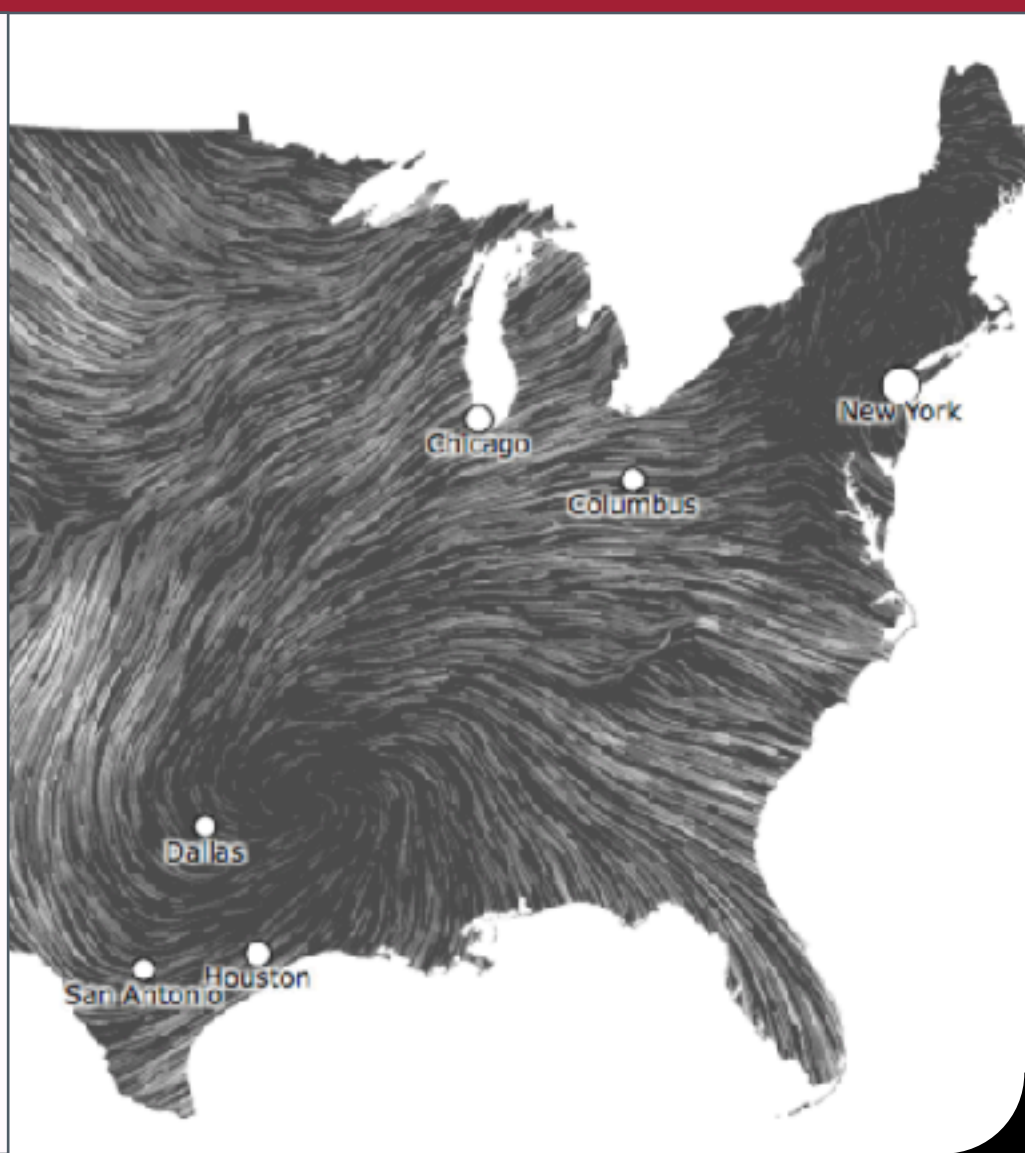
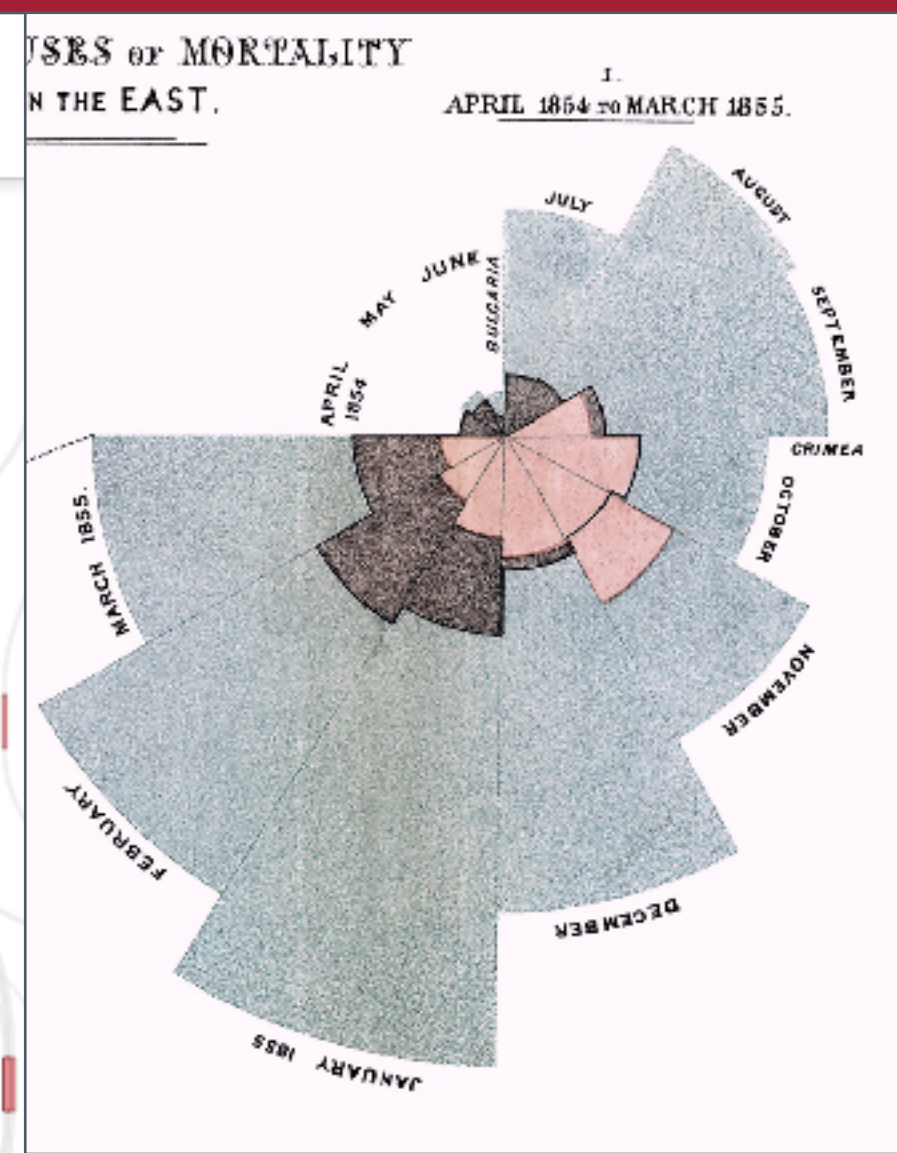
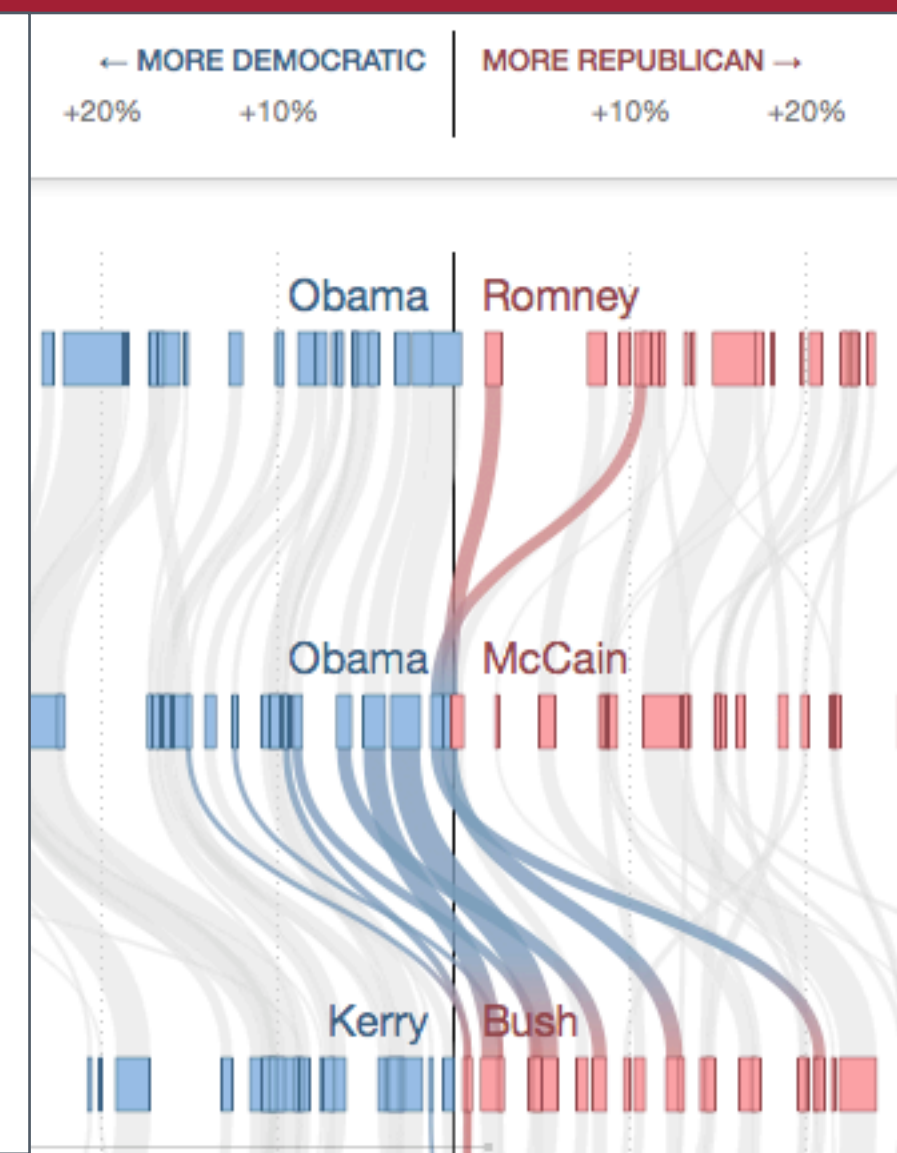
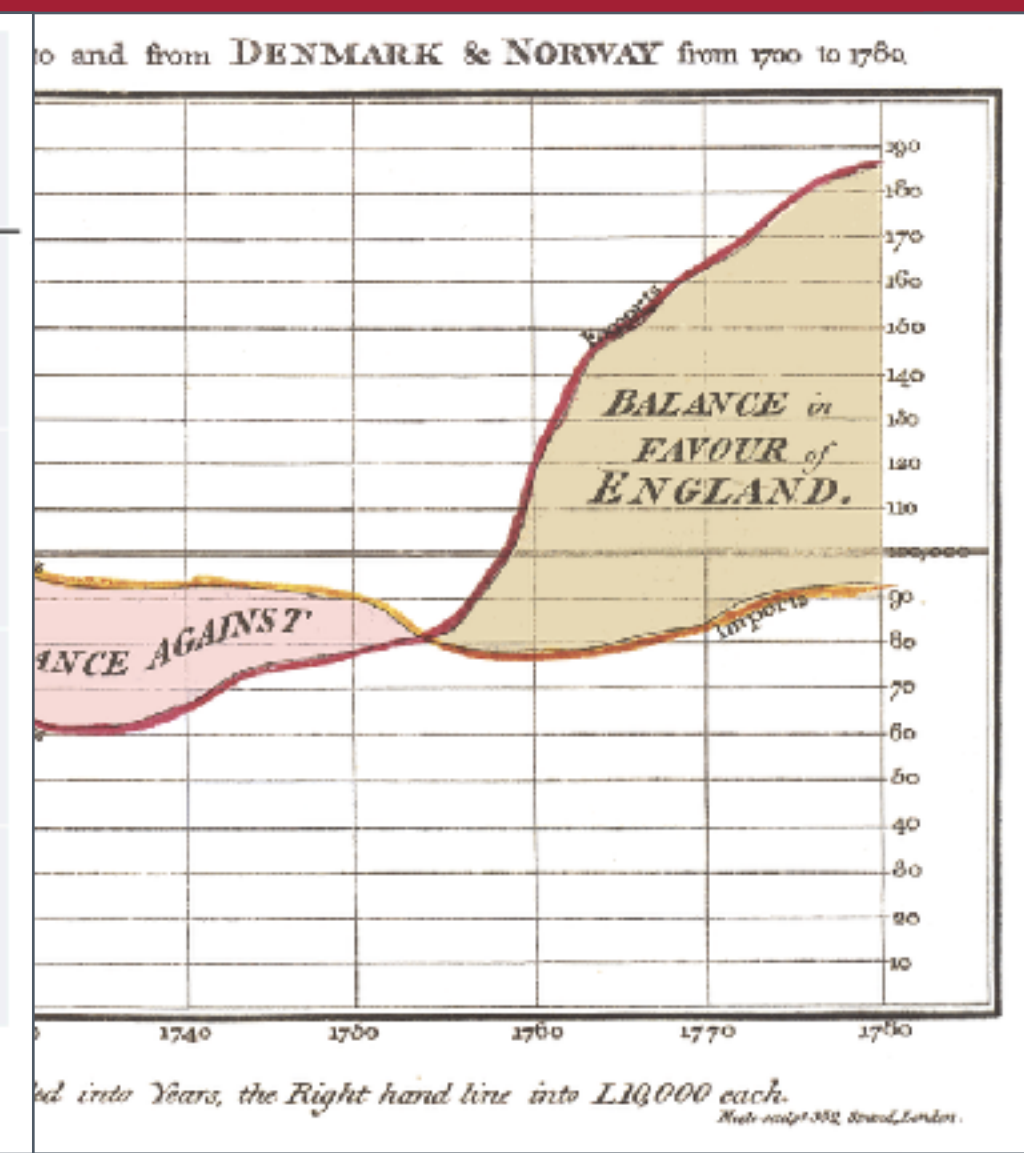
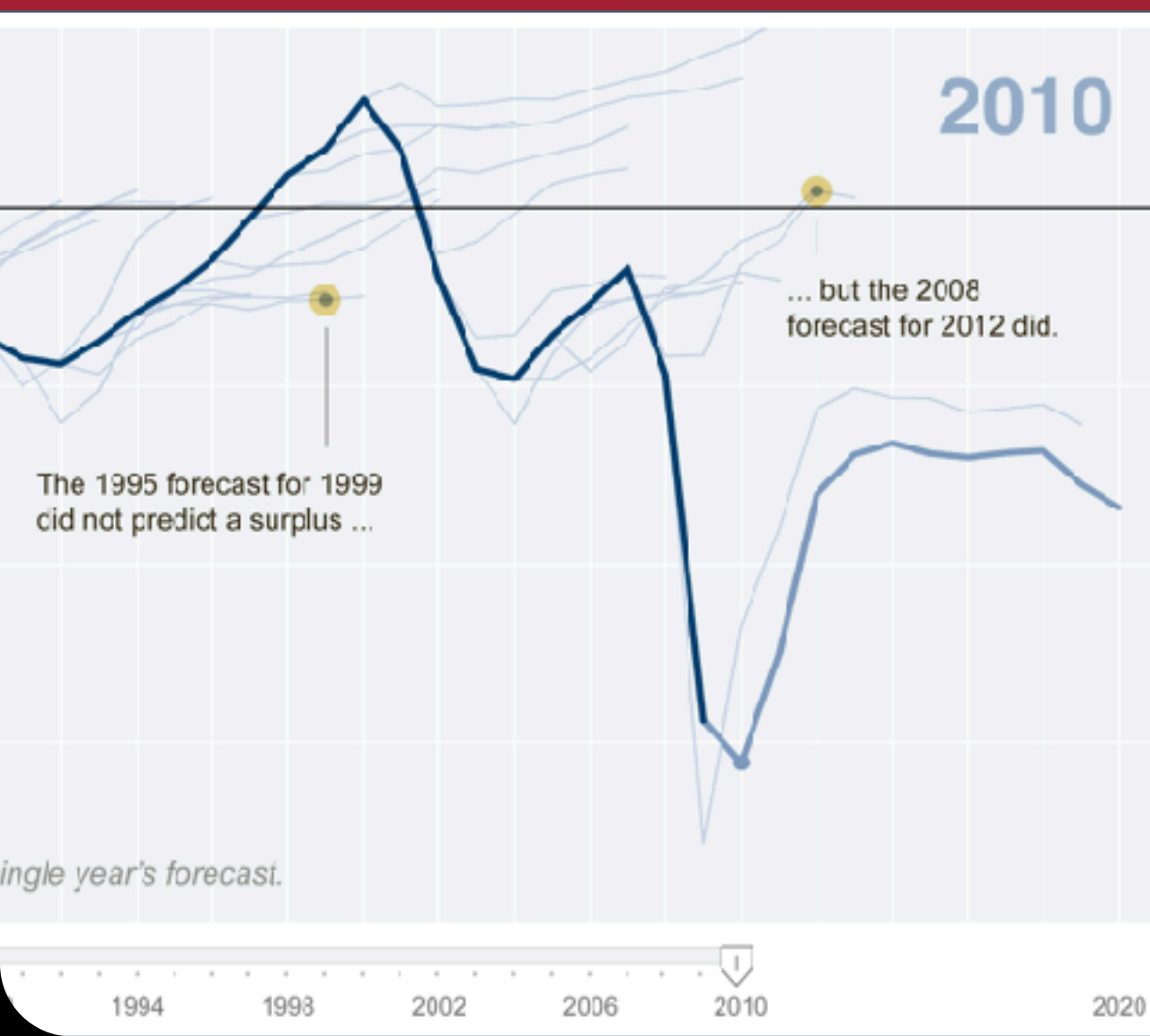


6.859: Interactive Data Visualization

Exploratory Data Analysis

Arvind Satyanarayan



Course Grading

Class Participation	5%
Reading Commentaries	5%
A0: Sketching Visualizations	2%
A1: Visualization Design	3%
A2: Exploratory Data Analysis	10%
A3: White/Black Hat Visualization	15%
A4: Interactive Narratives	20%
Final Project	40%
Proposal	
MVP + Presentations	
Poster Session + Final Deliverables	

5 slack days which can be used as you wish for assignments.

Slack days should cover minor illnesses, special occasions (including religious holidays).

Additional extensions only granted for serious issues with a written note of support from S3 or GradSupport @ OGE.

Share your work on Slack to inspire your classmates + receive design feedback!



Expressiveness

A set of facts is *expressible* in a visual language if the sentences (i.e. the visualizations) in the language express *all the facts in the set of data, and only the facts in the data.*

Data models give us a way of talking about this.

Effectiveness

A visualization is more *effective* than another if the information it conveys *is more readily perceived* than the information in the other visualization

Image models give us a way of talking about this.

[Mackinlay 1986]

Channels: Expressiveness Types and Effectiveness Ranks

➔ Magnitude Channels: O or Q attributes

Position on common scale



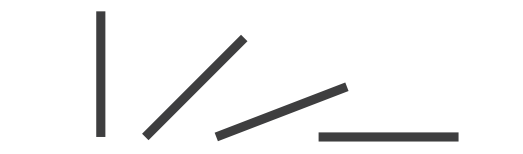
Position on unaligned scale



Length (1D size)



Tilt/angle



Area (2D size)



Depth (3D position)



Color luminance



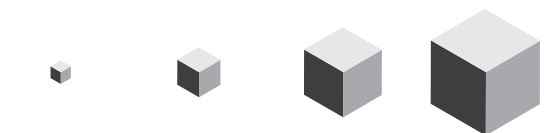
Color saturation



Curvature



Volume (3D size)



Same

Same

Same

Most Effectiveness Least

➔ Identity Channels: N attributes

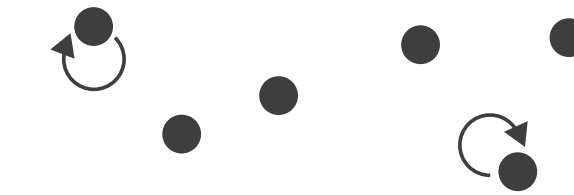
Spatial region



Color hue



Motion



Shape



Tamara Munzner, *Visualization Analysis and Design* (2014).

Visualization Critique

What is this a visualization of?

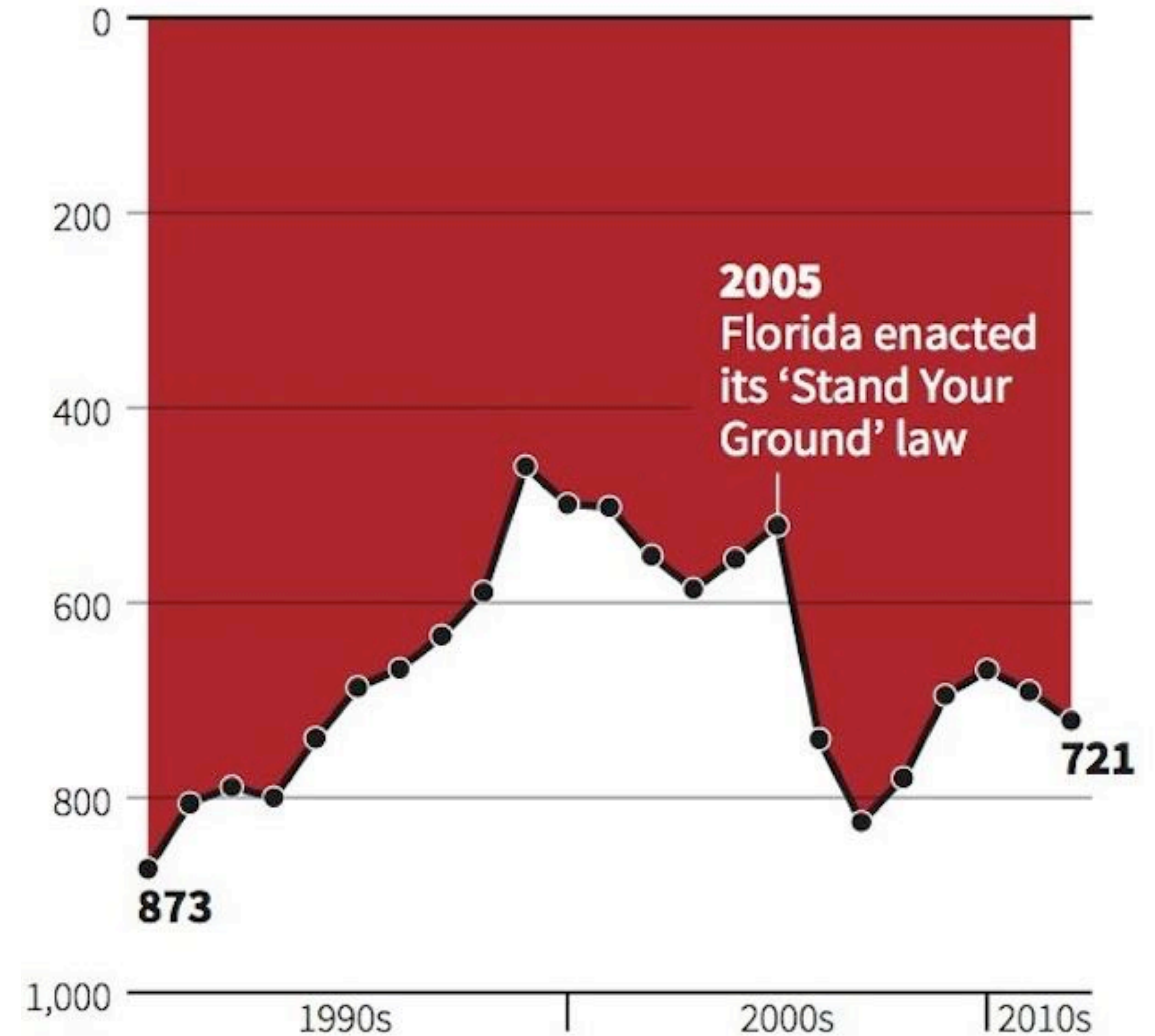
Which of the visual encoding techniques we've discussed this week are being used?
How effective or ineffective are they?

