5 946.2 2958 922 953 2874.1 914.4 935.4 2780.5 786.7 894.2 2605.3 587.8	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
Report	ed crime in Arkar	isas			
Year 2004 2005 2006 2007 2008	Population 2750000 4033.1 2775708 4068 2810872 4021.6 2834797 3945.5 2855390 3843.7	Property crime rate 1096.4 2699.7 237 1085.1 2720 262 1154.4 2596.7 270.4 1124.4 2574.6 246.5 1182.7 2433.4 227.6	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
Report	ed crime in Calif	ornia			
Year 2004 2005 2006 2007 2008	Population 35842038 36154147 36457549 36553215 36756666	Property crime rate 3423.9 686.1 2033.1 3321 692.9 1915 3175.2 676.9 1831.5 3032.6 648.4 1784.1 2940.3 646.8 1769.8	Burglary rate 704.8 712 666.8 600.2 523.8	Larceny-theft rate	Motor vehicle theft rate
Report	ed crime in Color	ado			
Year 2004	Population 4601821 3918.5	Property crime rate 717.3 2679.5 521.6	Burglary rate	Larceny-theft rate	Motor vehicle theft rate

## **DataWrangler**

uggestions	rows: 408 prev next	
	The Year	Property_crime_rate
Delete rows 8.10	1 Reported crime in Alabama	
	2	
Delete empty rows	3 2004	4029.3
belete empty rous	4 2005	3900
Delete rows where Property_crime_rate	5 2006	3937
is null	6 2007	3974.9
	7 2008	4081.9
Delete rows where Year is null	8	
	9 Reported crime in Alaska	
cript Expor	10	
Split data repeatedly on newline into	11 2004	3370.9
rows	12 2005	3615
Calls data associatedly an 11	13 2006	3582
Split data repeatedly on "	14 2007	3373.9

with Sean Kandel, Philip Guo, Andreas Paepcke & Joe Hellerstein

# Wrangler in 2 Parts...

Data transformation language
Map Operations – split, merge, extract, drop
Reshaping – fold, pivot (cross-tabulate)
Lookups & Joins – e.g., FIPS code to US state
Sorting, Aggregation, etc.

# Wrangler in 2 Parts...

1. Data transformation language

 Mixed-initiative interface for data transforms User: Selects data elements of interest System: Suggests applicable transforms via search over the space of viable transforms Enable visual preview and refinement





Extract

Impute

Reshape

Extrac	t										
Impute	е										
Reshape	e										
	0	1	2	3	4	5	6	7	8	9	10
	User	Study	Task C	Comple	tion Tir	ne (mii	nutes)	•	Wrangle	er 🛑 Exc	cel





Median completion time for Wrangler at least **twice as fast** in all tasks (*p* < 0.001).

Suggestions and visual previews used heavily.

## DataWrangler

Suggestions	rows: 408 pr
	#
Delete rows 910	1 Reported crim
Delete Tows 0,10	2
Delete empty rows Delete rows where Property_crime_rate is null	3 2004
	4 2005
	5 2006
	6 2007
	7 2008
Delete rows where Year is null	8
	9 Reported crim
Script Expo	10
Split data repeatedly on newline into	11 2004
rows	12 2005
b Colit data repeatedly on U	13 2006
split data repeatedly on "	14 2007

	rows: 408 prev next	
	🇰 Year 🛔	Property_crime_rate
1	Reported crime in Alabama	
2		
3	2004	4029.3
4	2005	3900
5	2006	3937
6	2007	3974.9
7	2008	4081.9
8		
9	Reported crime in Alaska	
10		
11	2004	3370.9
12	2005	3615
13	2006	3582
14	2007	3373.9

http://vis.stanford.edu/wrangler






## 

Jeffrey Heer http://idl.cs.washington.edu



## Effective statistical models for syntactic and semantic disambiguation

Student: Kristina Nikolova Toutanova Advisor: Christopher D. Manning

Computer Science (2005)

Keywords: Syntactic, Semantic, Tree kernels, Parsing

Abstract:

This thesis focuses on building effective statistical models for disambiguation of sophisticated syntactic and semantic natural language (NL) structures. We advance the state of the art in several domains by (i) choosing representations that encode domain knowledge more effectively and (ii) developing machine learning algorithms that deal with the specific properties of NL disambiguation tasks--sparsity of training data and large, structured spaces of hidden labels. For the task of syntactic disambiguation, we propose a novel representation of parse trees that connects the words of the sentence with the hidden syntactic structure in a direct way. Experimental evaluation on parse selection for a Head Driven Phrase Structure Grammar shows the new representation achieves superior performance compared to previous models. For the task of disambiguating the semantic role structure of verbs, we build a more accurate model, which captures the knowledge that the semantic frame of a verb is a joint structure with strong dependencies between arguments. We achieve this using a Conditional Random Field without Markov independence assumptions on the sequence of semantic role labels. To address the sparsity problem in machine learning for NL, we develop a method for incorporating many additional sources of information, using Markov chains in the space of words. The Markov chain framework makes it possible to combine multiple knowledge sources, to learn how much to trust each of them, and to chain inferences together. It achieves large gains in the task of disambiguating prepositional phrase attachments.

