# Outline

- Nan: Fundamentals and State-of-the-art (25-30 minutes)
- **Eugene:** Efficient, Effective and Interactive Visualizations (60-65 minutes)
- **Guoliang:** Recommendation (~60 minutes)
- Nan: Uncertainty, collaborative, and immersive data visualizations (~30 minutes)
   *uncertainty* and "*cleaning*" bad data visualizations
  - collaborative data visualization
  - immersive data visualization

## In (Data) Science, The Only Certainty is

# In (Data) Science, The Only Certainty is Uncertainty

# In (Data) Science, The Only Certainty is Uncertainty

#### Doubt



#### Bad

ordenings (neary, samon sinekr), or no props at a
 **5% Spening Joke** Remember the one about the shoe salesmen wh
 wort to Africa in the 1900a? That's how Benjami

boot classical music.
 5%
Spontaneous Moment
Don't overprepare. Tease the guy in the front row
"You could light us a village with this guy's eyes").
Commend the stagehand who handles the human
with work investi

• 5% Statement of Utter Certainty People come for answers –-give 'em what they want, as Shawa Achor dd: "By training your brain ... we can reverse the formula for happiness and success."

• 12% Snappy Refrain The TED equivalent of "I have a dream." Exam "People don't buy what you do; they buy why y it." Repeat 7x.

23%
 Personal Failure
 Be relatable. We want to know about that nervo breakdown. Or at least the time you didn't fit in summer camp.
 49%

Contrarian Thesis Wait a sec --we should be playing more videogames? The more choices we have, the worse off we are? TED is where conventional wisdom goes to fee



### Error



#### incomplete data, lack of knowledge, variability ...

# In (Data) Science, The Only Certainty is Uncertainty

### Doubt





### Error



### incomplete data, lack of knowledge, variability ...

• What does uncertainty mean?

1%

12%

49%

- How should I visualize it?
- What can go wrong?

- Diversified meanings.
- It depends.
- Everything.



*In Pursuit of Error: A Survey of Uncertainty Visualization Evaluation.* **Jessica Hullman et al.**, InfoVis 2018. Uncertainty. CSE 442 - Data Visualization. <u>https://courses.cs.washington.edu/courses/cse442/17au/lectures/CSE442-Uncertainty.pdf</u>

💭 Jupyter	USer1 (unsaved changes)
File Edit	View Insert Cell Kernel Widgets Help
8 + % 2	] 🗈 🛧 🖌 № Run 🔳 C 🕨 Code 💠 📼
In [2]:	<pre>import warnings import sys</pre>
	<pre>warnings.filterwarnings('ignore') sys.path.append('/')</pre>
In [ ]:	
In [3]:	import pandas as pd
	<pre>df = pd.read_csv('pima-indians-diabetes.csv')</pre>
In [6]:	<pre>df['Body mass index'].plot.hist()</pre>
Out[6]:	<matplotlib.axessubplots.axessubplot 0x113c364e0="" at=""></matplotlib.axessubplots.axessubplot>
	250 -
	200 -
	2 150 -
	100 -
	50 -
	0 10 20 30 40 50 60 70



sensor, light, voltage, humidity, temperature

#### 54 sensors 3.2k readings/hour





*Scorpion: Explaining Away Outliers in Aggregate Queries*, **Eugene Wu et al.**, VLDB 2013 A Demonstration of DBWipes: Clean as You Query, **Eugene Wu et al.**, VLDB demo 2012

sensor, light, voltage, humidity, temperature

#### 54 sensors 3.2k readings/hour



#### **Influential Predicates Problem**

Given a select-project-group-by query, user inputs **O**, hold-out results **H**, it is find the predicate **p** that maximizes the influence on **O**, and minimizes the influence on **H**.



*Scorpion: Explaining Away Outliers in Aggregate Queries*, **Eugene Wu et al.**, VLDB 2013 A Demonstration of DBWipes: Clean as You Query, **Eugene Wu et al.**, VLDB demo 2012

sensor, light, voltage, humidity, temperature

#### 54 sensors 3.2k readings/hour

DBWipes + Scorpion! toggle scorpion

incheck a

date 🗸

epoch

num 🗸

humidity

hr timestamr

id

num 🔽

light

Facets [?]