To help structure your critique:

- > "I like..."
- > "I wish..."
- > "What if...?"

Gun deaths in Florida

Number of murders committed using firearms



Source: Florida Department of Law Enforcement

How to Lie with Statistics: Stand Your Ground and Gun Deaths

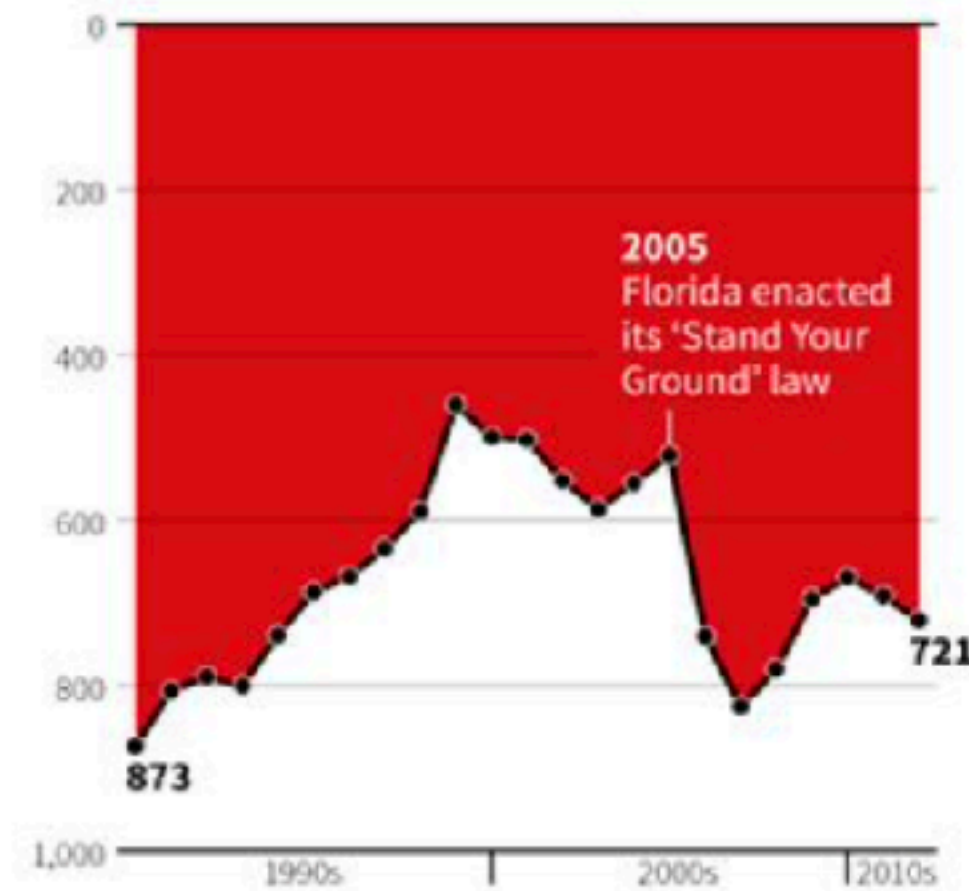
Lisa Wade, PhD on December 28, 2014

BEST OF 2014

At [Junk Charts](#), Kaiser Fung drew my attention to a graph released by Reuters. It is so deeply misleading that I loathe to expose your eyeballs to it. So, I offer you this:

Gun deaths in Florida

Number of murders committed using firearms



Source: Florida Department of Law Enforcement

C. Chan 16/02/2014

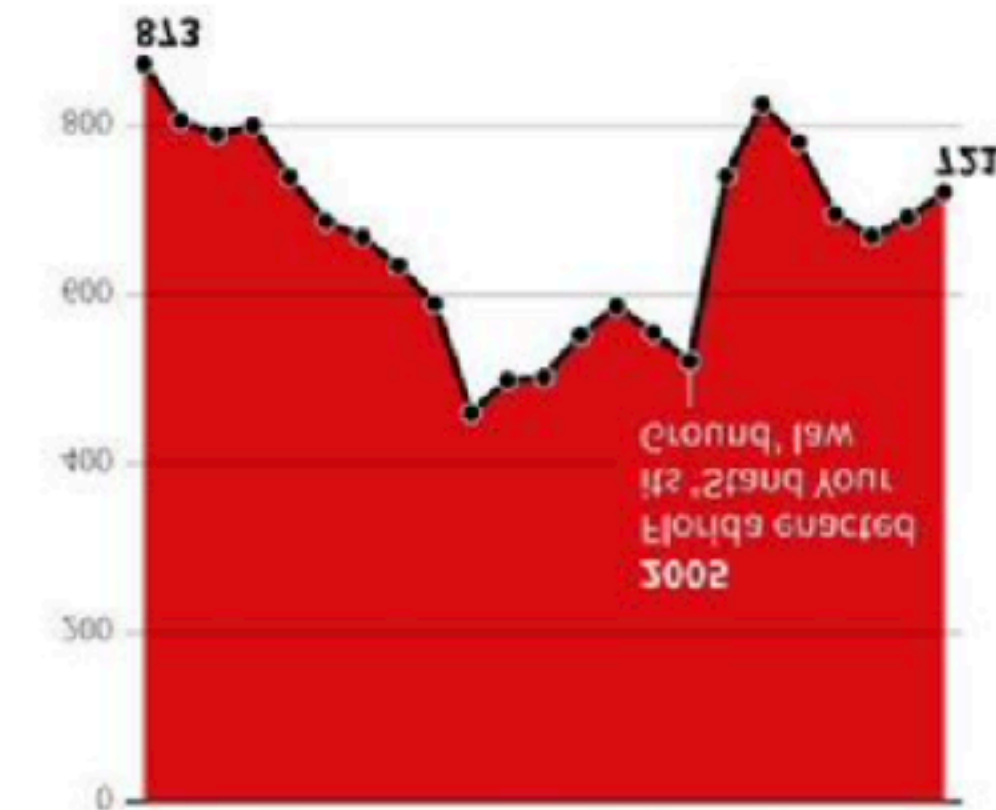
© REUTERS

C. Chan 16/02/2014

© REUTERS

Source: Florida Department of Law Enforcement

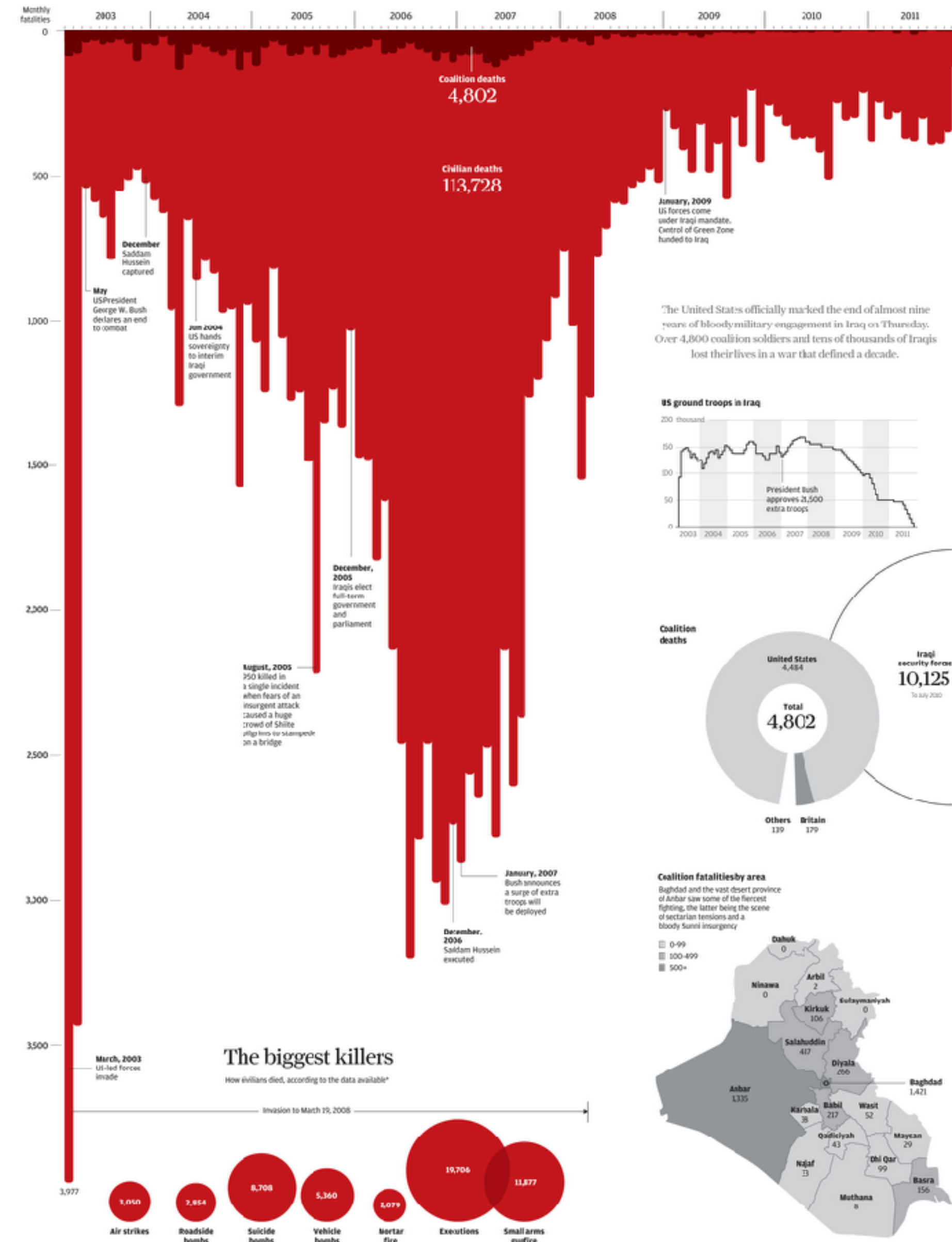
1,000 800 600 400 200 0



Number of murders committed using firearms

Gun deaths in Florida

Iraq's bloody toll



The United States officially marked the end of almost nine years of bloody military engagement in Iraq on Thursday. Over 4,800 coalition soldiers and tens of thousands of Iraqis lost their lives in a war that defined a decade.

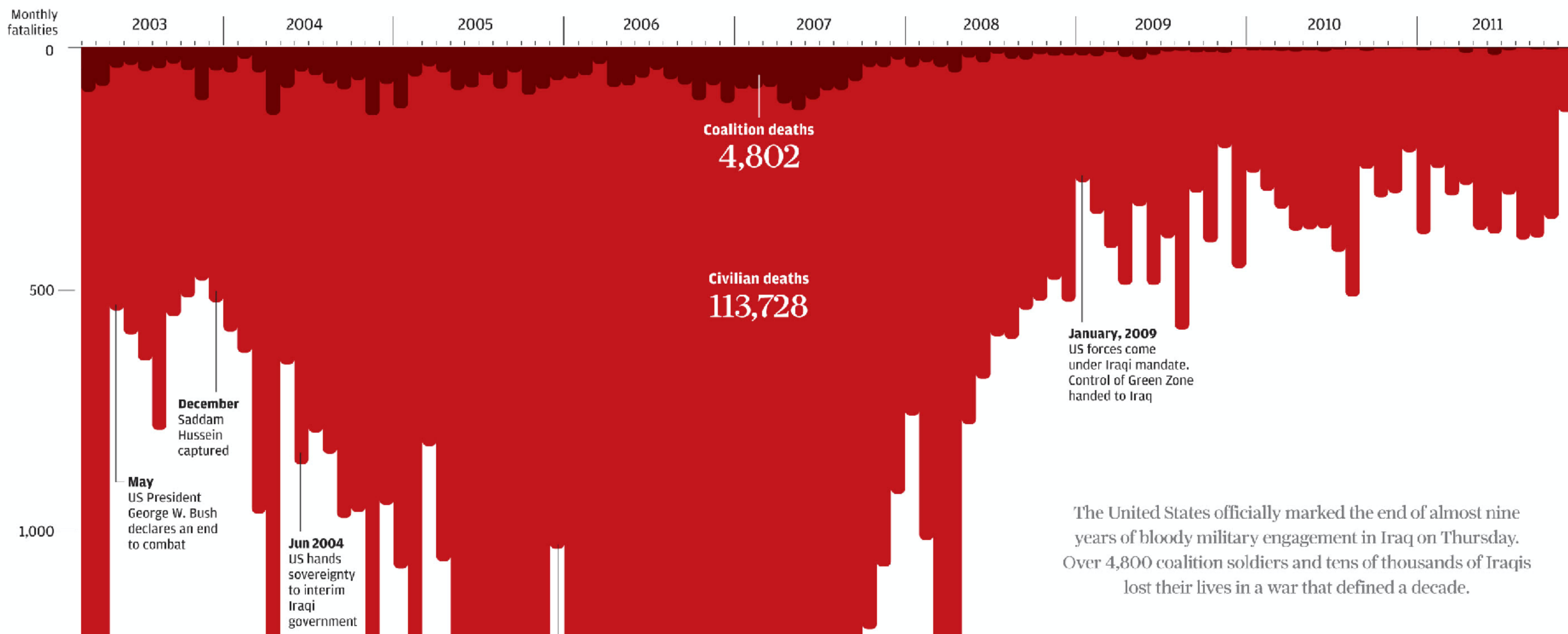
Baghdad and the vast desert province of Anbar saw some of the fiercest fighting, the latter being the scene of sectarian tensions and a bloody Sunni insurgency

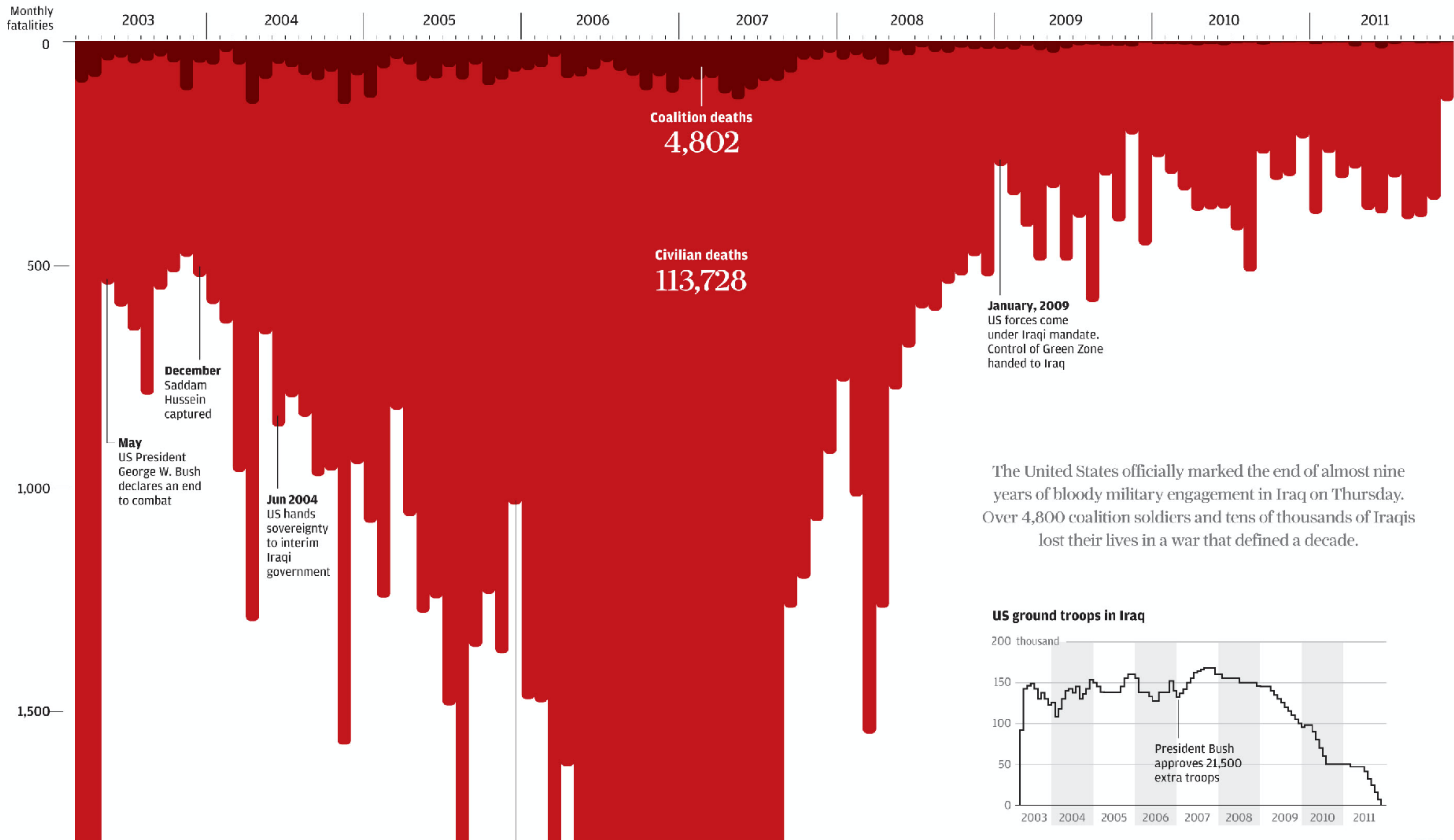
SCMP graphics: Simon Scott

*Deaths from unknown causes are not included. Causes of death accounting for less than 0.5 per cent of killings also not included.

Sources: Casualties, Iraq Body Count, New England Journal of Medicine, Global Security, Brookings Institution

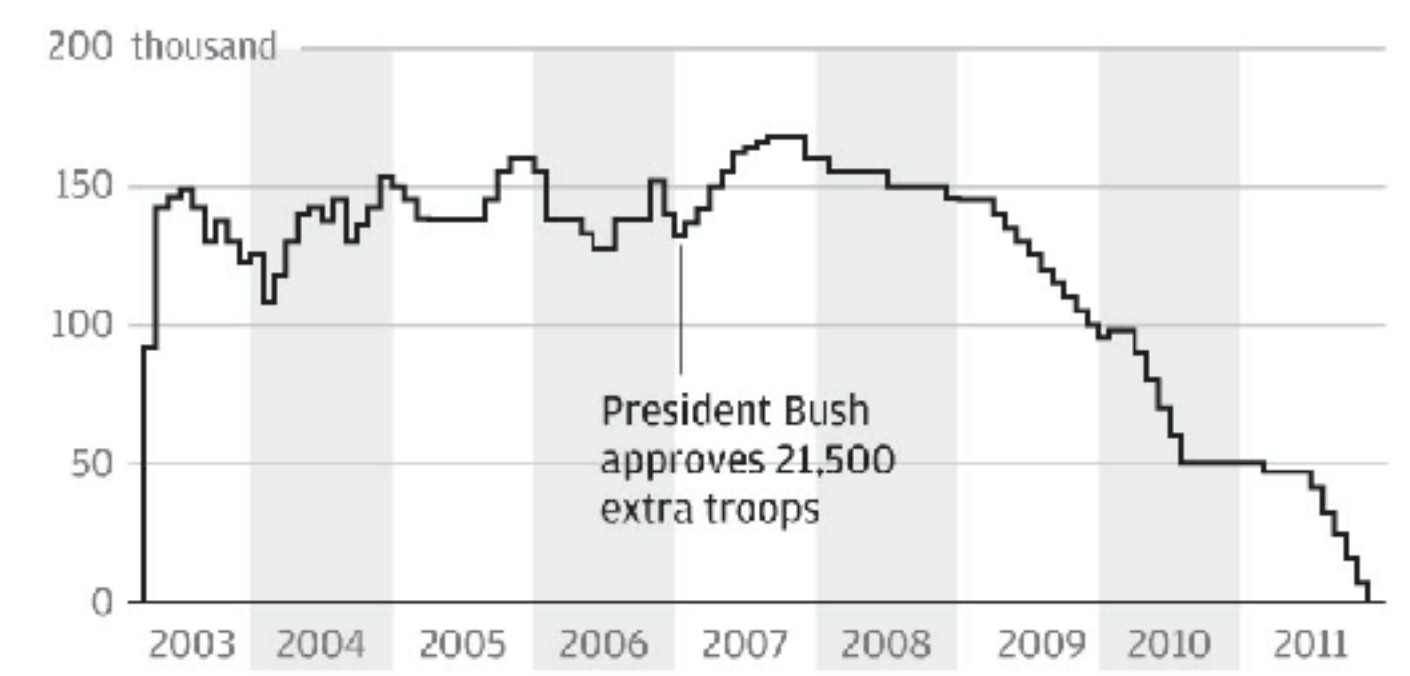
Iraq's bloody toll

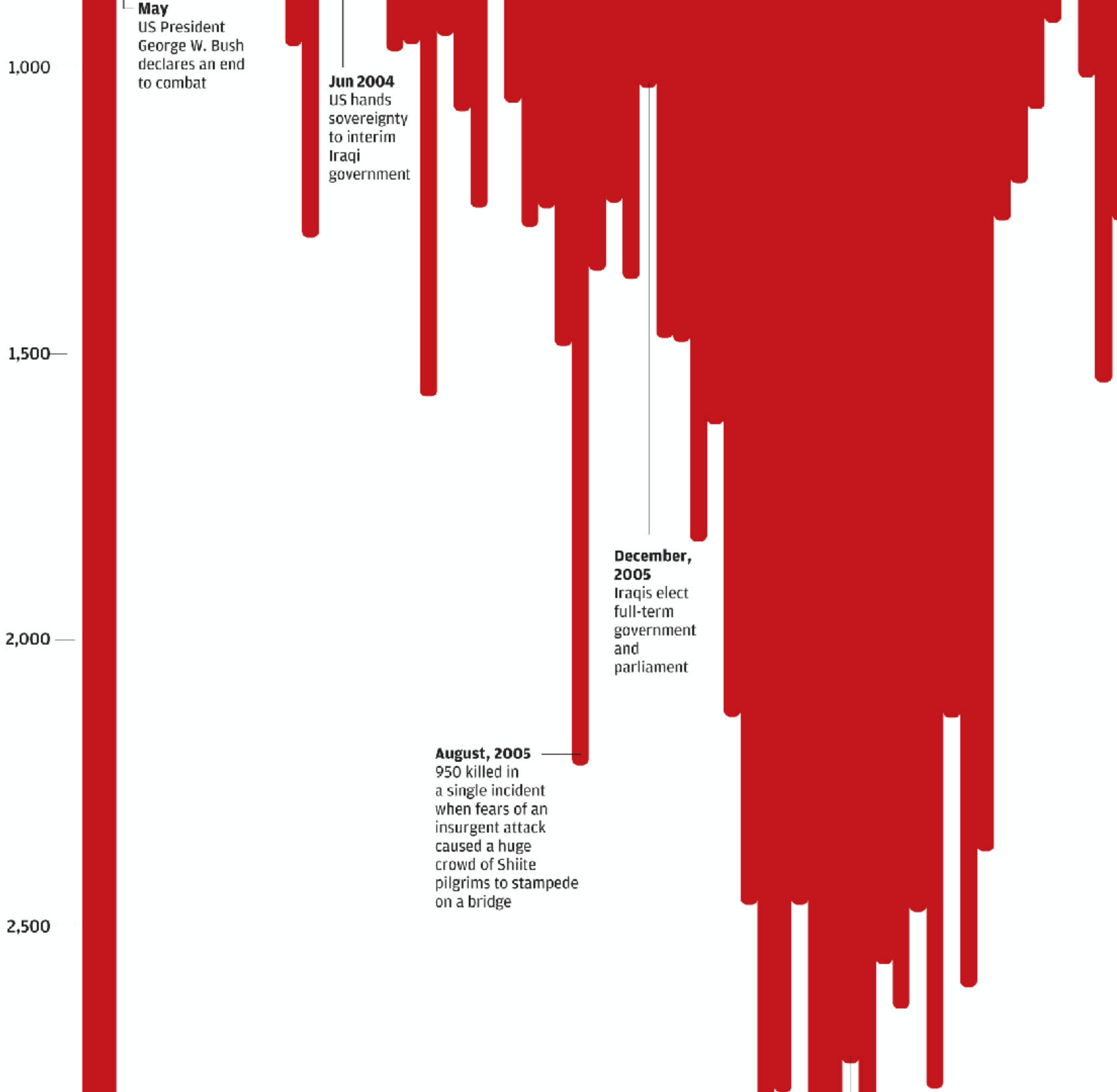




The United States officially marked the end of almost nine years of bloody military engagement in Iraq on Thursday. Over 4,800 coalition soldiers and tens of thousands of Iraqis lost their lives in a war that defined a decade.

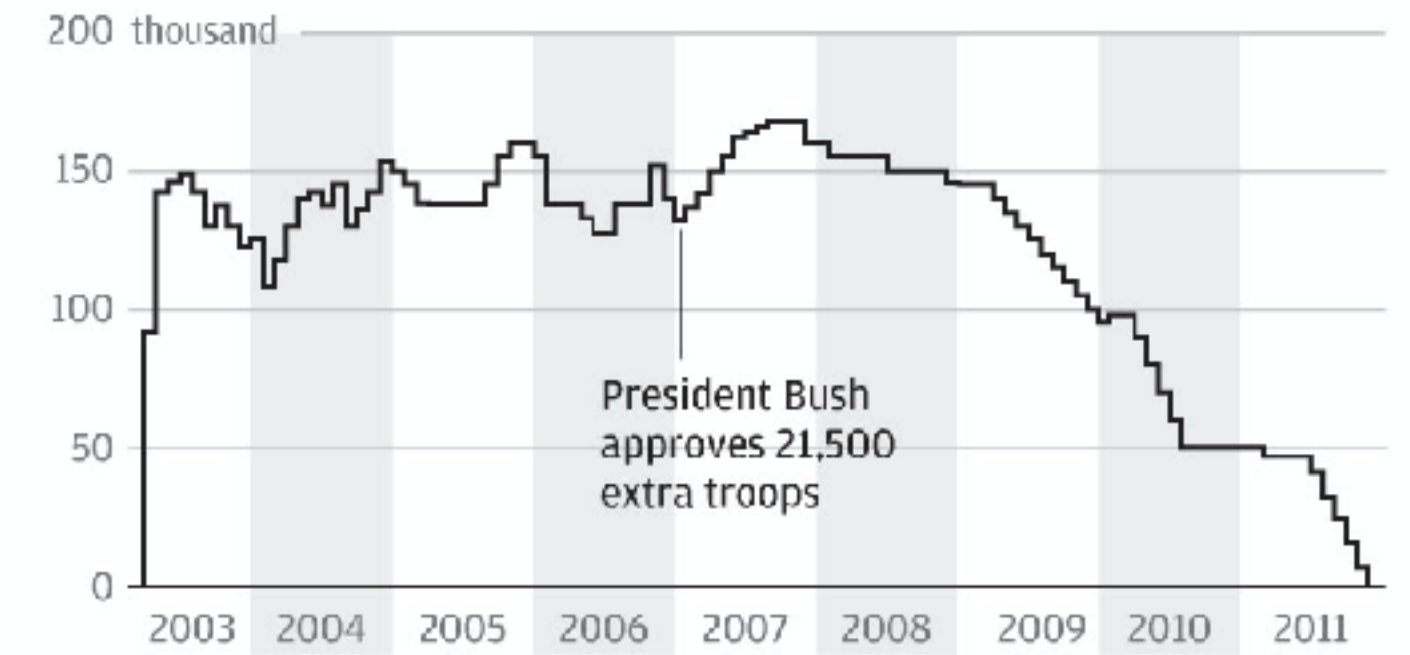
US ground troops in Iraq



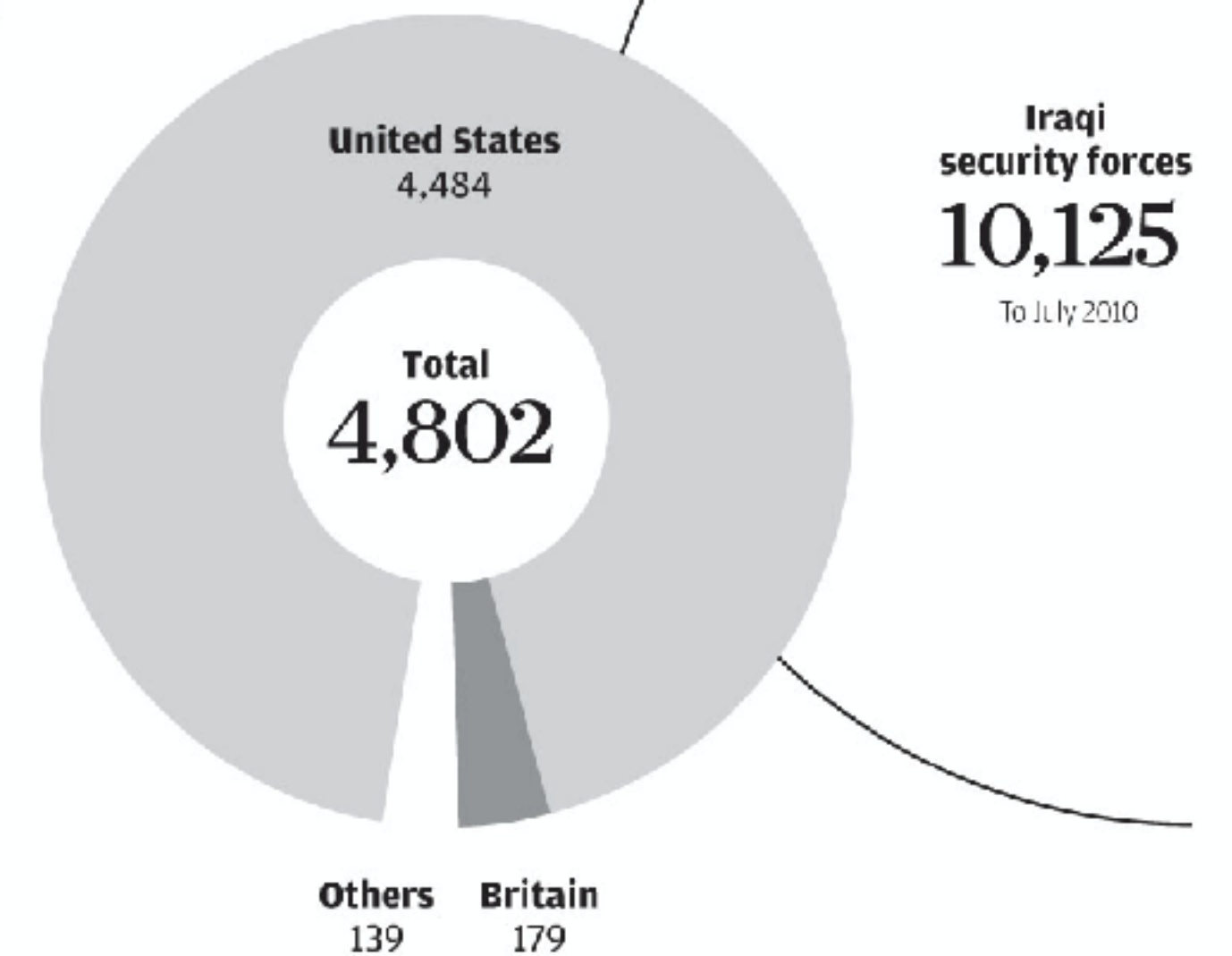


The United States officially marked the end of almost nine years of bloody military engagement in Iraq on Thursday. Over 4,800 coalition soldiers and tens of thousands of Iraqis lost their lives in a war that defined a decade.

US ground troops in Iraq



Coalition deaths

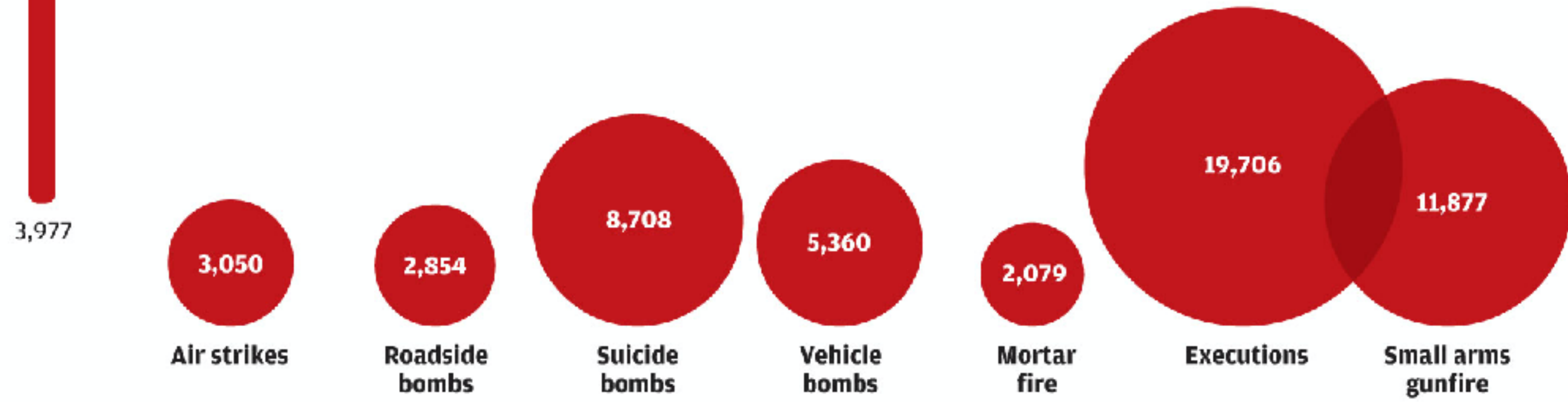




The biggest killers

How civilians died, according to the data available*

Invasion to March 19, 2008

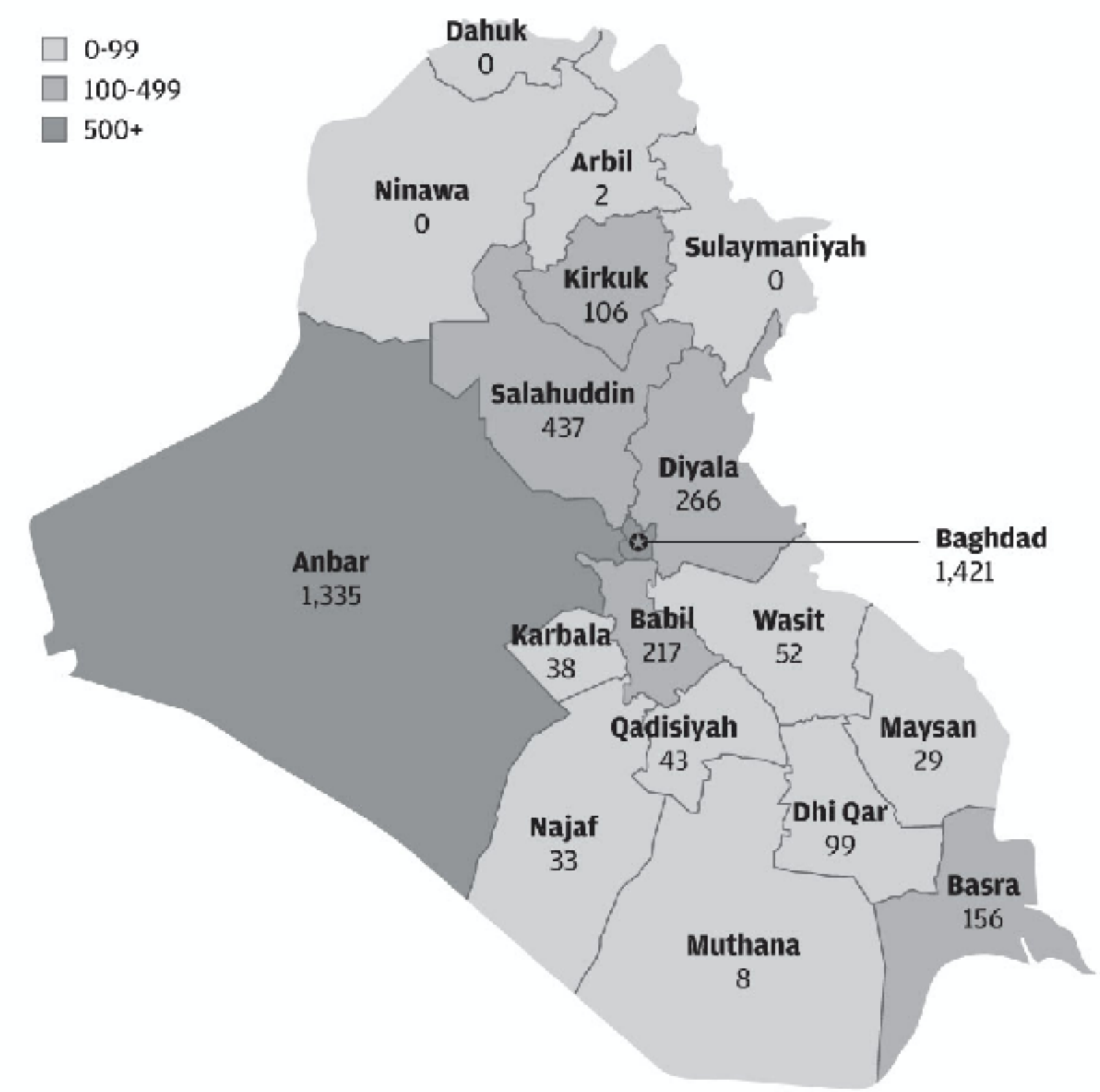


Others 139
Britain 179

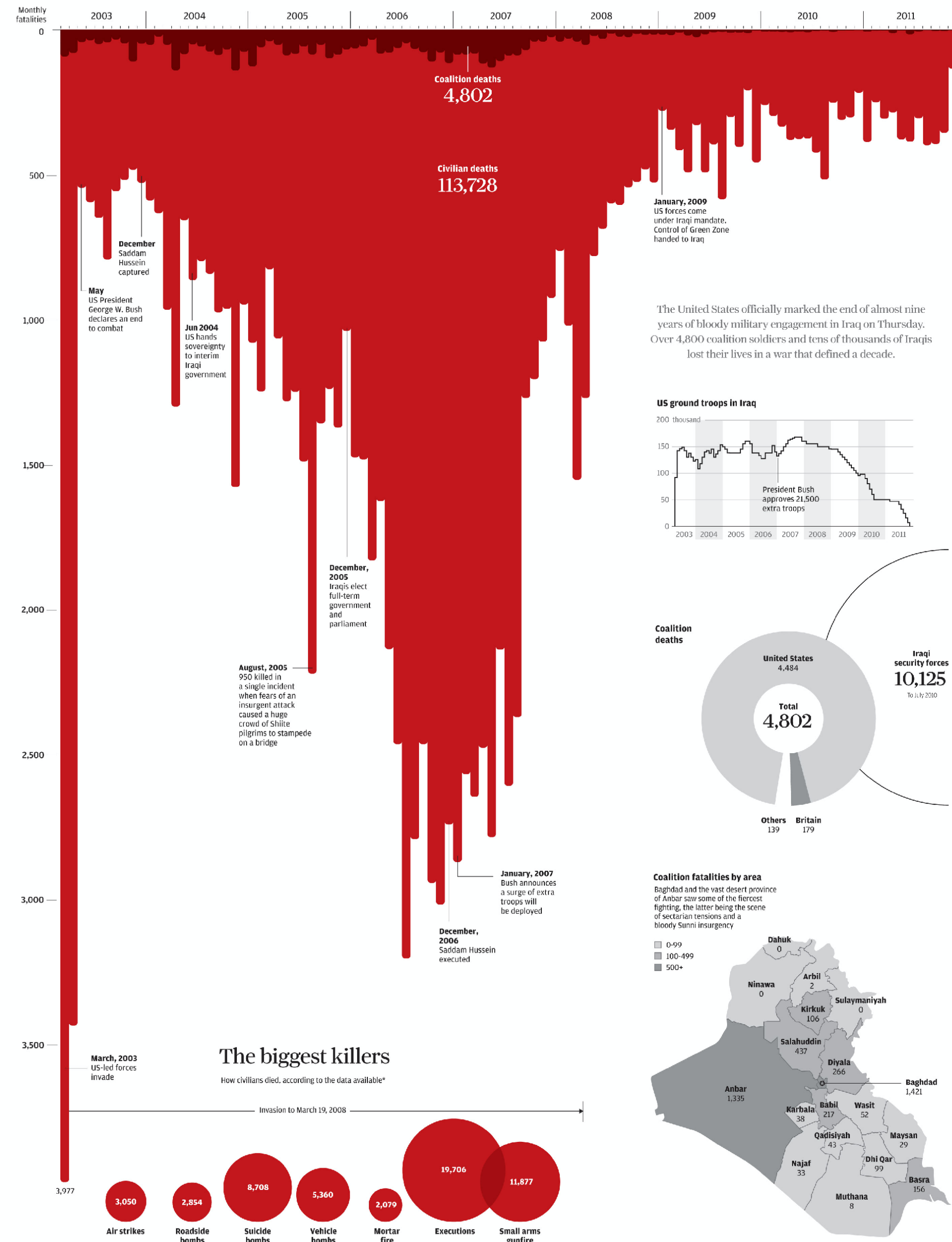
Coalition fatalities by area

Baghdad and the vast desert province of Anbar saw some of the fiercest fighting, the latter being the scene of sectarian tensions and a bloody Sunni insurgency