#### **Influential Predicates Problem**

Given a select-project-group-by query, user inputs **O**, hold-out results **H**, it is find the predicate **p** that maximizes the influence on **O**, and minimizes the influence on **H**.



*Scorpion: Explaining Away Outliers in Aggregate Queries*, **Eugene Wu et al.**, VLDB 2013 A Demonstration of DBWipes: Clean as You Query, **Eugene Wu et al.**, VLDB demo 2012

sensor, light, voltage, humidity, temperature

#### 54 sensors 3.2k readings/hour



### "sensors with low voltage"





ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_1$	2013	NADEEF	ACM SIGMOD	QCRI	174.0
$t_2$	2013	NADEEF	SIGMOD Conf.	QCRI, HBKU	1740
$t_3$	2013	NADEEF	SIGMOD	QCRI HBKU	174.0
$t_4$	2013	KuaFu	ICDE 2013	Microsoft	15.0
$t_5$	2013	TsingNUS	SIGMOD'13	Tsinghua	13.0
$t_6$	2013	TsingNUS	SIGMOD'13	THU	13.0
$t_7$	2014	SeeDB	VLDB	Stanford Univ.	N.A.
$t_8$	2014	SeeDB	Very Large Data Bases	Stanford	55.0
$t_9$	2015	Elaps	ICDE	NUS	42.0
$t_{10}$	2015	Elaps	IEEE ICDE Conf. 2015	CS@NUS	44.0

#### TABLE I

AN EXCERPT OF PUBLICATIONS (DIRTY)

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_{123}$	2013	NADEEF	SIGMOD	QCRI	174.0
$t_4$	2013	KuaFu	ICDE	Microsoft	15.0
$t_{56}$	2013	TsingNUS	SIGMOD	Tsinghua	13.0
$t_{78}$	2014	SeeDB	VLDB	Stantford Univ.	55.0
$t_{910}$	2015	Elaps	ICDE	NUS	43.0

TABLE II AN EXCERPT OF PUBLICATIONS (GROUND TRUTH)



ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_1$	2013	NADEEF	ACM SIGMOD	QCRI	174.0
$t_2$	2013	NADEEF	SIGMOD Conf.	QCRI, HBKU	1740
$t_3$	2013	NADEEF	SIGMOD	QCRI HBKU	174.0
$t_4$	2013	KuaFu	ICDE 2013	Microsoft	15.0
$t_5$	2013	TsingNUS	SIGMOD'13	Tsinghua	13.0
$t_6$	2013	TsingNUS	SIGMOD'13	THU	13.0
$t_7$	2014	SeeDB	VLDB	Stanford Univ.	N.A.
$t_8$	2014	SeeDB	Very Large Data Bases	Stanford	55.0
$t_9$	2015	Elaps	ICDE	NUS	42.0
$t_{10}$	2015	Elaps	IEEE ICDE Conf. 2015	CS@NUS	44.0

#### TABLE I

#### AN EXCERPT OF PUBLICATIONS (DIRTY)

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_{123}$	2013	NADEEF	SIGMOD	QCRI	174.0
$t_4$	2013	KuaFu	ICDE	Microsoft	15.0
$t_{56}$	2013	TsingNUS	SIGMOD	Tsinghua	13.0
$t_{78}$	2014	SeeDB	VLDB	Stantford Univ.	55.0
$t_{910}$	2015	Elaps	ICDE	NUS	43.0