0-99
100-499
500+

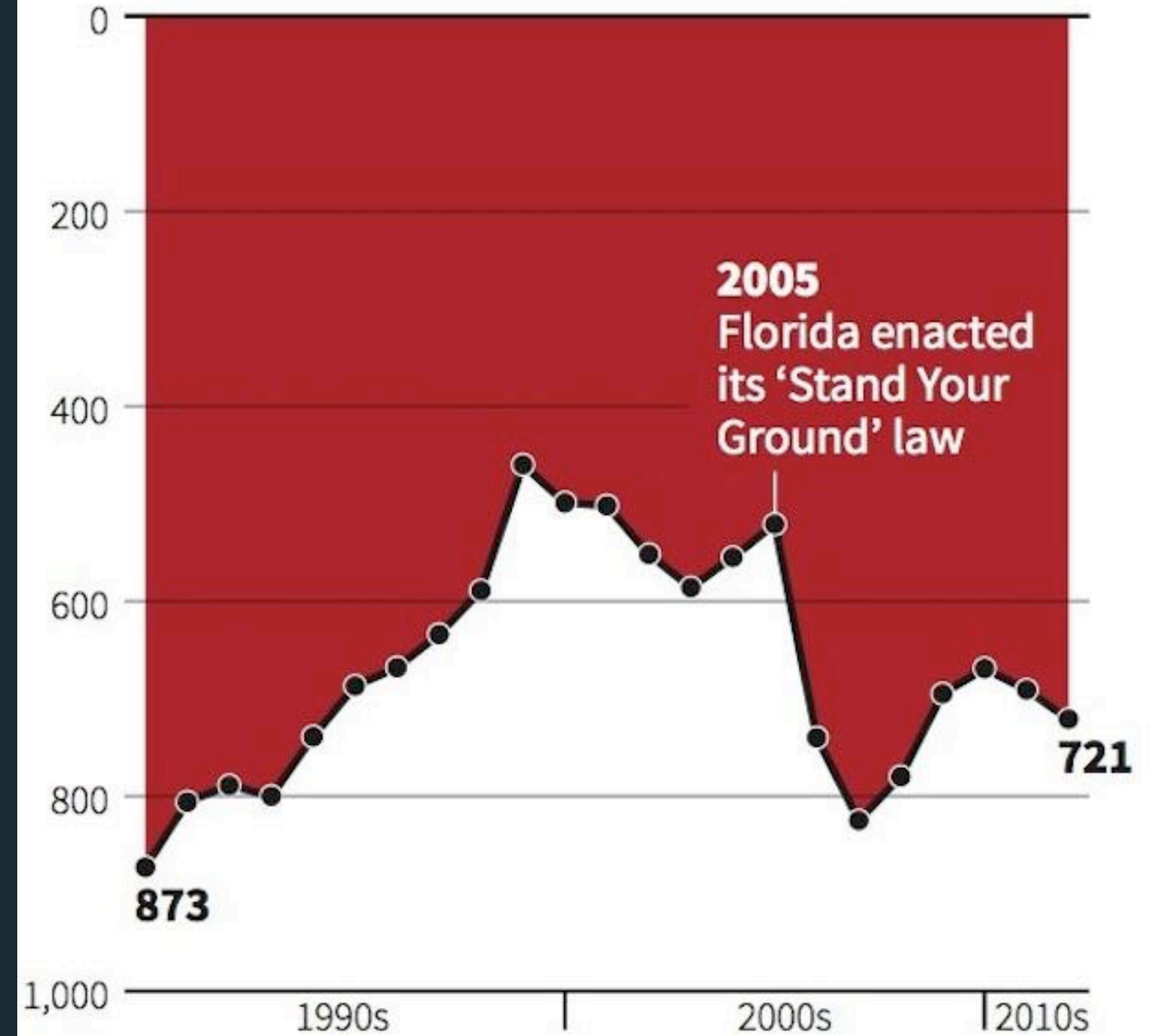


Iraq's bloody toll



Gun deaths in Florida

Number of murders committed using firearms

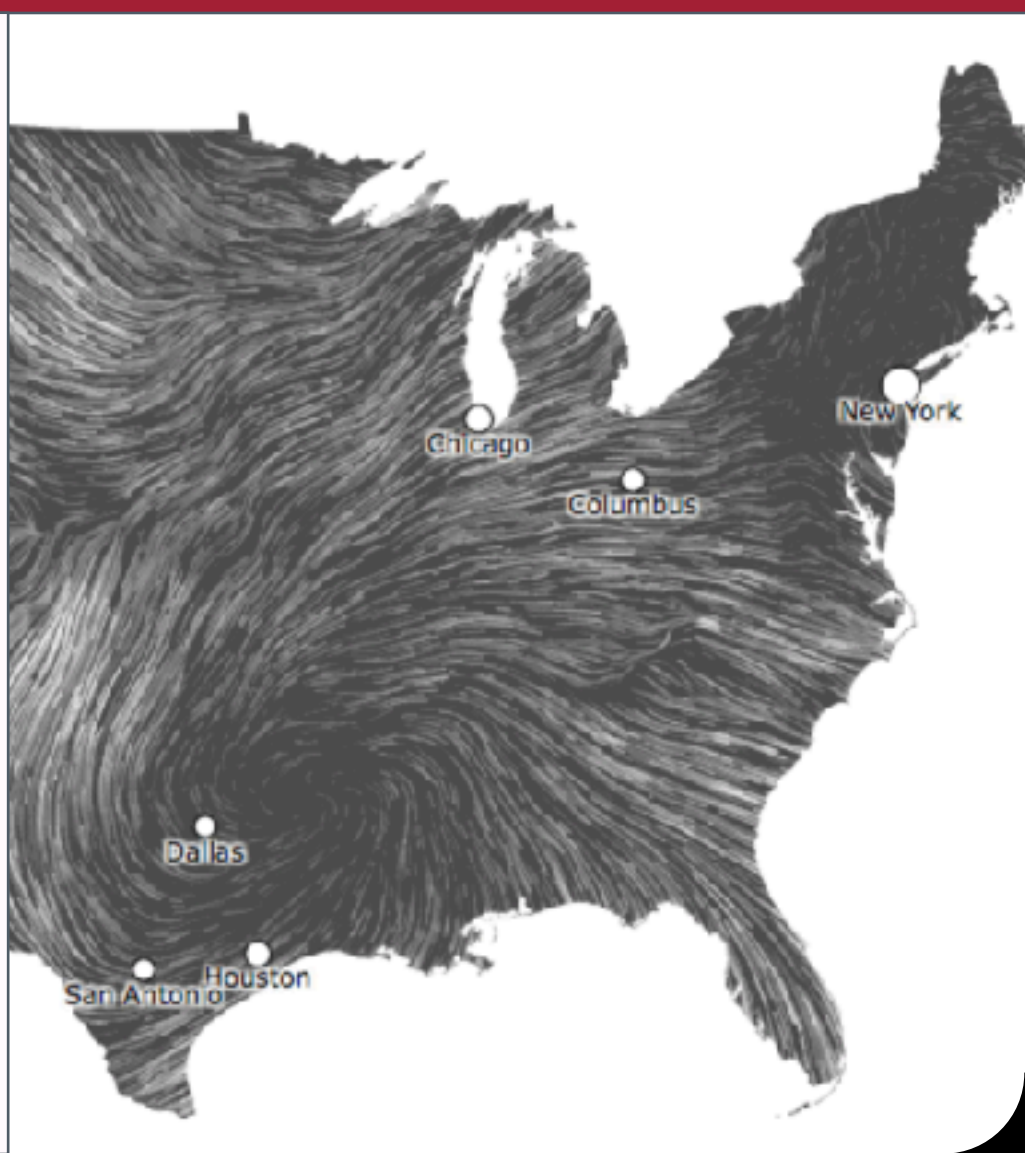
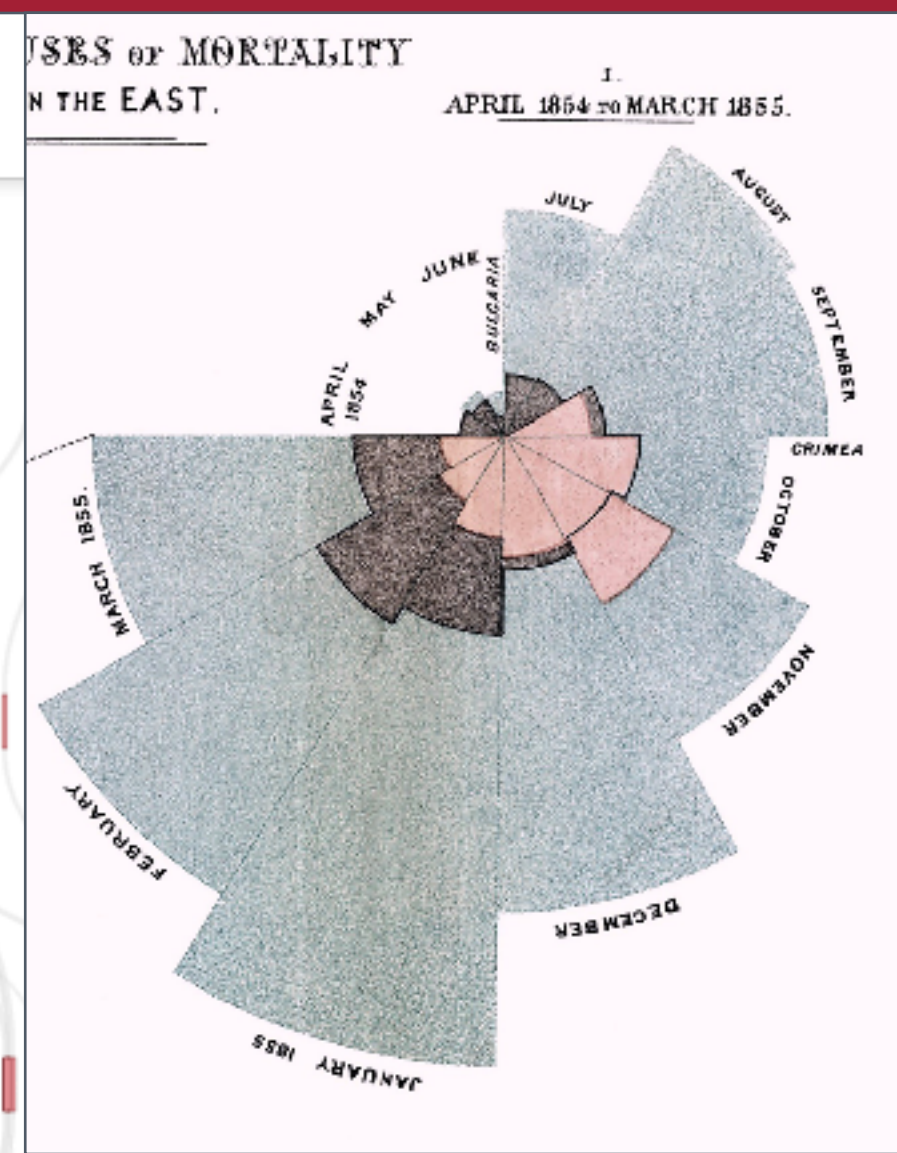
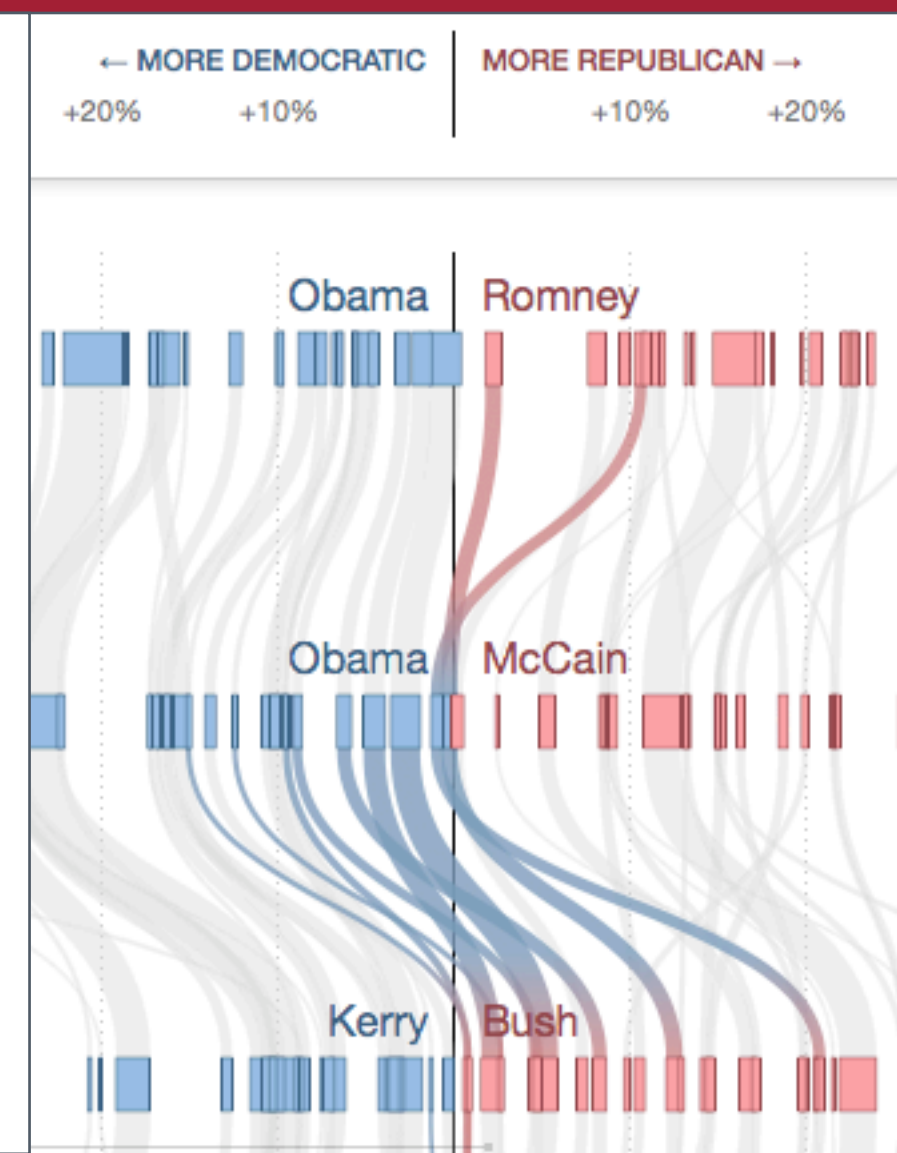
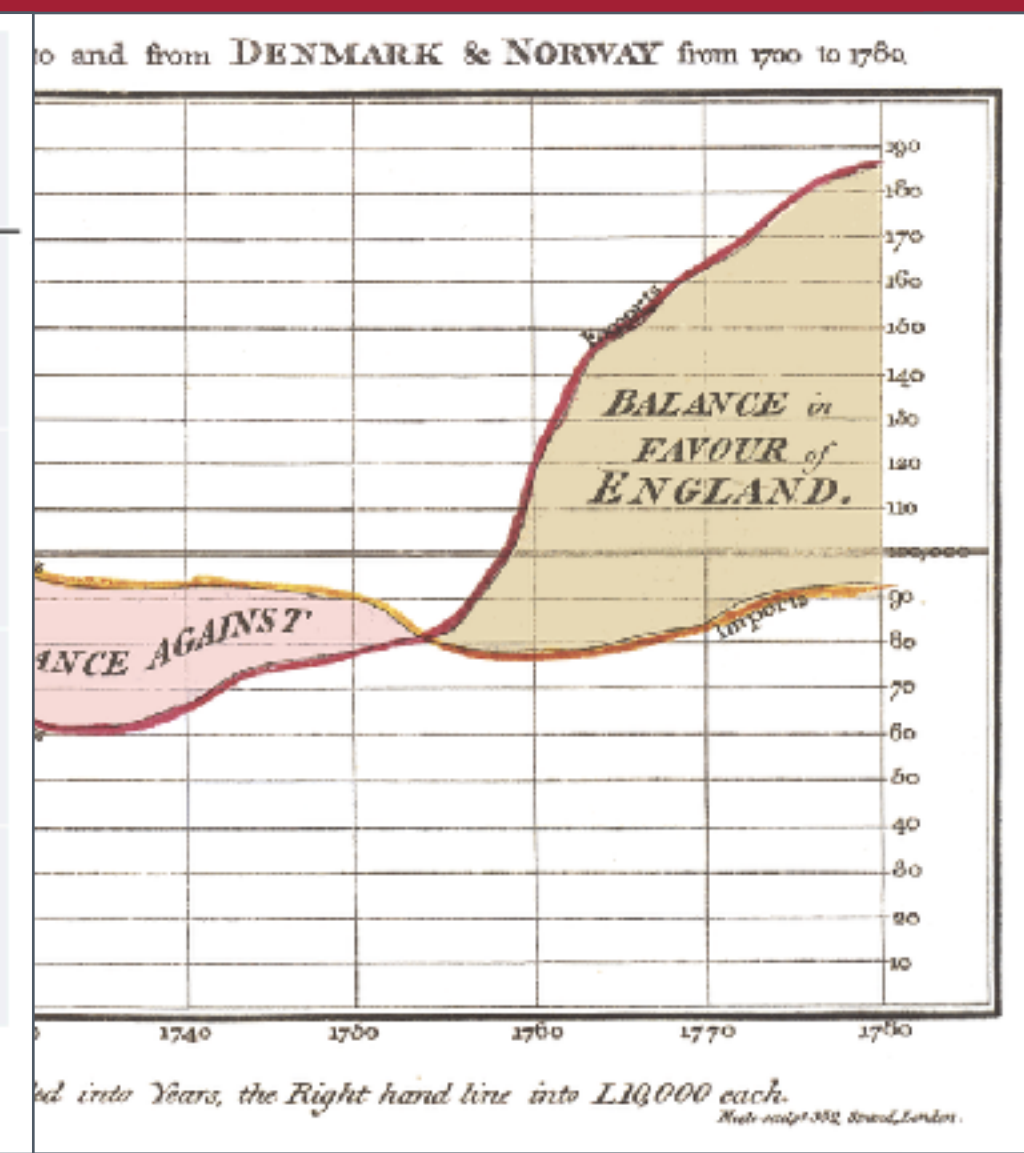
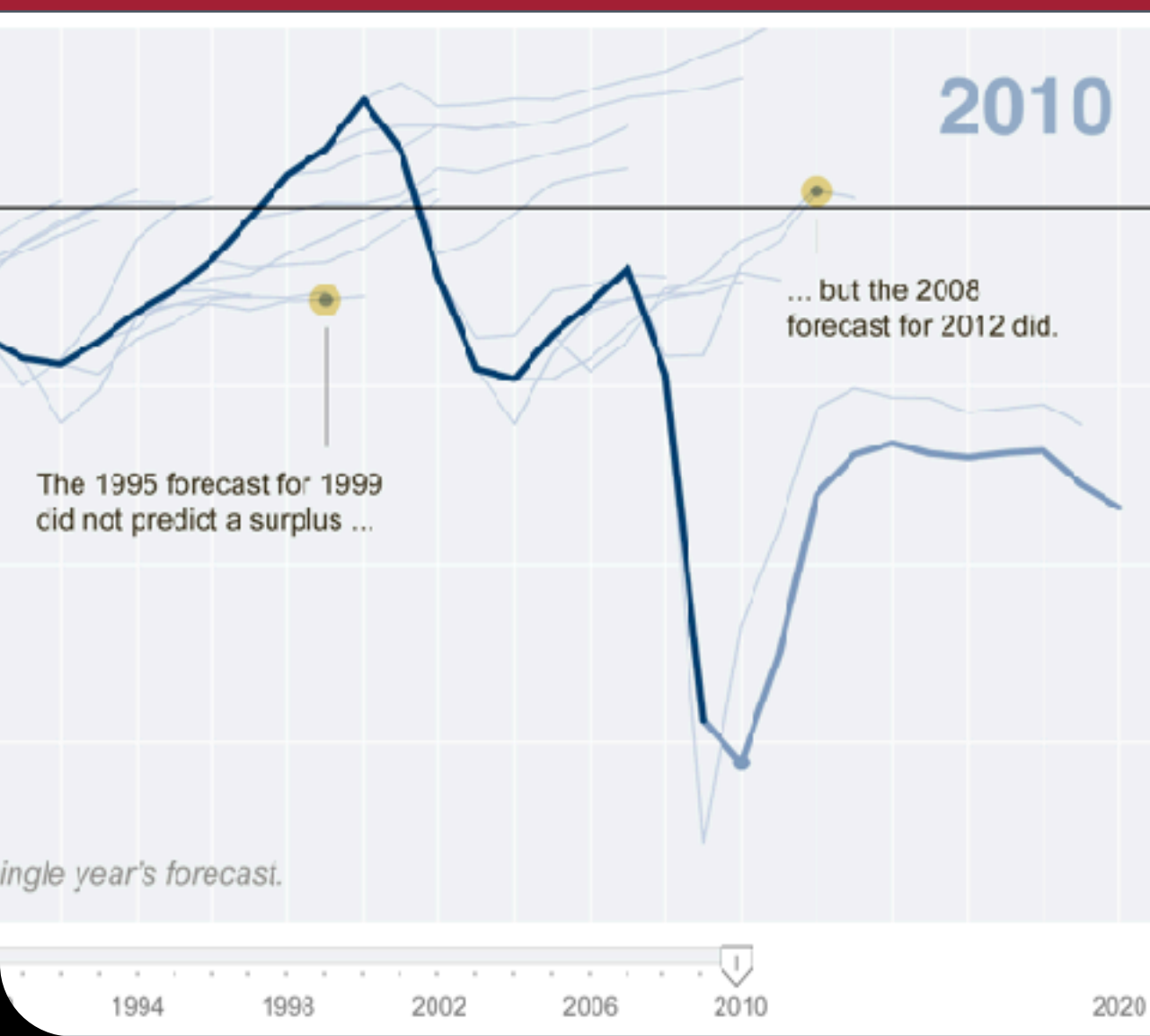


Source: Florida Department of Law Enforcement

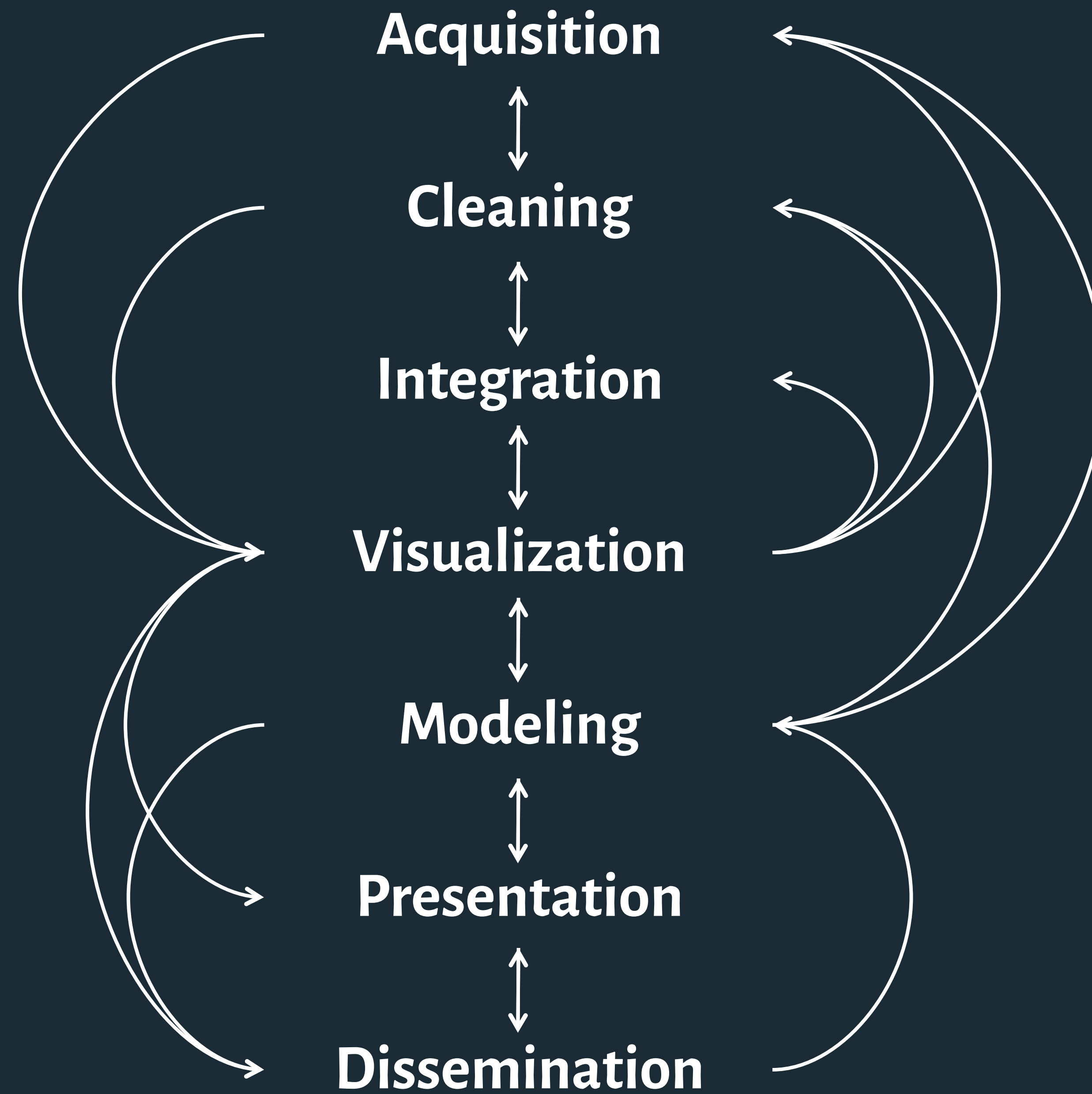
6.859: Interactive Data Visualization

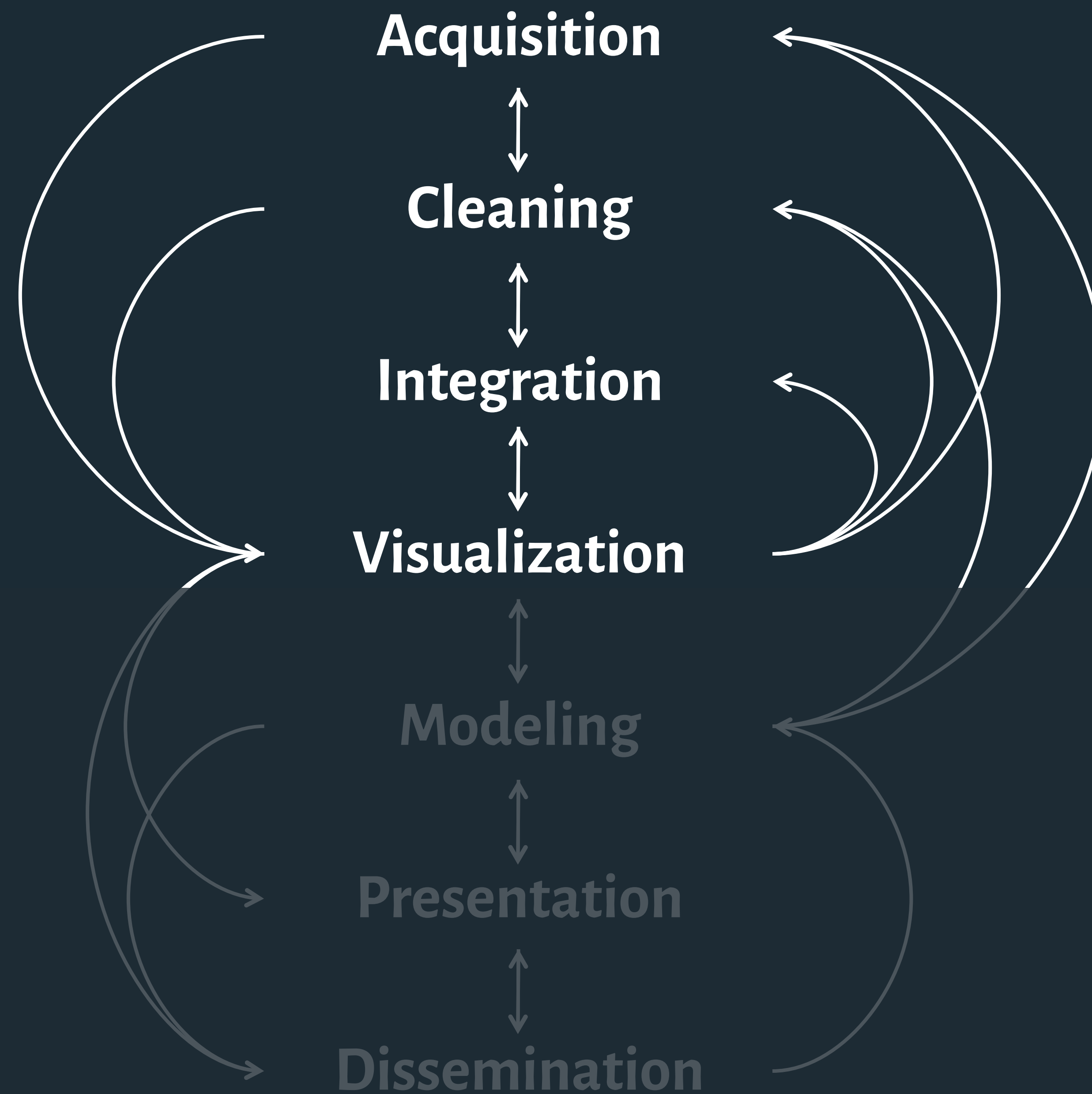
Exploratory Data Analysis

Arvind Satyanarayan



Data Visualization





Data Provenance

Who collected or produced it?

What was their intent?

Is it a reputable source?

What are the motives of the data producer (are they an advocate or lobbyist?)

Data Provenance

Who collected or produced it?

What was their intent?

Is it a reputable source?

What are the motives of the data producer (are they an advocate or lobbyist?)

The screenshot shows the Global Terrorism Database (GTD) website. At the top, there is a navigation bar with the GTD logo and links for 'ABOUT GTD', 'USING GTD', 'FAQ', 'TERMS OF USE', 'CONTACT', and a 'START HOME PAGE' button. Below the navigation bar is a search section with a red background. It features a search input field, a 'SEARCH' button, a link for 'I'm a New User', an 'ADVANCED SEARCH' button, and a 'Browse by:' dropdown menu with a 'Go' button. To the right of the search section is a grey box with the heading 'Information on more than 180,000 Terrorist Attacks'. The text below explains that the GTD is an open-source database covering terrorist events from 1970 to 2017, including domestic and international incidents. It includes a 'Learn more' link and a link to 'Global Terrorism in 2017'. Below the search section are three columns of content. The first column, 'GTD DATA VISUALIZATIONS', features a world map titled '45 Years of Terrorism' showing global terrorism from 1970 to 2015. Below the map are links for world maps from 2017, 2016, 2015, 2014, 2013, and 2012. The second column, 'THIS DATE IN TERRORISM', highlights 'February 9' with two entries: '2015 Bantacan, Philippines' and '2015 Logo district, Nigeria'. Each entry includes a brief description of the event and a 'Learn more' link. The third column, 'FEATURED', contains a 'Message from the Global Terrorism Database Manager' and a paragraph explaining the database's funding by the U.S. State Department and the impact of the loss of that funding in 2018. It also includes a 'Continue Reading' link.

Data Provenance

When was the data collected?

Measurements can change over time.

Definitions/interpretations in quantification can change over time.

Is the data recent, and how much does that matter to the insight you wish to convey?

The Three-Year Plunge

To help gauge each city's overall crime level, the FBI tracks eight "index crimes." From 1993 to 2010, Chicago's annual total dropped by 47 percent. But from 2010 to 2013, it dropped a stunning 56 percent, or nearly 19 percent per year, according to data from the Chicago Police Department.



Graph Viewer

Roll-up by:

All

Visualization:

Node-Link

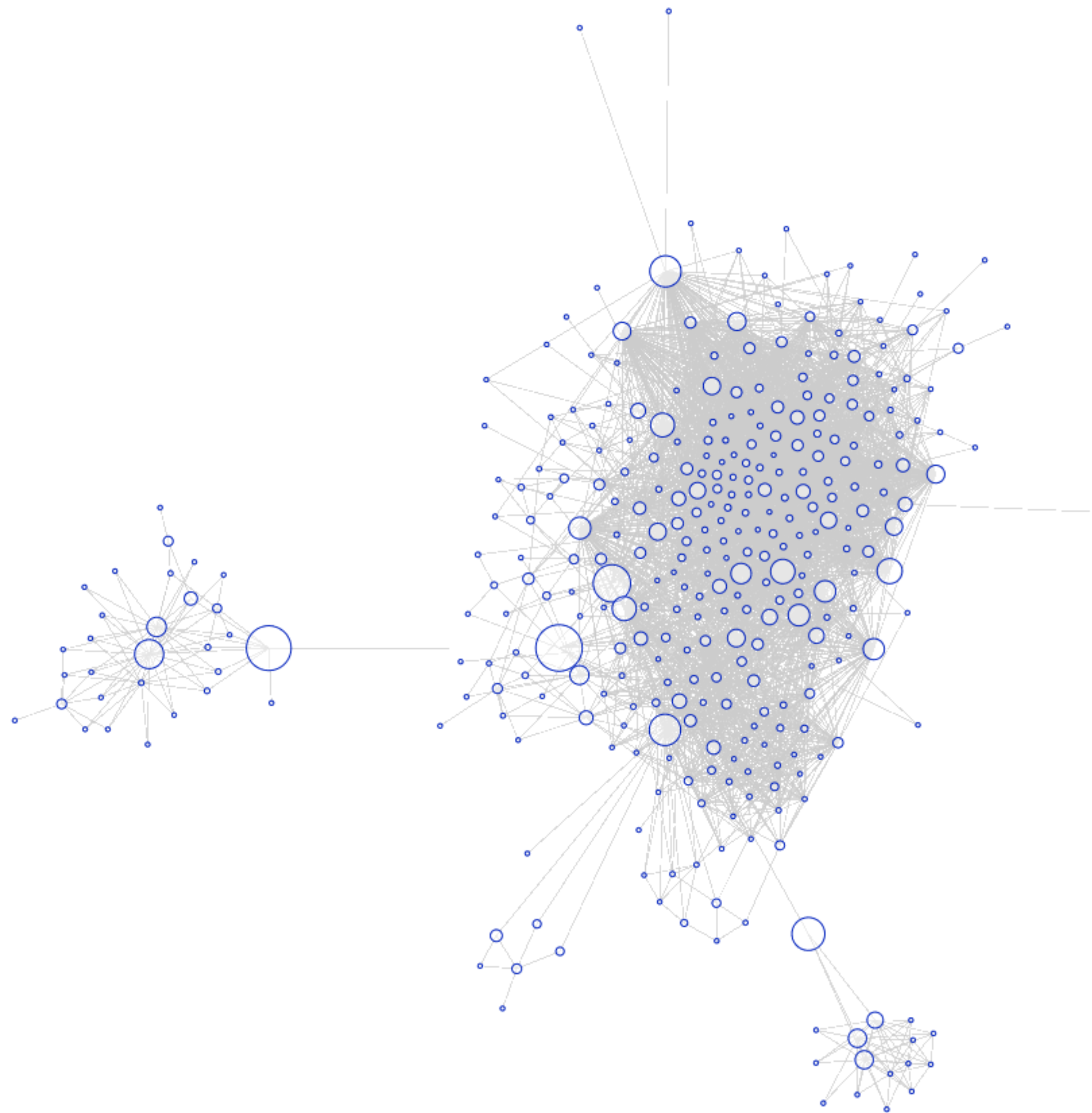
Sort by:

None

Edge centrality filters:

Two horizontal sliders for edge centrality filters, both currently set to the minimum value.

- Images
- Animate





Graph Viewer

Roll-up by:

All

Visualization:

Matrix

Sort by:

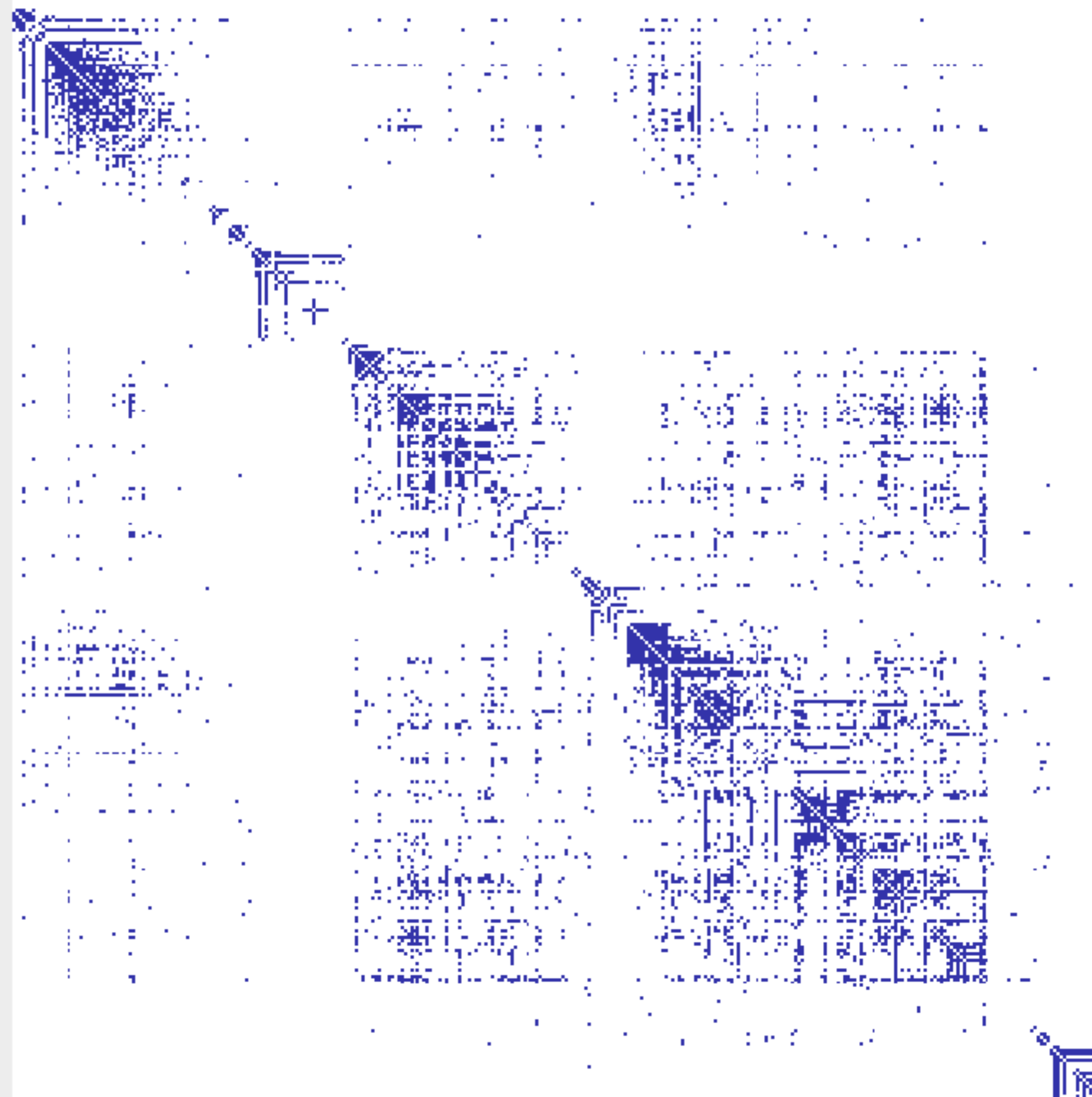
Linkage

Edge centrality filters:

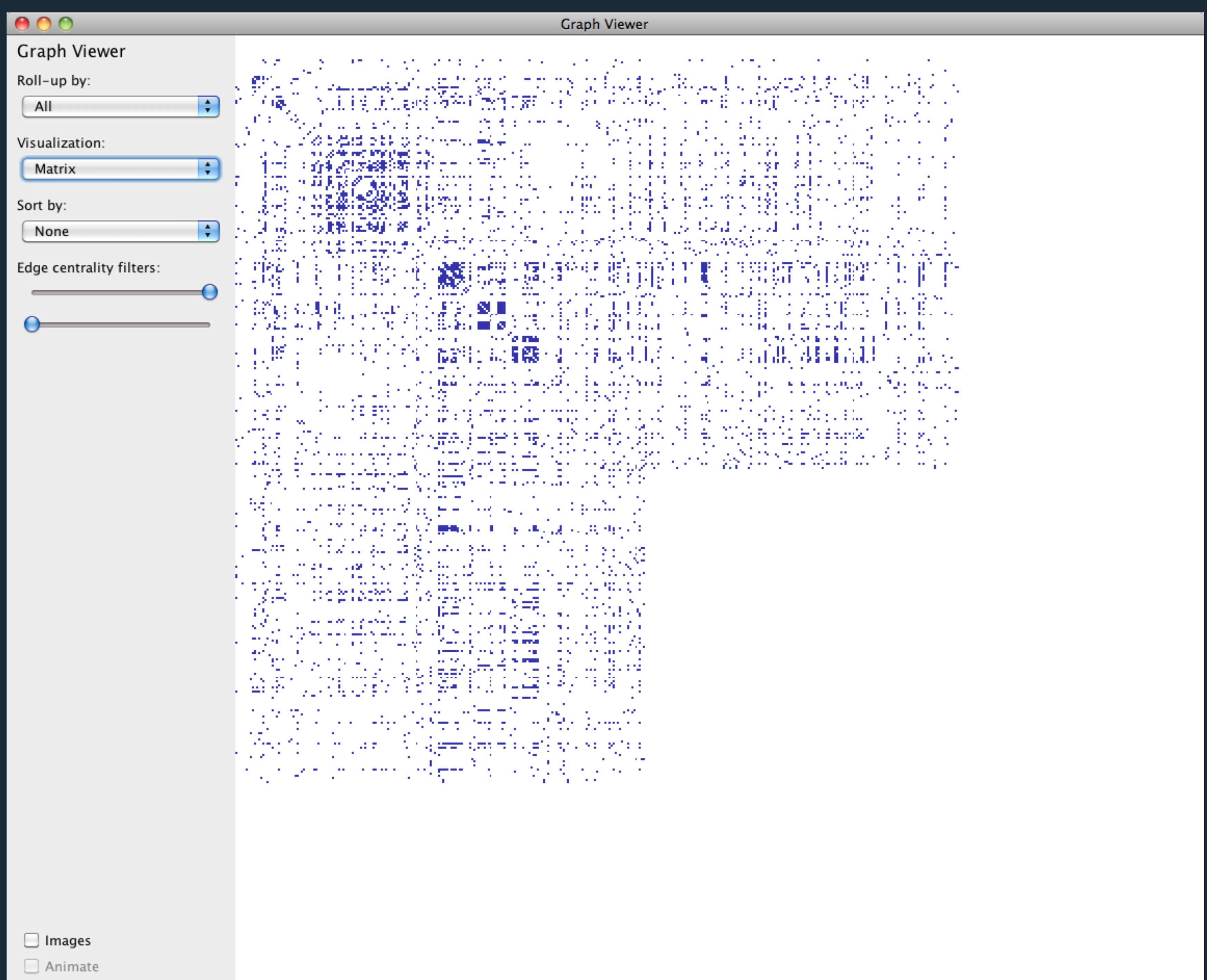
Two horizontal sliders for edge centrality filters, one above the other.

Images

Animate



Missing Values



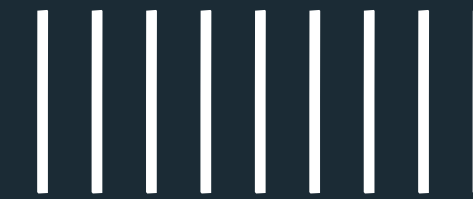
Berkeley



Cornell



Harvard



Harvard University



Stanford



Stanford University



UC Berkeley



UC Davis



University of California at Berkeley



University of California, Berkeley



University of California, Davis



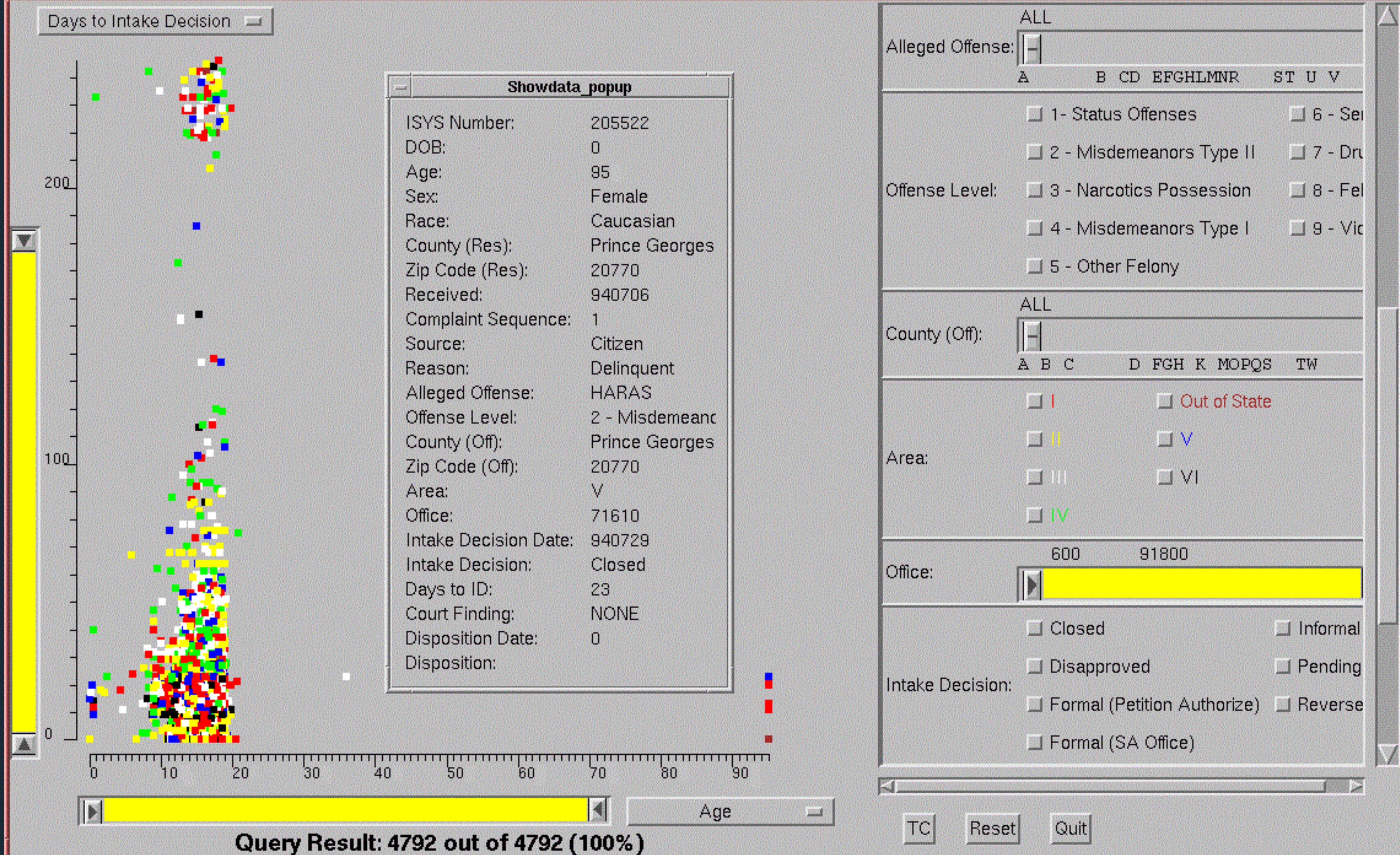
“The first sign that a visualization is good is that **it shows you a problem in your data**. Every successful visualization that I've been involved with has had this stage where you realize, **"Oh my God, this data is not what I thought it would be!"** So already, you've discovered something.”