#### TABLE II

#### AN EXCERPT OF PUBLICATIONS (GROUND TRUTH)



(b) A Correct Pie Chart





0 450 900 1350 1800

(a) An Incorrect Bar Chart

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_1$	2013	NADEEF	ACM SIGMOD	QCRI	174.0
$t_2$	2013	NADEEF	SIGMOD Conf.	QCRI, HBKU	1740
$t_3$	2013	NADEEF	SIGMOD	QCRI HBKU	174.0
$t_4$	2013	KuaFu	ICDE 2013	Microsoft	15.0
$t_5$	2013	TsingNUS	SIGMOD'13	Tsinghua	13.0
$t_6$	2013	TsingNUS	SIGMOD'13	THU	13.0
$t_7$	2014	SeeDB	VLDB	Stanford Univ.	N.A.
$t_8$	2014	SeeDB	Very Large Data Bases	Stanford	55.0
$t_9$	2015	Elaps	ICDE	NUS	42.0
$t_{10}$	2015	Elaps	IEEE ICDE Conf. 2015	CS@NUS	44.0

#### TABLE I

AN EXCERPT OF PUBLICATIONS (DIRTY)

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_{123}$	2013	NADEEF	SIGMOD	QCRI	174.0
$t_4$	2013	KuaFu	ICDE	Microsoft	15.0
$t_{56}$	2013	TsingNUS	SIGMOD	Tsinghua	13.0
$t_{78}$	2014	SeeDB	VLDB	Stantford Univ.	55.0
$t_{910}$	2015	Elaps	ICDE	NUS	43.0

TABLE II







#### **Data errors**

duplicates synonyms outliers missing values

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_1$	2013	NADEEF	ACM SIGMOD	QCRI	174.0
$t_2$	2013	NADEEF	SIGMOD Conf.	QCRI, HBKU	1740
$t_3$	2013	NADEEF	SIGMOD	QCRI HBKU	174.0
$t_4$	2013	KuaFu	ICDE 2013	Microsoft	15.0
$t_5$	2013	TsingNUS	SIGMOD'13	Tsinghua	13.0
$t_6$	2013	TsingNUS	SIGMOD'13	THU	13.0
$t_7$	2014	SeeDB	VLDB	Stanford Univ.	N.A.
$t_8$	2014	SeeDB	Very Large Data Bases	Stanford	55.0
$t_9$	2015	Elaps	ICDE	NUS	42.0
$t_{10}$	2015	Elaps	IEEE ICDE Conf. 2015	CS@NUS	44.0

#### TABLE I

AN EXCERPT OF PUBLICATIONS (DIRTY)

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_{123}$	2013	NADEEF	SIGMOD	QCRI	174.0
$t_4$	2013	KuaFu	ICDE	Microsoft	15.0
$t_{56}$	2013	TsingNUS	SIGMOD	Tsinghua	13.0
$t_{78}$	2014	SeeDB	VLDB	Stantford Univ.	55.0
$t_{910}$	2015	Elaps	ICDE	NUS	43.0

TABLE II AN EXCERPT OF PUBLICATIONS (GROUND TRUTH)





8

#### **Data errors**

*duplicates synonyms outliers missing values* 

#### **Distance between visualizations**

Earth Mover Distance (EMD)

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_1$	2013	NADEEF	ACM SIGMOD	QCRI	174.0
$t_2$	2013	NADEEF	SIGMOD Conf.	QCRI, HBKU	1740
$t_3$	2013	NADEEF	SIGMOD	QCRI HBKU	174.0
$t_4$	2013	KuaFu	ICDE 2013	Microsoft	15.0
$t_5$	2013	TsingNUS	SIGMOD'13	Tsinghua	13.0
$t_6$	2013	TsingNUS	SIGMOD'13	THU	13.0
$t_7$	2014	SeeDB	VLDB	Stanford Univ.	N.A.
$t_8$	2014	SeeDB	Very Large Data Bases	Stanford	55.0
$t_9$	2015	Elaps	ICDE	NUS	42.0
$t_{10}$	2015	Elaps	IEEE ICDE Conf. 2015	CS@NUS	44.0

#### TABLE I

AN EXCERPT OF PUBLICATIONS (DIRTY)