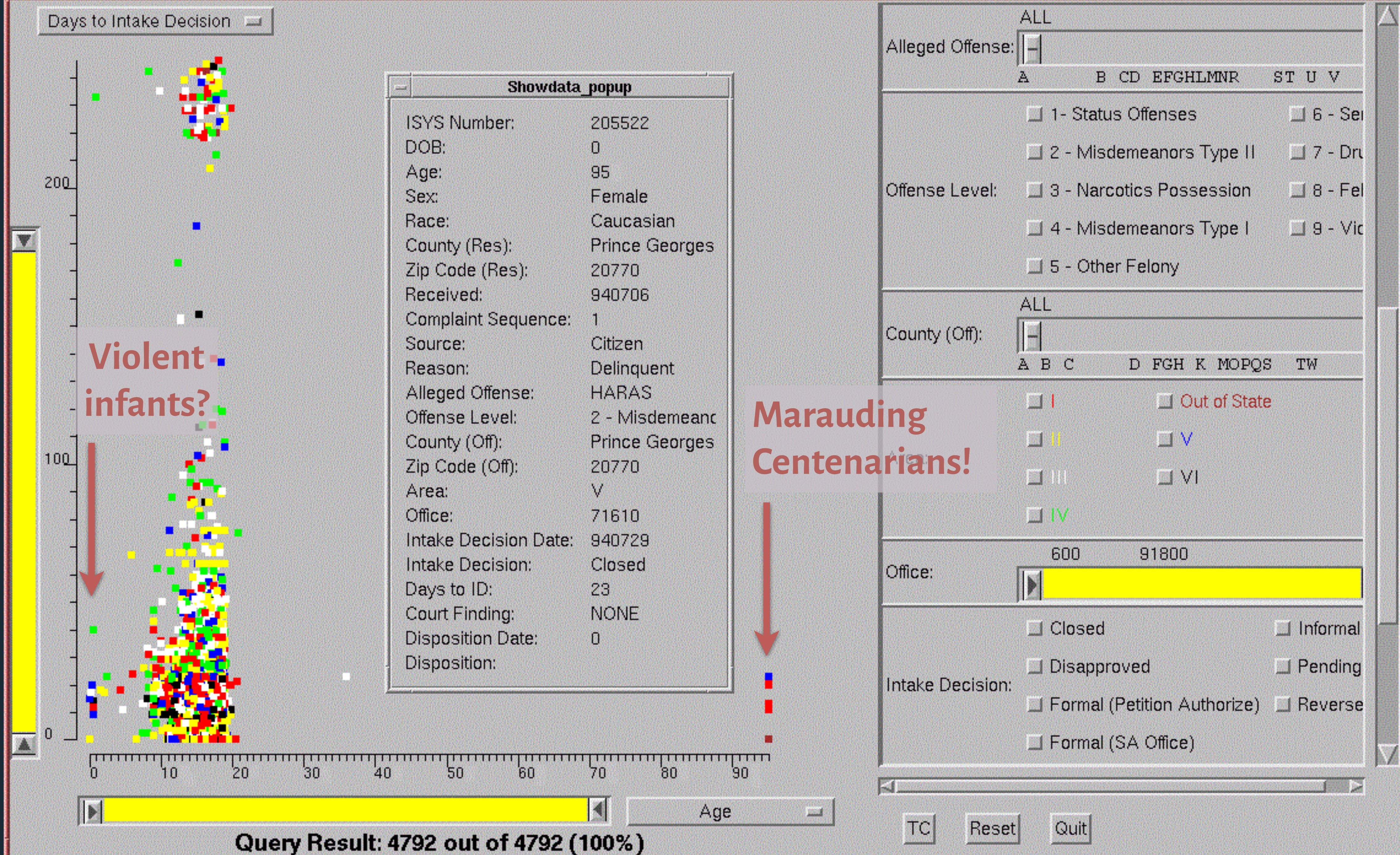
– **Martin Wattenberg**

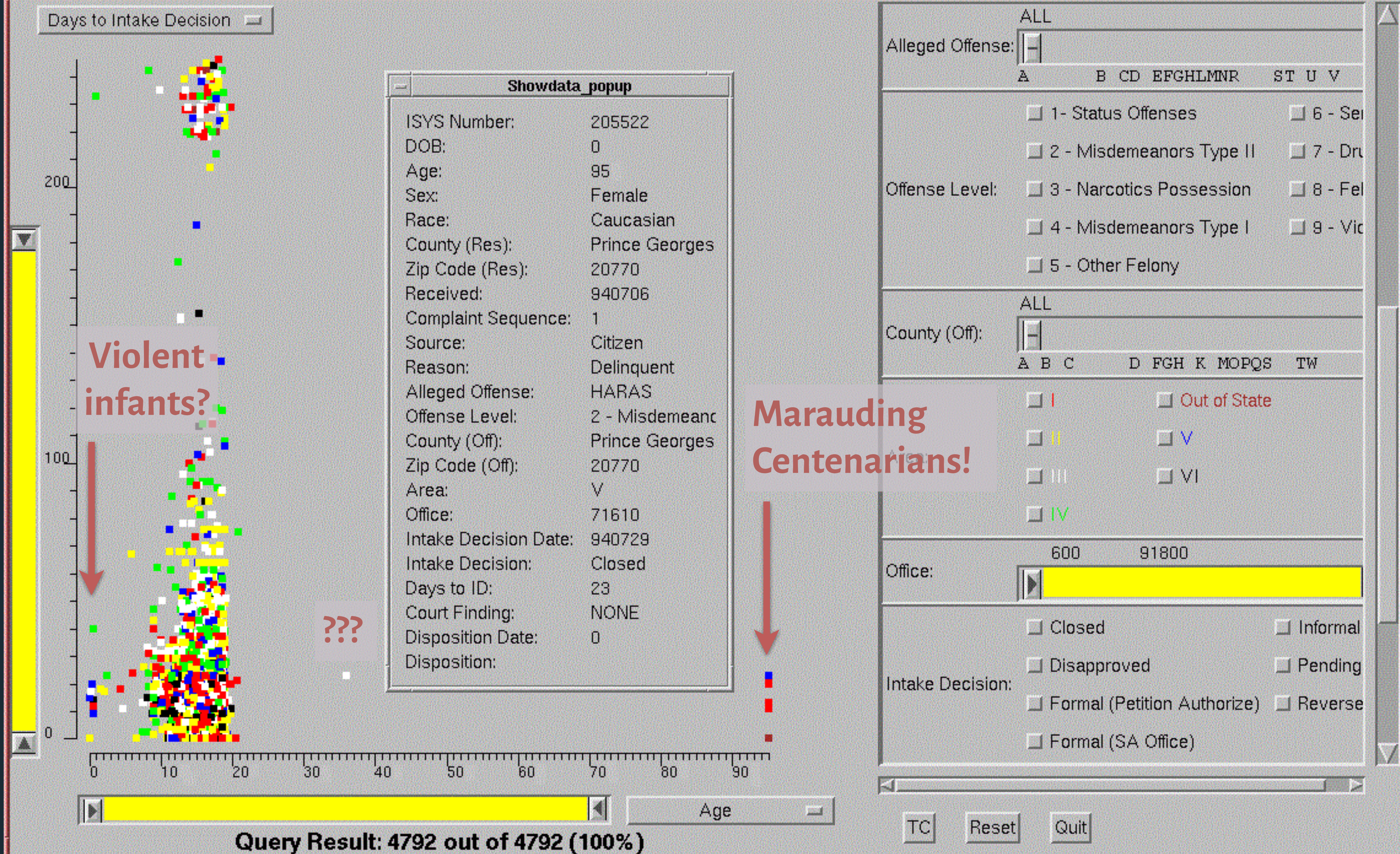
Co-lead of Google's People + AI Initiative

ACM Queue, Mar 2010













Big Data Borat

@BigDataBorat

Follow



In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

6:47 PM - 26 Feb 2013

540 Retweets 343 Likes



 12

 540

 343



“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**

[Kandel et al. VAST 2012]

Reported crime in Alabama

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	4525375	4029.3	987	2732.4	309.9
2005	4548327	3900	955.8	2656	289
2006	4599030	3937	968.9	2645.1	322.9
2007	4627851	3974.9	980.2	2687	307.7
2008	4661900	4081.9	1080.7	2712.6	288.6

Reported crime in Alaska

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	657755	3370.9	573.6	2456.7	340.6
2005	663253	3615	622.8	2601	391
2006	670053	3582	615.2	2588.5	378.3
2007	683478	3373.9	538.9	2480	355.1
2008	686293	2928.3	470.9	2219.9	237.5

Reported crime in Arizona

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	5739879	5073.3	991	3118.7	963.5
2005	5953007	4827	946.2	2958	922
2006	6166318	4741.6	953	2874.1	914.4
2007	6338755	4502.6	935.4	2780.5	786.7
2008	6500180	4087.3	894.2	2605.3	587.8

Reported crime in Arkansas

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	2750000	4033.1	1096.4	2699.7	237
2005	2775708	4068	1085.1	2720	262
2006	2810872	4021.6	1154.4	2596.7	270.4
2007	2834797	3945.5	1124.4	2574.6	246.5
2008	2855390	3843.7	1182.7	2433.4	227.6

Reported crime in California

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	35842038	3423.9	686.1	2033.1	704.8
2005	36154147	3321	692.9	1915	712
2006	36457549	3175.2	676.9	1831.5	666.8
2007	36553215	3032.6	648.4	1784.1	600.2
2008	36756666	2940.3	646.8	1769.8	523.8

Reported crime in Colorado

Year	Population	Property crime rate	Burglary rate	Larceny-theft rate	Motor vehicle theft rate
2004	4601821	3918.5	717.3	2679.5	521.6

DataWrangler

Suggestions

- Delete rows 8,10
- Delete empty rows
- Delete rows where Property_crime_rate is null
- Delete rows where Year is null

Script Export

- ▶ Split data repeatedly on newline into rows
- ▶ Split data repeatedly on ','

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

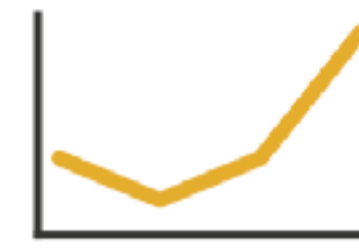
Wrangler: Interactive Visual Specification of Data Transformation Scripts. Sean Kandel et al., ACM CHI 2011.

Exploratory Visual Analysis

Process

1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...

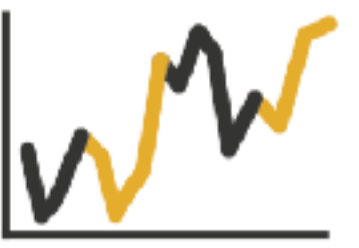
→ Trends



→ Outliers



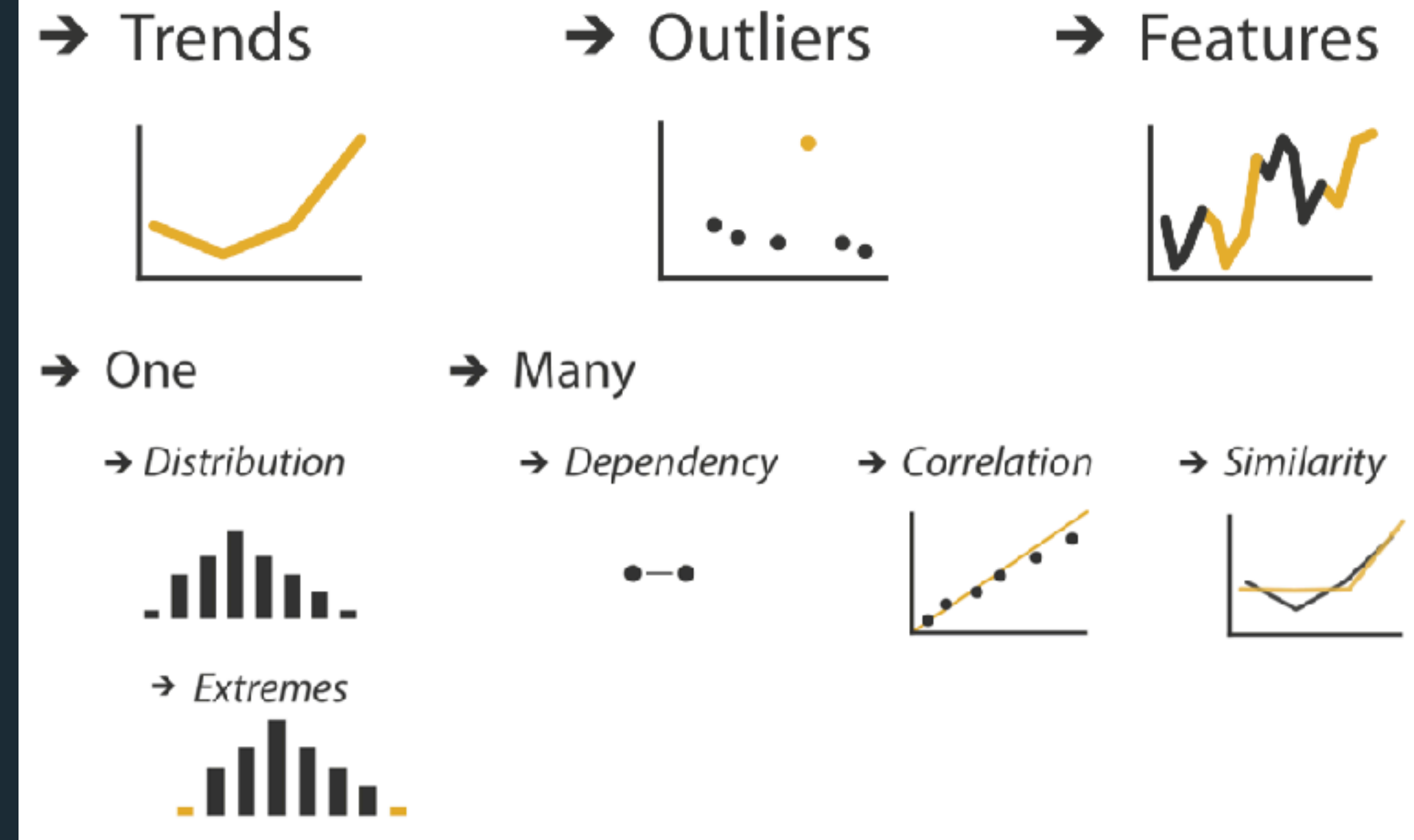
→ Features



Exploratory Visual Analysis

Process

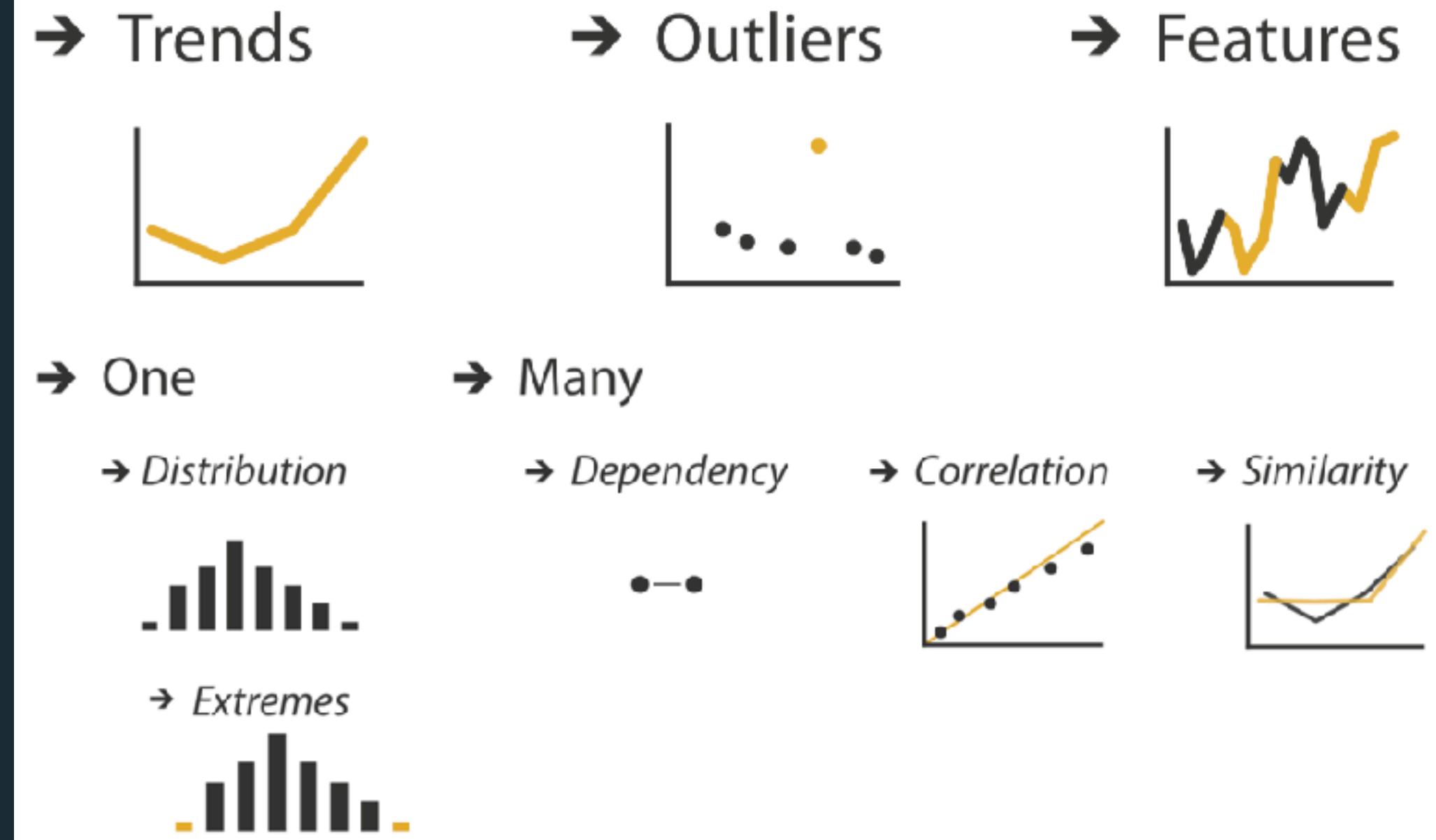
1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...



Exploratory Visual Analysis

Process

1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...



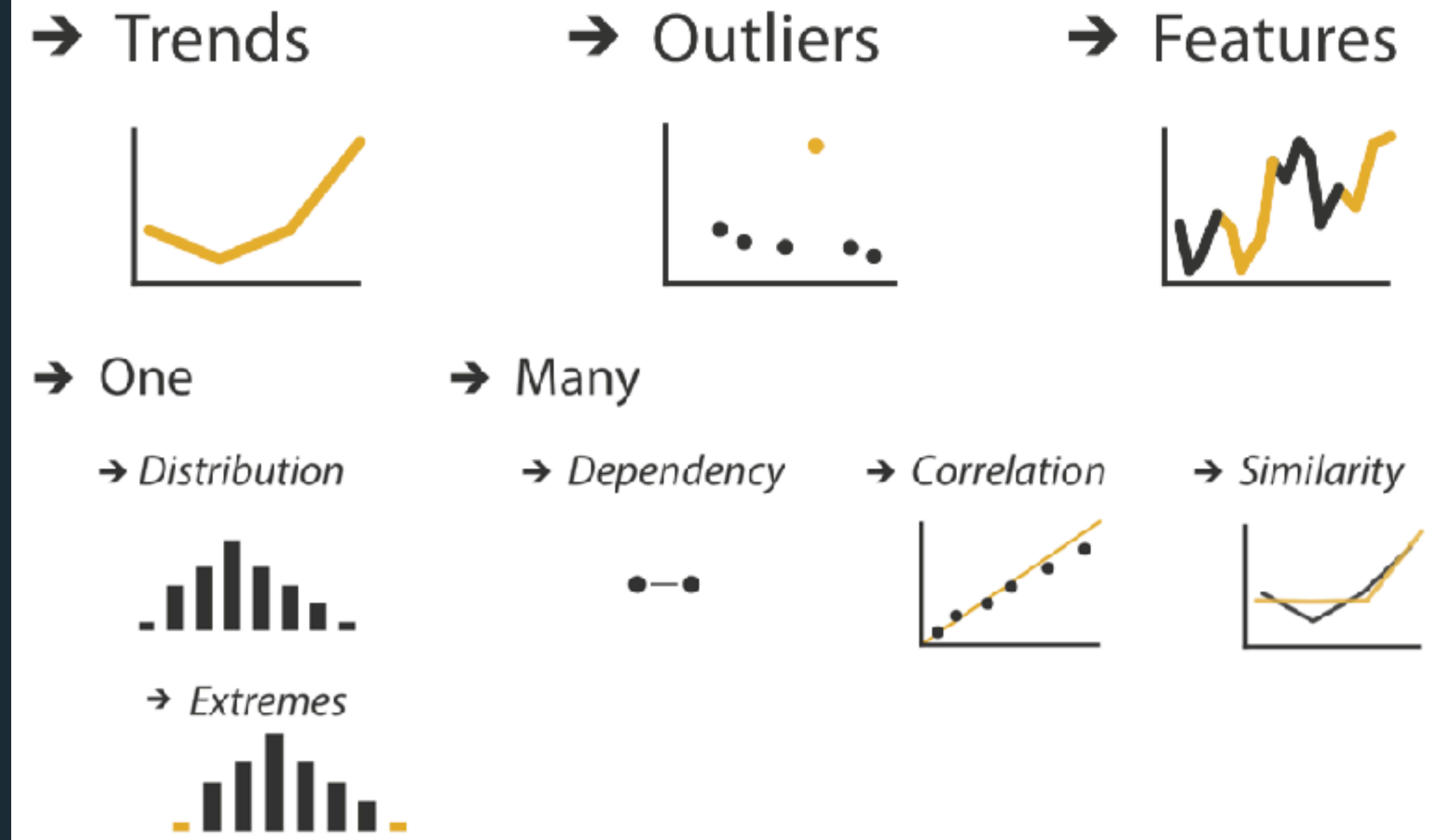
→ Search

	Target known	Target unknown
Location known	Lookup	Browse
Location unknown	Locate	Explore

Exploratory Visual Analysis

Process

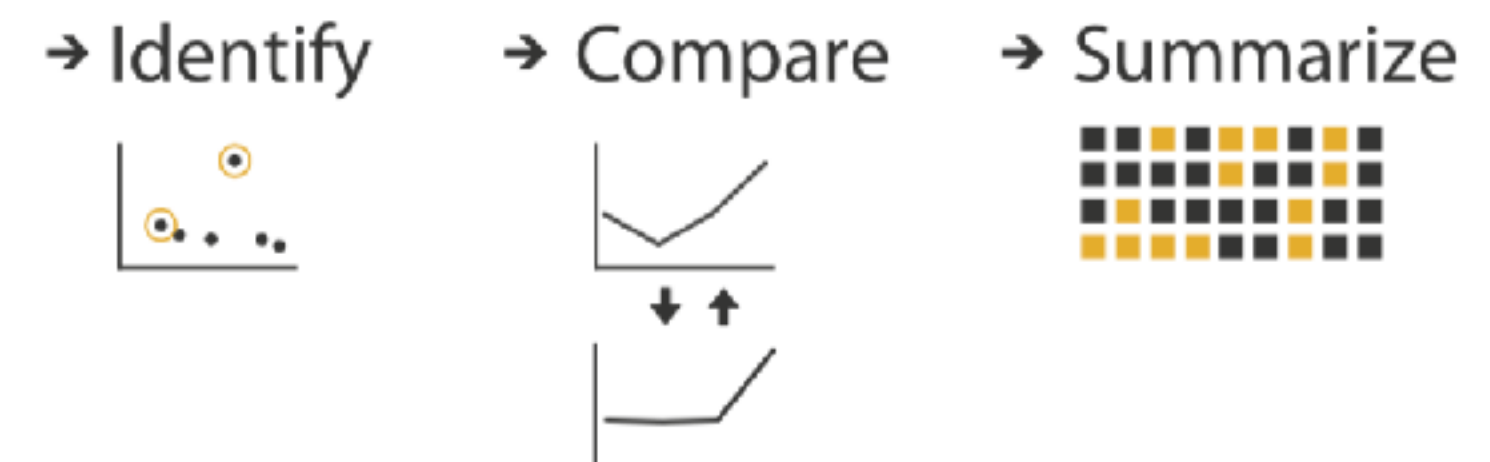
1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...