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_{123}$	2013	NADEEF	SIGMOD	QCRI	174.0
$t_4$	2013	KuaFu	ICDE	Microsoft	15.0
$t_{56}$	2013	TsingNUS	SIGMOD	Tsinghua	13.0
$t_{78}$	2014	SeeDB	VLDB	Stantford Univ.	55.0
$t_{910}$	2015	Elaps	ICDE	NUS	43.0

TABLE II AN EXCERPT OF PUBLICATIONS (GROUND TRUTH)





#### **Data errors**

*duplicates synonyms outliers missing values* 

#### **Distance between visualizations**

Earth Mover Distance (EMD)

#### Error repair graph (ERG)

Compile heterogenous errors into one graph

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_1$	2013	NADEEF	ACM SIGMOD	QCRI	174.0
$t_2$	2013	NADEEF	SIGMOD Conf.	QCRI, HBKU	1740
$t_3$	2013	NADEEF	SIGMOD	QCRI HBKU	174.0
$t_4$	2013	KuaFu	ICDE 2013	Microsoft	15.0
$t_5$	2013	TsingNUS	SIGMOD'13	Tsinghua	13.0
$t_6$	2013	TsingNUS	SIGMOD'13	THU	13.0
$t_7$	2014	SeeDB	VLDB	Stanford Univ.	N.A.
$t_8$	2014	SeeDB	Very Large Data Bases	Stanford	55.0
$t_9$	2015	Elaps	ICDE	NUS	42.0
$t_{10}$	2015	Elaps	IEEE ICDE Conf. 2015	CS@NUS	44.0

#### TABLE I

AN EXCERPT OF PUBLICATIONS (DIRTY)

ld	Year	Title (abbr.)	Venue	Affiliation	Citations
$t_{123}$	2013	NADEEF	SIGMOD	QCRI	174.0
$t_4$	2013	KuaFu	ICDE	Microsoft	15.0
$t_{56}$	2013	TsingNUS	SIGMOD	Tsinghua	13.0
$t_{78}$	2014	SeeDB	VLDB	Stantford Univ.	55.0
$t_{910}$	2015	Elaps	ICDE	NUS	43.0

TABLE II AN EXCERPT OF PUBLICATIONS (GROUND TRUTH)







#### **Data errors**

duplicates synonyms outliers missing values

#### **Distance between visualizations**

Earth Mover Distance (EMD)

#### Error repair graph (ERG)

Compile heterogenous errors into one graph

#### **Interactive cleaning for progressive visualization (ICPV)**

*Given a bad visualization V, it is to obtain a "cleaned" visualization V whose distance is far from V, under a given (small) budget w.r.t. human cost.* 



**Interactive Cleaning** 



(b) Example of Label Vertex t1















### **Collaborative Data Cleaning (GUI)**



#### Welcome to the Collaborative Data Cleaning!

mashaal2 -

#### Share this dataset with other collaborators:

#### http://localhost:3000/dataset/iCLkGtYybGSi5emeZ

	Number of times pregnant 💌	Plasma glucose concentration 💌	Diastolic blood pressure 💌	Tric	eps skin fold thickness 💌	2-Hour serum insulin 💌	Body mass index 💌	Diabetes pedigree fu
1	6	148	72	35		0	33.6	0.627
2	1	85	66	29		0	26.6	0.351
3	8	183	64	0		0	23.3	0.672
4	1	89	66	23		94	28.1	0.167
5	0	137	40	35		168	43.1	2.288
6	5	116	74	0		0	25.6	0.201
7	3	78	50	32		88	31	0.248
8	10	115	0	0		0	35.3	0.134
9	2	197	70	45	Flag as dirty		30.5	0.158
10	8	125	96	0	Values from other collaborators		0	0.232
11	4	110	92	0			37.6	0.191
12	10	168	74	0	55		38	0.537
13	10	139	80	0		0	27.1	1.441
14	1	189	60	23		846	30.1	0.398
15	5	166	72	19		175	25.8	0.587
16	7	100	0	0		0	30	0.484
17	0	118	84	47		230	45.8	0.551
18	7	107	74	0		0	29.6	0.254
19	1	103	30	38		83	43.3	0.183
20	1	115	70	30		96	34.6	0.529
21	3	126	88	41		235	39.3	0 704





missing data may influence analysts'

perceptions of data quality and their confidence in their conclusions?







missing data may influence analysts'

perceptions of data quality and their confidence in their conclusions?