Search

	Target known	Target unknown
Location known	Lookup	Browse
Location unknown	Locate	Explore

Query



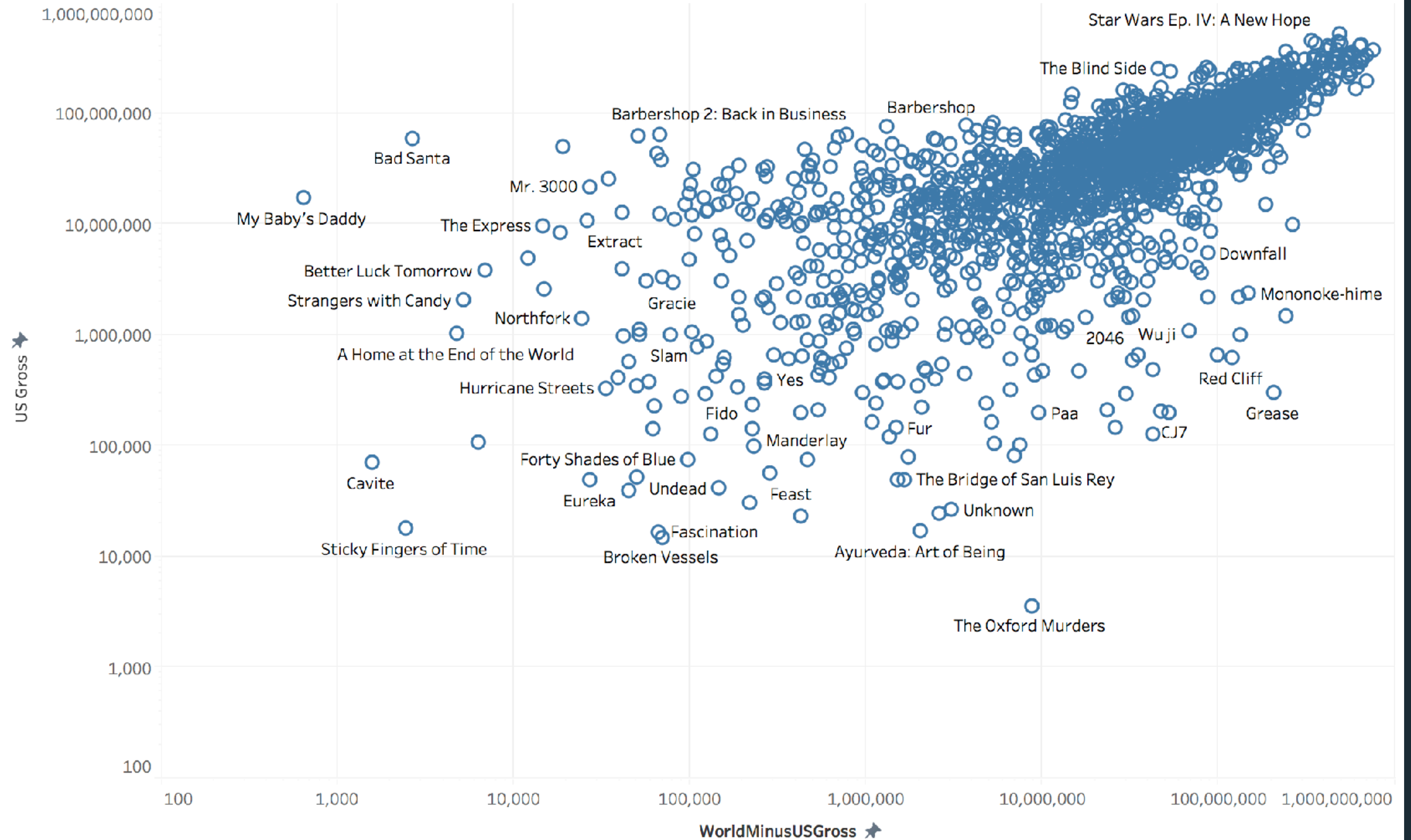
Analysis Example: Motion Pictures Data

A sample of 3,201 movies collected in 2010.

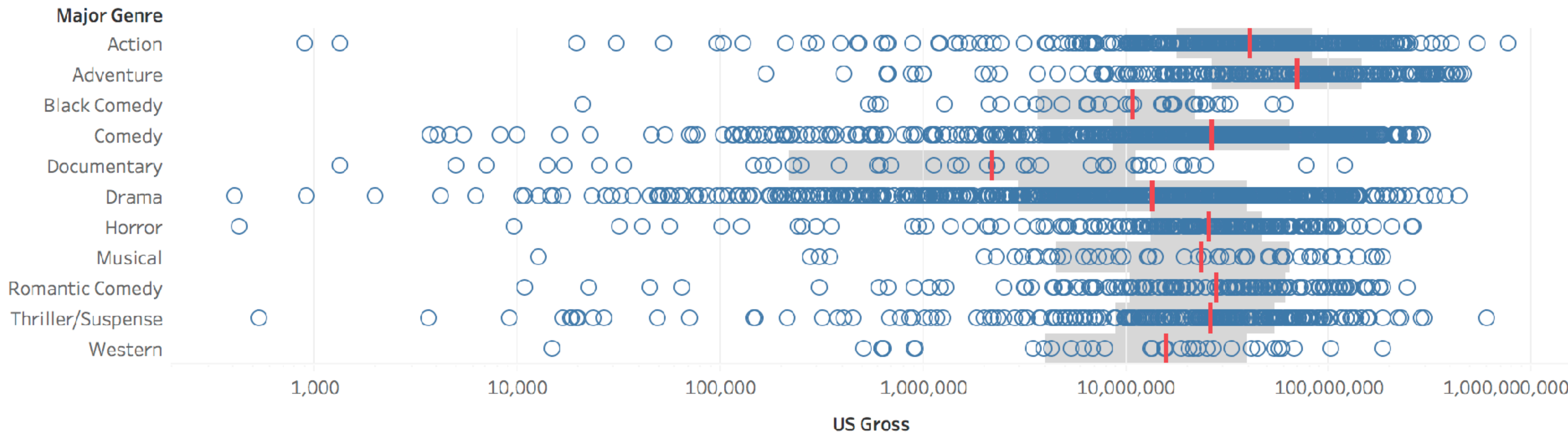
Analysis Example: Motion Pictures Data

Title	String (N)
IMDB Rating	Number (Q)
Rotten Tomatoes Rating	Number (Q)
Genre	String (N)
Release Date	Date (T)
US Gross	Number (Q)
Worldwide Gross	Number (Q)

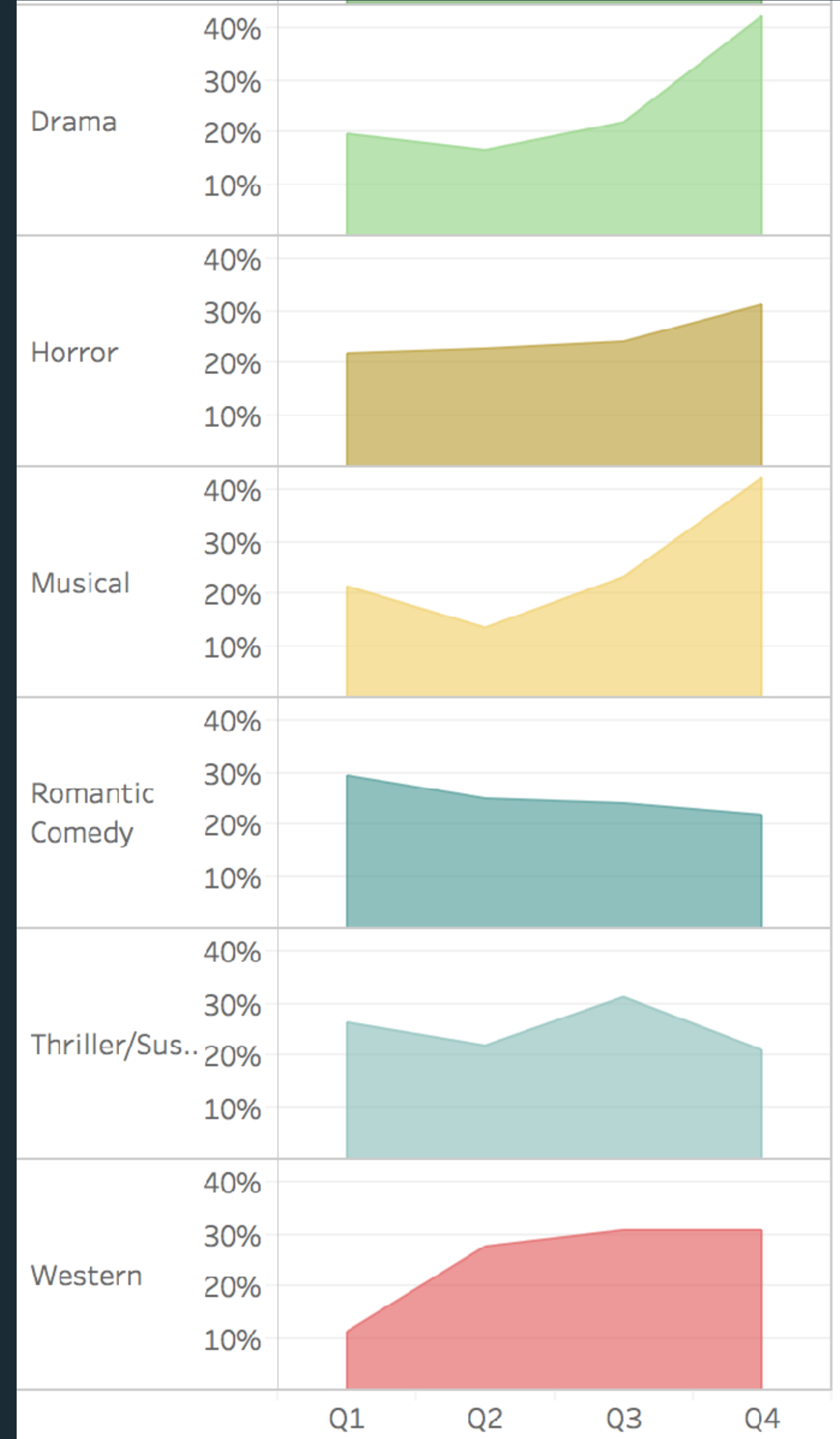
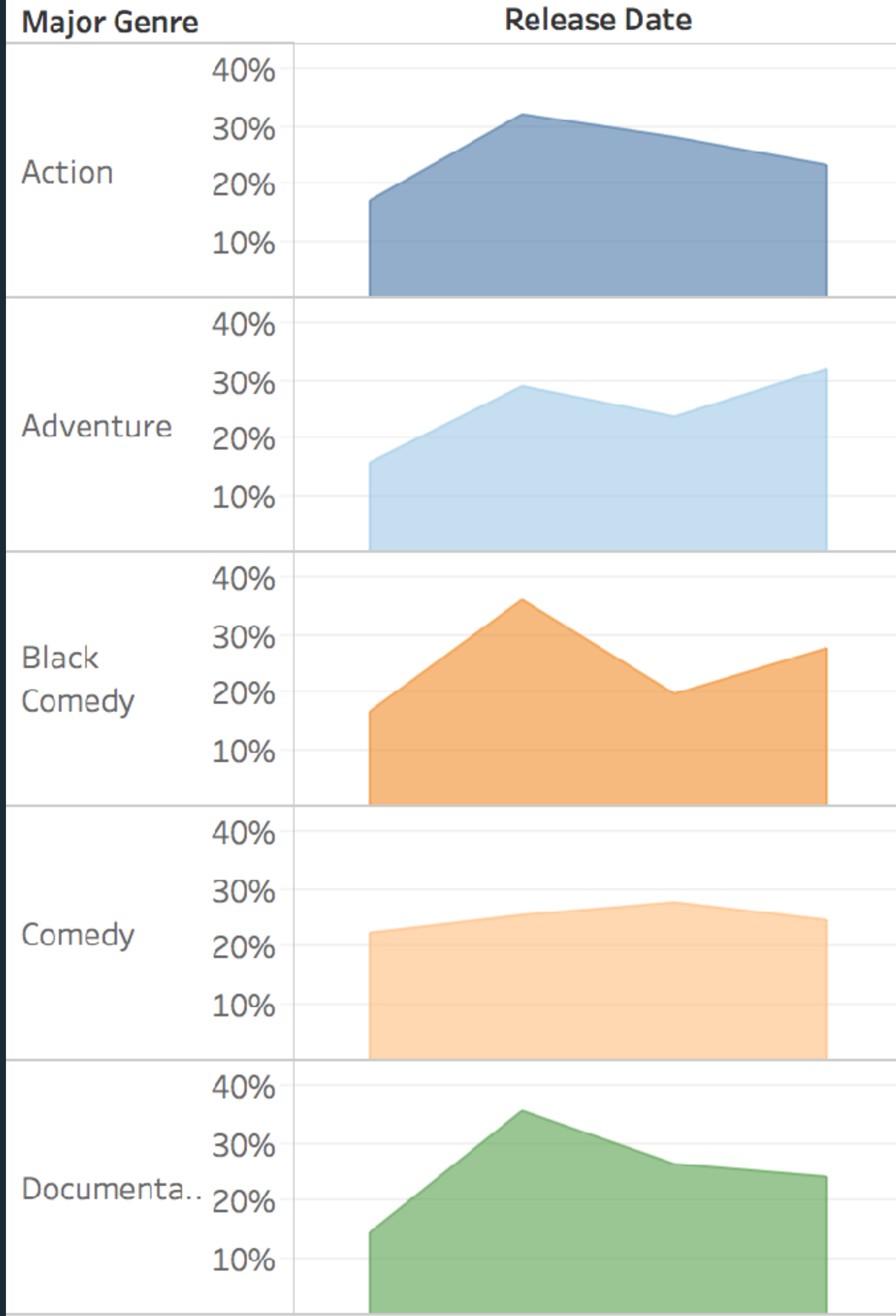
US vs WW