#### **Imputation Methods**







missing data may influence analysts'

perceptions of data quality and their confidence in their conclusions?



#### **Imputation Methods**



#### **Encode Imputed Values**



Where is My Data? Evaluating Visualizations with Missing Data. Hayeong Song et al., InfoVis 2018.

# Missing Data



### Questions

- Is the overall rate of change larger in the first or second half-hour?
- How confident are you in your response?
- How complete is this data?
- How reliable is this data?

H1– Perceived data quality will degrade as the amount of missing data increases.

H2– Highlighting will achieve higher perceived data quality than downplaying and information removal

H3– Linear interpolation will lead to higher perceived
confidence and data quality than marginal means or zero-filling
H4– Imputed values will lead to higher perceived data quality

than removed values.



A Framework for Externalizing Implicit Error Using Visualization, Nina McCurday et al., InfoVis 2018.

Measurement error: the difference between a measured value and the true value as it exists

**Implicit error:** measure error that is inherent to a dataset, assumed to be present and prevalent, but not explicitly defined or accounted for

- exists as tacit knowledge in the minds of experts
- is rarely quantifiable
- is accounted for subjectively during expert interpretation of the data

Measurement error: the difference between a measured value and the true value as it exists

**Implicit error:** measure error that is inherent to a dataset, assumed to be present and prevalent, but not explicitly defined or accounted for

- exists as tacit knowledge in the minds of experts
- is rarely quantifiable

• is accounted for subjectively during expert interpretation of the data processing

detection recording collection a processing reporting distributed heterogeneous data generation pipelines

Measurement error: the difference between a measured value and the true value as it exists

**Implicit error:** measure error that is inherent to a dataset, assumed to be present and prevalent, but not explicitly defined or accounted for **detection detection detection** 

- exists as tacit knowledge in the minds of experts
- is rarely quantifiable
- is accounted for subjectively during expert interpretation of the data



recording

collection

distributed heterogeneous data generation pipelines



17

Amazon Mechanical Turk People who go to saunas are more likely to know that Mike Stonebraker is not a character in "The Simpsons"

#### *Towards Sustainable Insights,* **Carsten Binning et al.**, CIDR 2017.

Amazon Mechanical Turk People who go to saunas are more likely to know that Mike Stonebraker is not a character in "The Simpsons"

#### SeeDB

Discovery rate: three times higher than for manual exploration tools such as Tableau or Vizdom



Amazon Mechanical Turk People who go to saunas are more likely to know that Mike Stonebraker is not a character in "The Simpsons"

#### SeeDB

Discovery rate: three times higher than for manual exploration tools such as Tableau or Vizdom

### Survey Data

Reference views: 9,996

Target views: 2,078,608

Recommended views: 708,109



(a) Interesting visualization

(b) Uninteresting visualization

*Towards Sustainable Insights,* **Carsten Binning et al.**, CIDR 2017.

Amazon Mechanical Turk People who go to saunas are more likely to know that Mike Stonebraker is not a character in "The Simpsons"

#### SeeDB

Discovery rate: three times higher than for manual exploration tools such as Tableau or Vizdom

### Survey Data

Reference views: 9,996

Target views: 2,078,608

Recommended views: 708,109



(a) Interesting visualization (b) Uninter

#### (b) Uninteresting visualization

#### QUDE

Quantify the Uncertainty in Data Exploration

– False Discovery Rate –

*Towards Sustainable Insights,* **Carsten Binning et al.**, CIDR 2017.

### **Collaborative Data Visualization**



https://vimeo.com/yalongyang

**Enhances** the traditional visualization by bringing together **many experts** so that each can contribute toward the common goal of the **understanding** of the object, phenomenon, or data



**Synchronous** 



### Northstar: An Interactive Data Science System



*Northstar: An Interactive Data Science System*, **Tim Kraska**. VLDB 2018.

### Northstar: An Interactive Data Science System



- **Vizdom**: a visual data exploration environment pen and touch interface
- **IDEA**: an intelligent cache and string approximation engine
- **QUDE**: warn about common mistakes and problems
- **Alphine Meadow**: automatically orchestrate a ML pipeline (i.e., plan)

*Northstar: An Interactive Data Science System*, **Tim Kraska**. VLDB 2018.

# **Democratizing Data Visualization**



#### Immersion exploration inside-feeling

### Interaction

- gesture + voice
- real-time

### Intelligence reactive

proactive

### **IATK: Tableau for Virtual Reality**

Inspector     Inspector	pector							
👕 🗹 [IATK] New Brushing And Linking 🗌 Static 🔻								
Tag Untagged		E Laye	r Defa	ult		\$		
▼ 🙏 Transform 🔯								
Position	X 0	Y	0	Z	0			
Rotation	X 0	Y	0	Z	0			
Scale	X 1	Y	1	Z	1			
🔻 🕼 🗹 Brushing And Linking (Script) 🛛 🔯 🖏								
Script	💽 Brush		0					
Compute Shader	⊚MyCo		0					
My Render Material	RenderingMaterial					0		
Visualisations Materials								
Brushing Visualisation None (Visualisation)								
Brushed Visualisations								
Size	0							
Brushed Linking Visualisations								
Size	0							
Show Brush								
Brush Color								
Input 1	None (Transform)							
Input 2	ransfor	m)			0			
Radius Sphere 0.213								
Brush Button Controlle								
BRUSH_TYPE	SPHERE \$							
Debug Object Texture	✓ SP	HERE				0		
▶ Brushed Data	BO	X						



#### brushed visualisation



**Brushed Linking Visualisation** 



IATK: An Immersive Analytics Toolkit. Maxime Cordeil et al., IEEE Conference on Virtual Reality and 3D User Interfaces (VR), 2019.

### **IATK: Tableau for Virtual Reality**



### **Querying Real-World Data using Augmented Reality**





ARQuery: Hallucinating Analytics over Real-World Data using Augmented Reality, **Codi Burley et al.**, CIDR 2019. 24

### Future Work

- 1. Visualization-driven data preparation/curation
- 2. A generic system to support efficient/approximate/progressive data visualization
- 3. (Asynchronous) Collaborative data visualization
- 4. Data Visualization Benchmarks
- 5. Immersive data analytics for abstract data
- 6. Deep learning for data visualization

# Further Readings

- Stephen Macke, "*Adaptive Sampling for Rapidly Matching Histograms*", VLDB 2018
- Doris Jung-Lin Lee et al., "The Case for a Visual Discovery Assistant: A Holistic Solution for Accelerating Visual Data Exploration", IEEE Data Bulletin 2018
- Kelly Mack et al., "Characterizing Scalability Issues in Spreadsheet Software using Online Forums", CHI 2018
- Leilani Battle et al., "Beagle: Automated Extraction and Interpretation of Visualizations from the Web", CHI 2018
- Wenbo Tao et al., "*Kyrix: Interactive Visual Data Exploration at Scale*", CIDR 2019
- Kevin Zeng Hu et al., "VizNet: Towards A Large-Scale Visualization Learning and Benchmarking Repository", CHI 2019
- Codi j. Burley et al., "ARQuery: Hallucinating Analytics over Real-World Data using Augmented Reality", CIDR 2019
- Meraj Ahmed Khan et al., "*Data Tweening: Incremental Visualization of Data Transforms*", PVLDB 2017
- Ciro Donalek et al., "Immersive and Collaborative Data Visualization Using Virtual Reality Platforms", IEEE International Conference on Big Data, 2014.
- Fernanda B. Viegas et al., "ManyEyes: A Site for Visualization at Internet Scale", IEEE Transactions on Visualization and Computer Graphics, 2007.