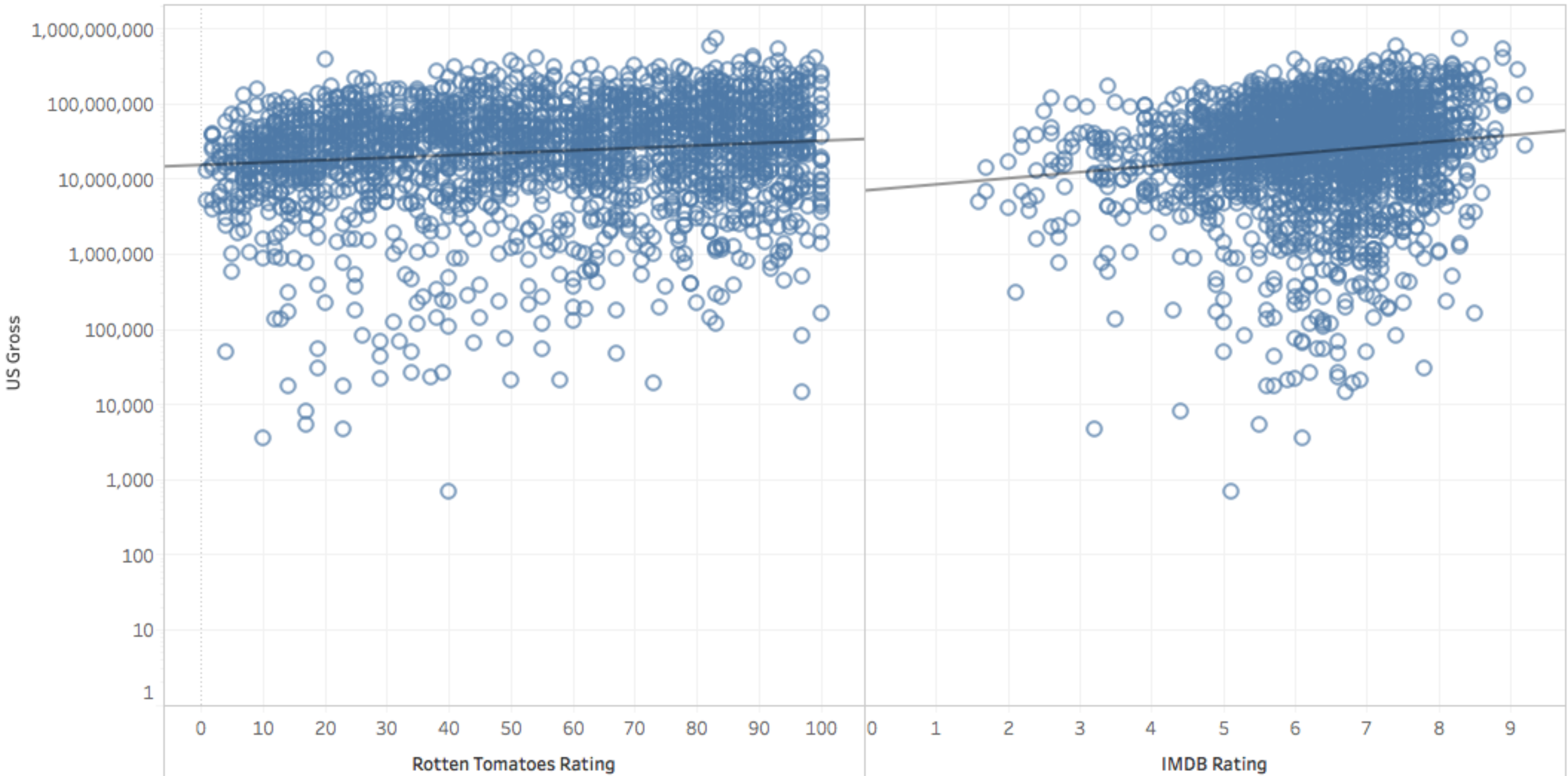
Distribution of US Gross by Genre



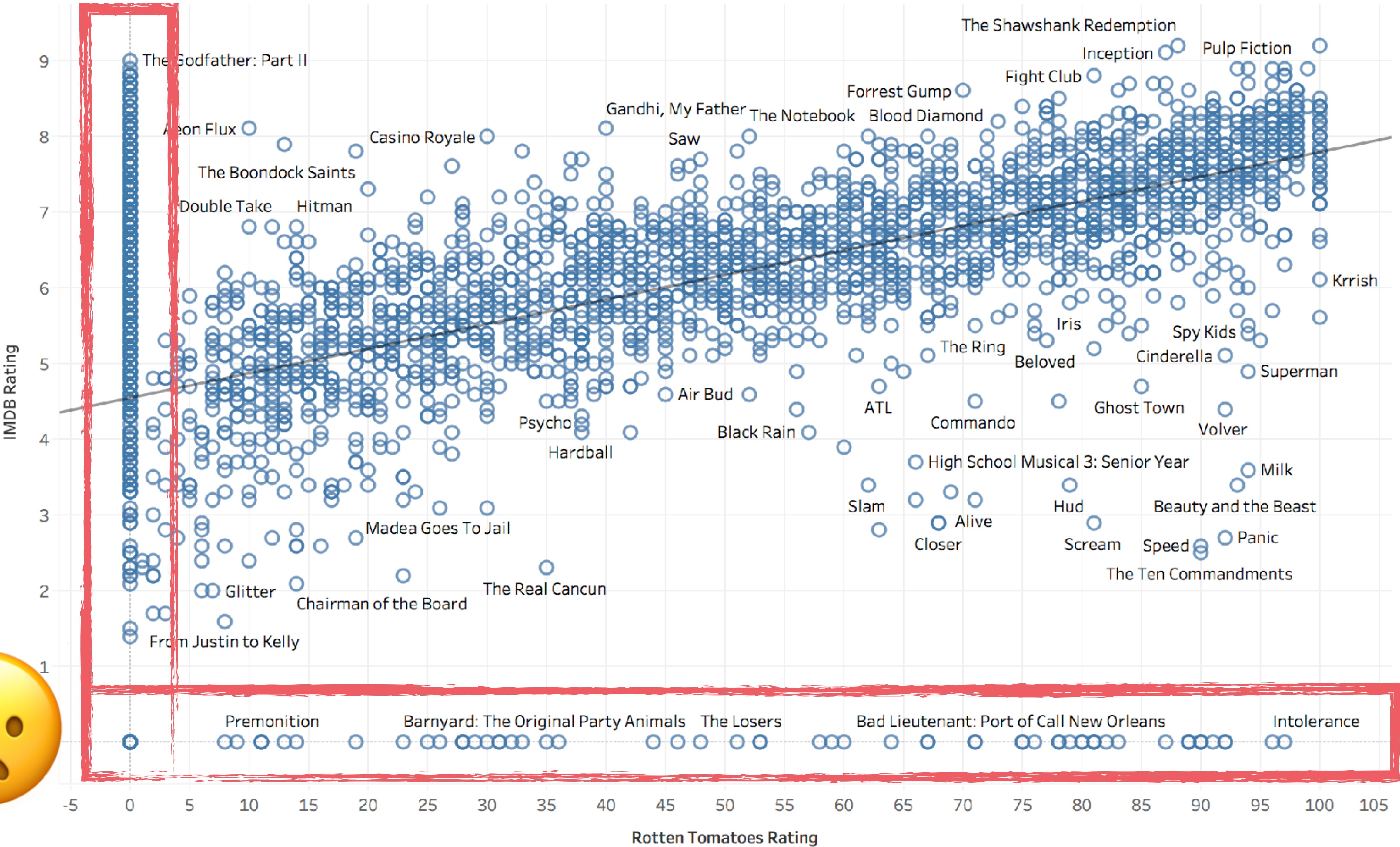
% Releases per Quarter per Genre

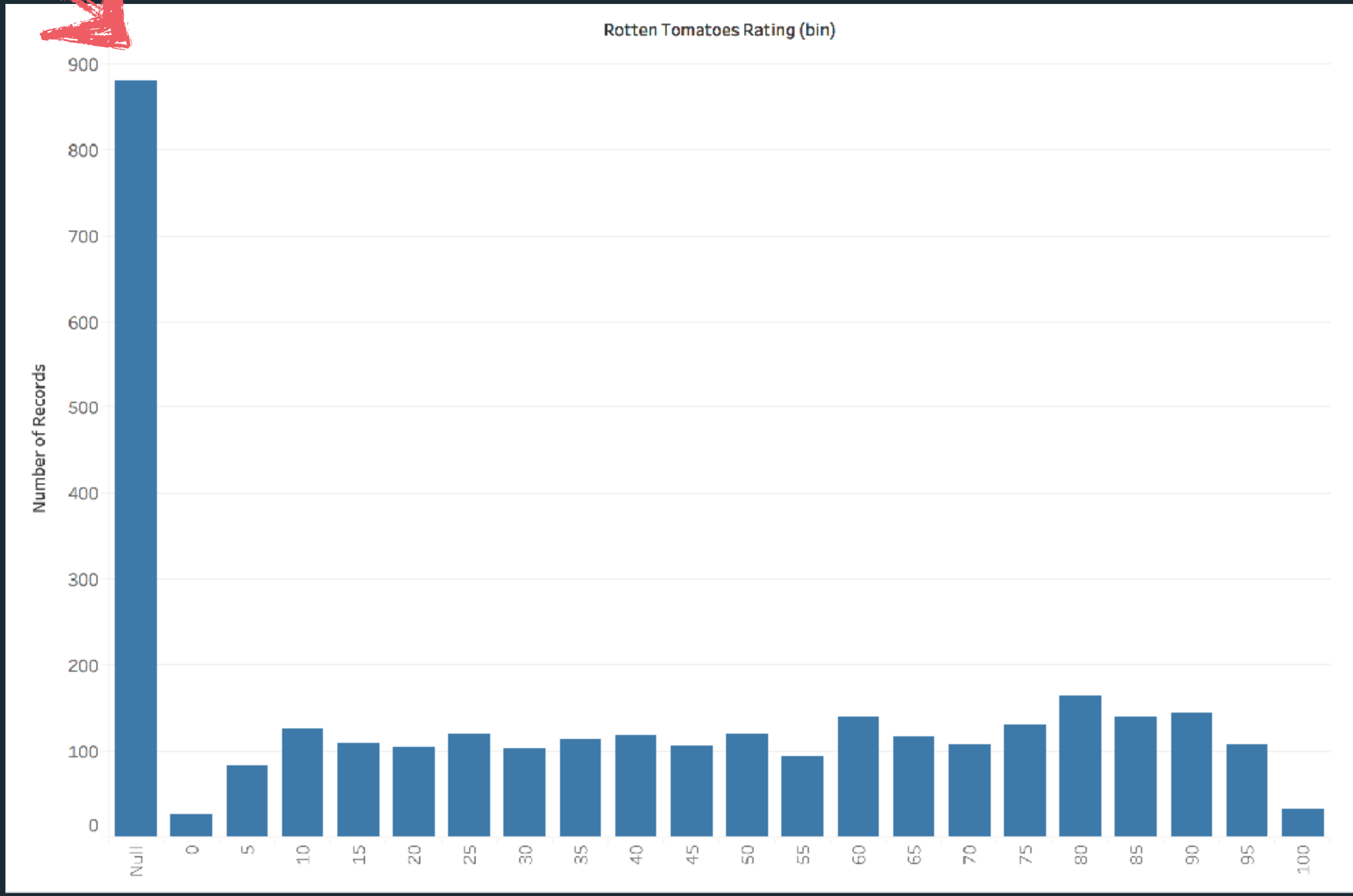
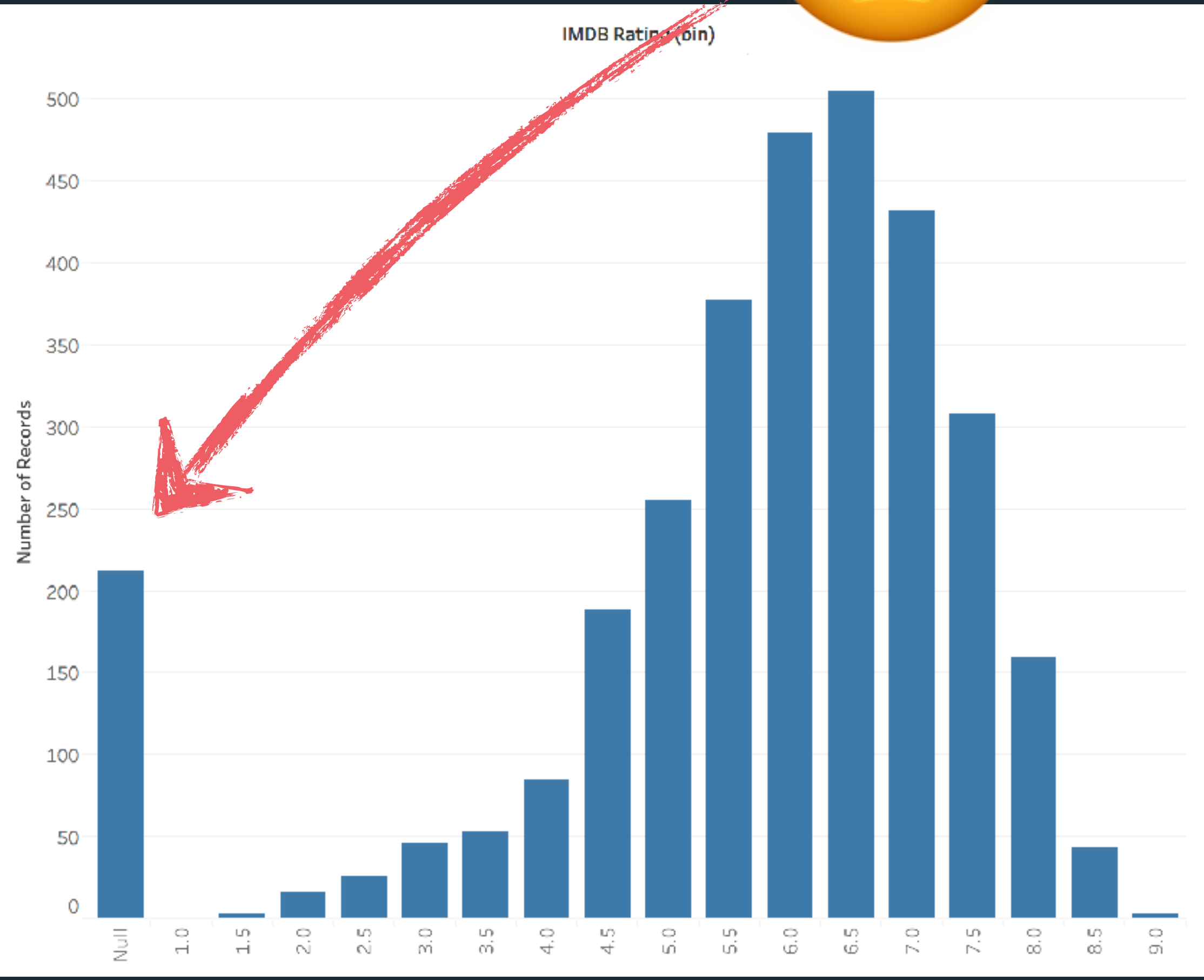


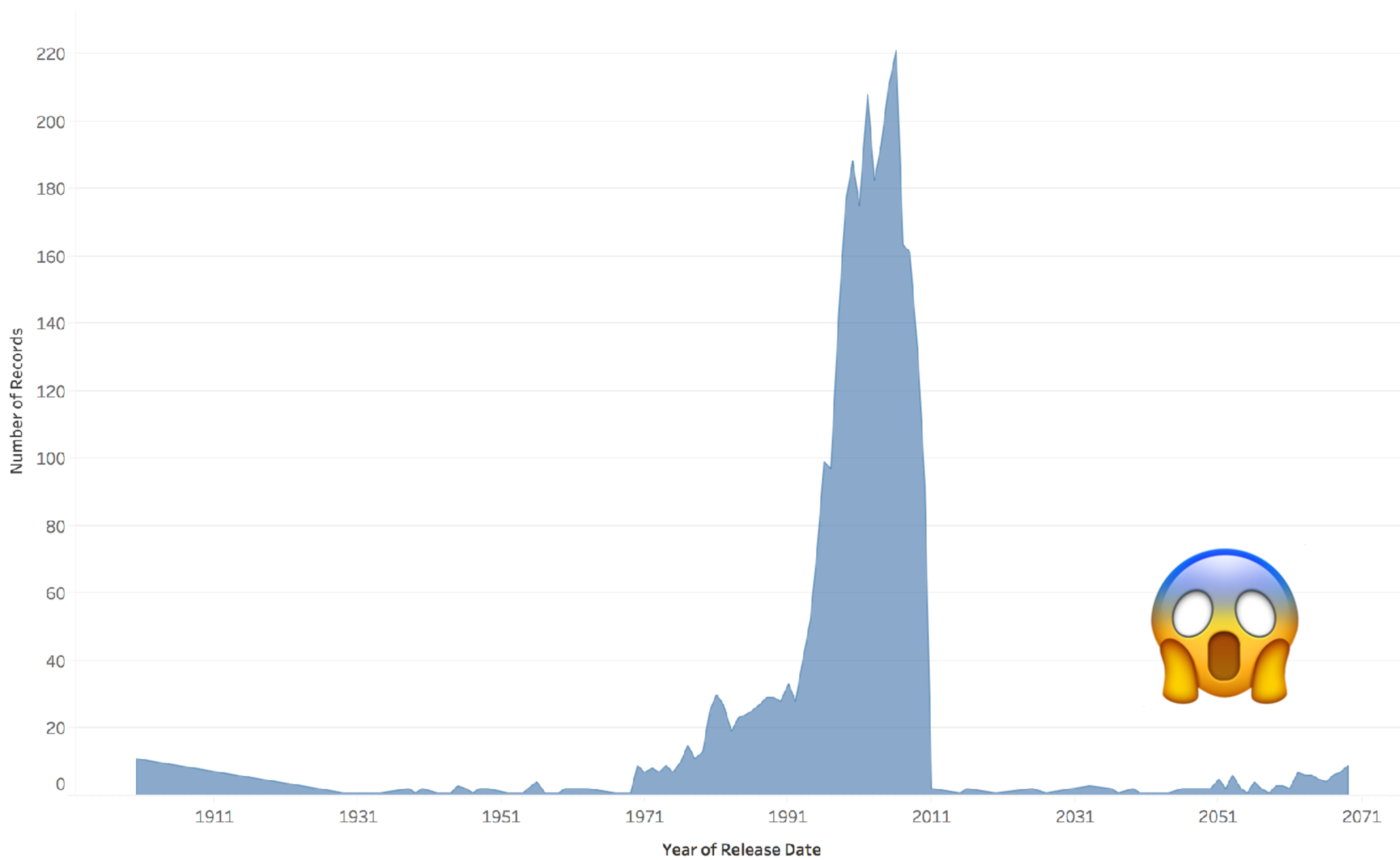
US Gross by Ratings



Audience vs. Critic Ratings





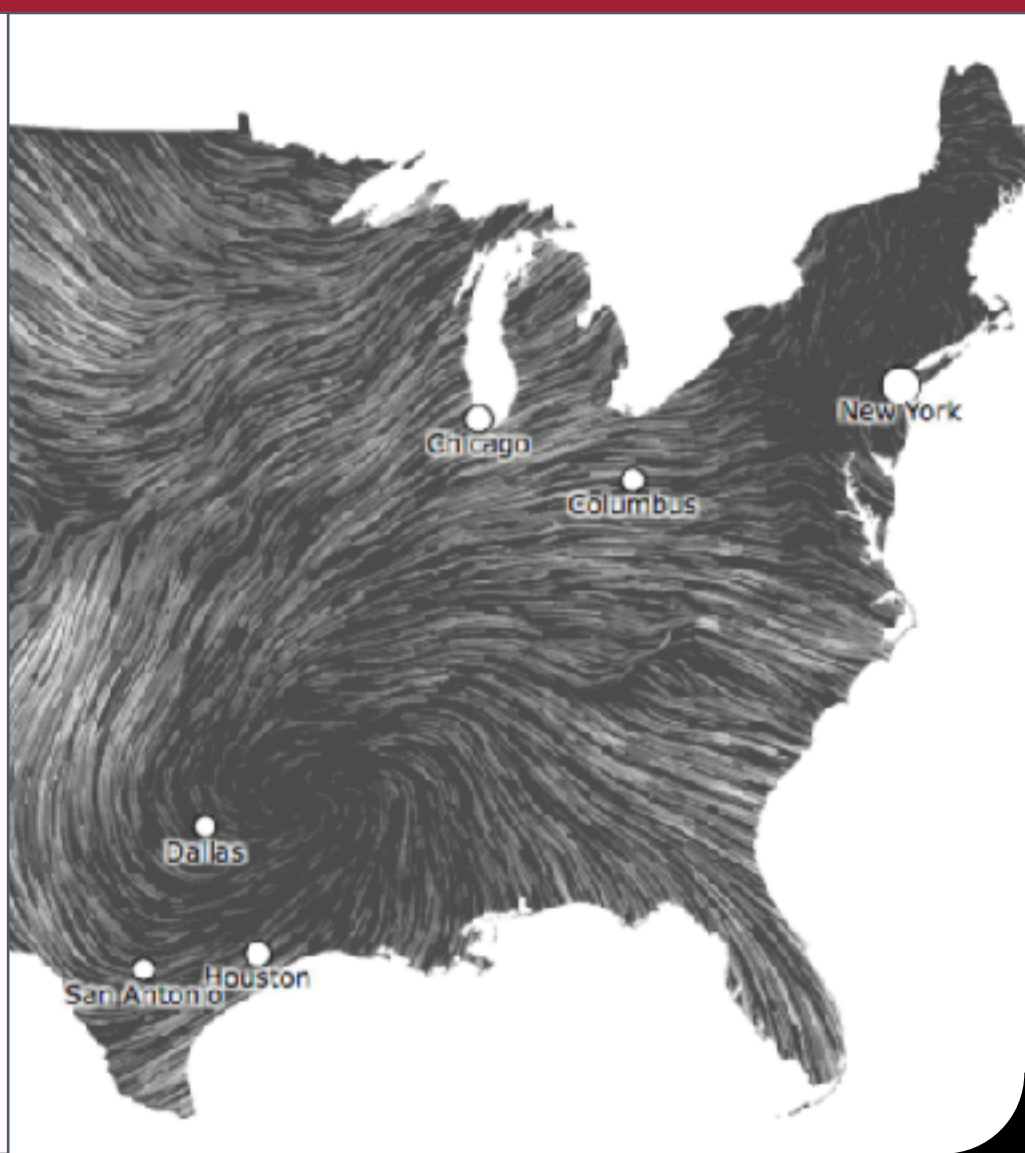
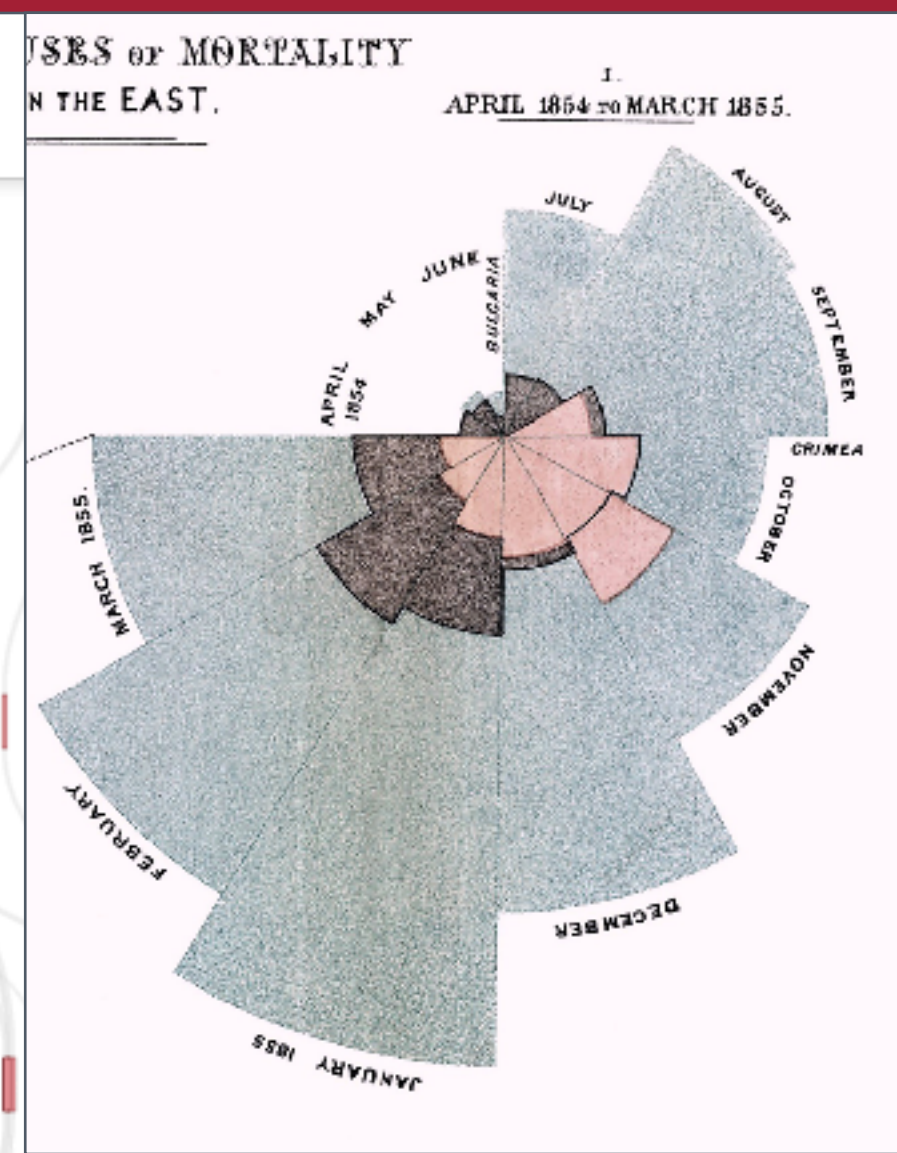
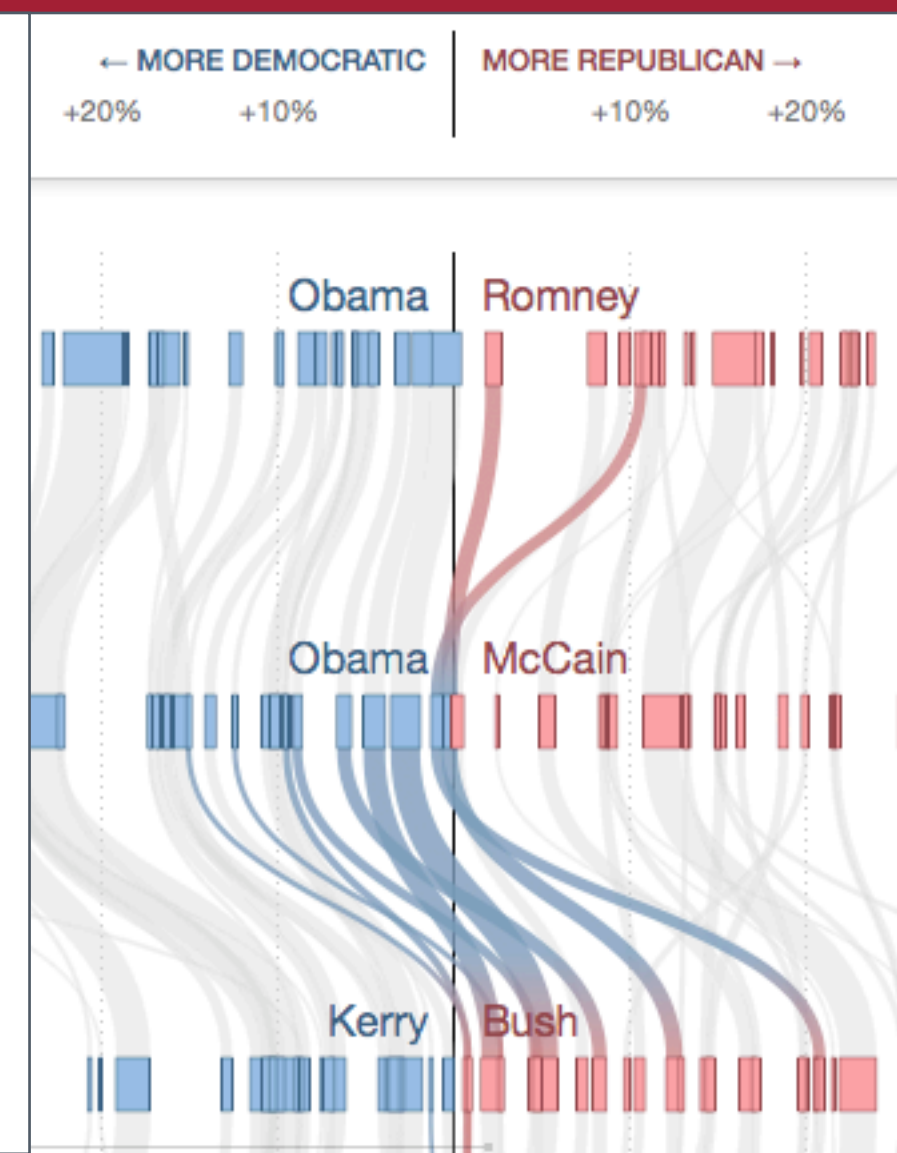
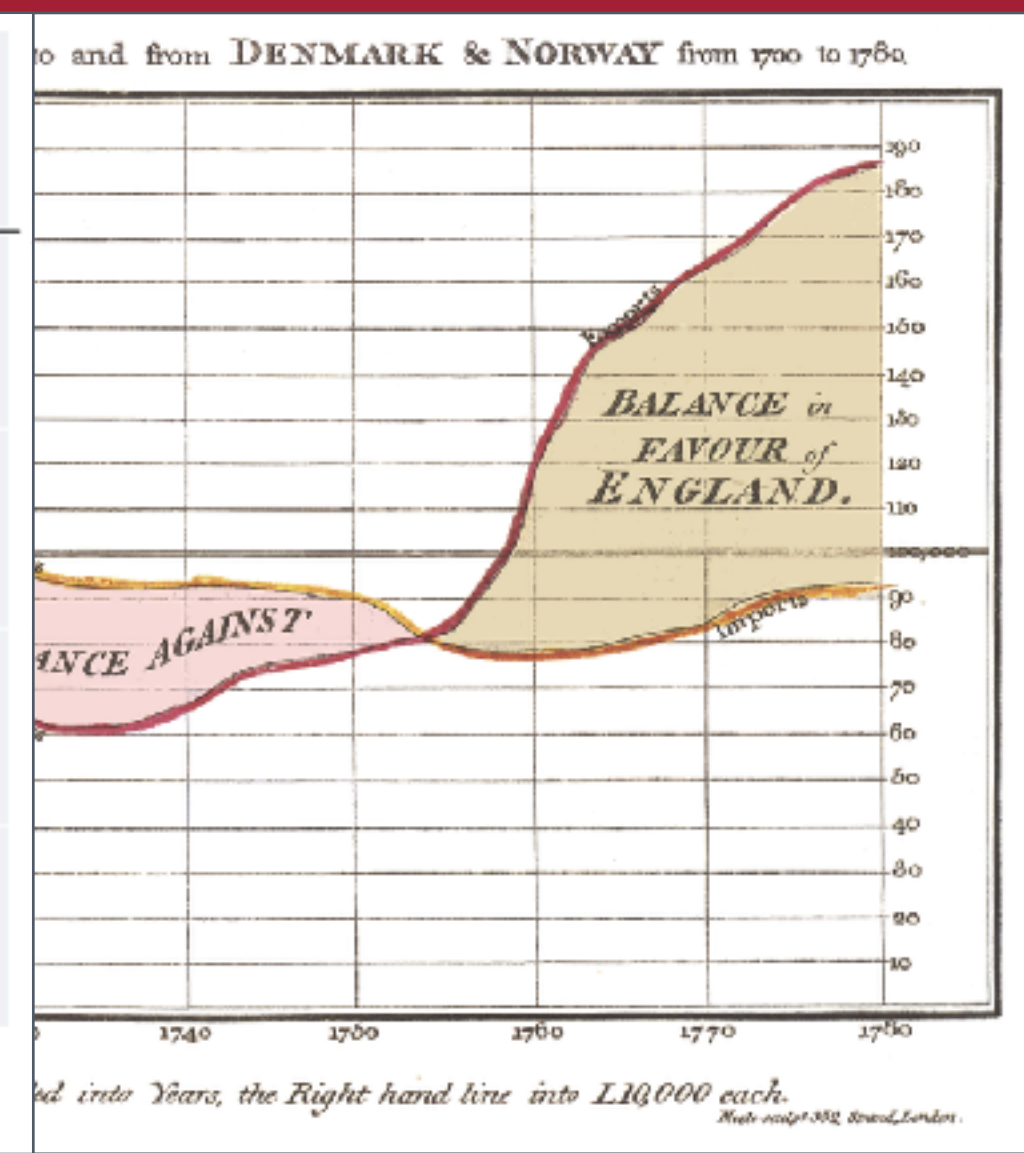
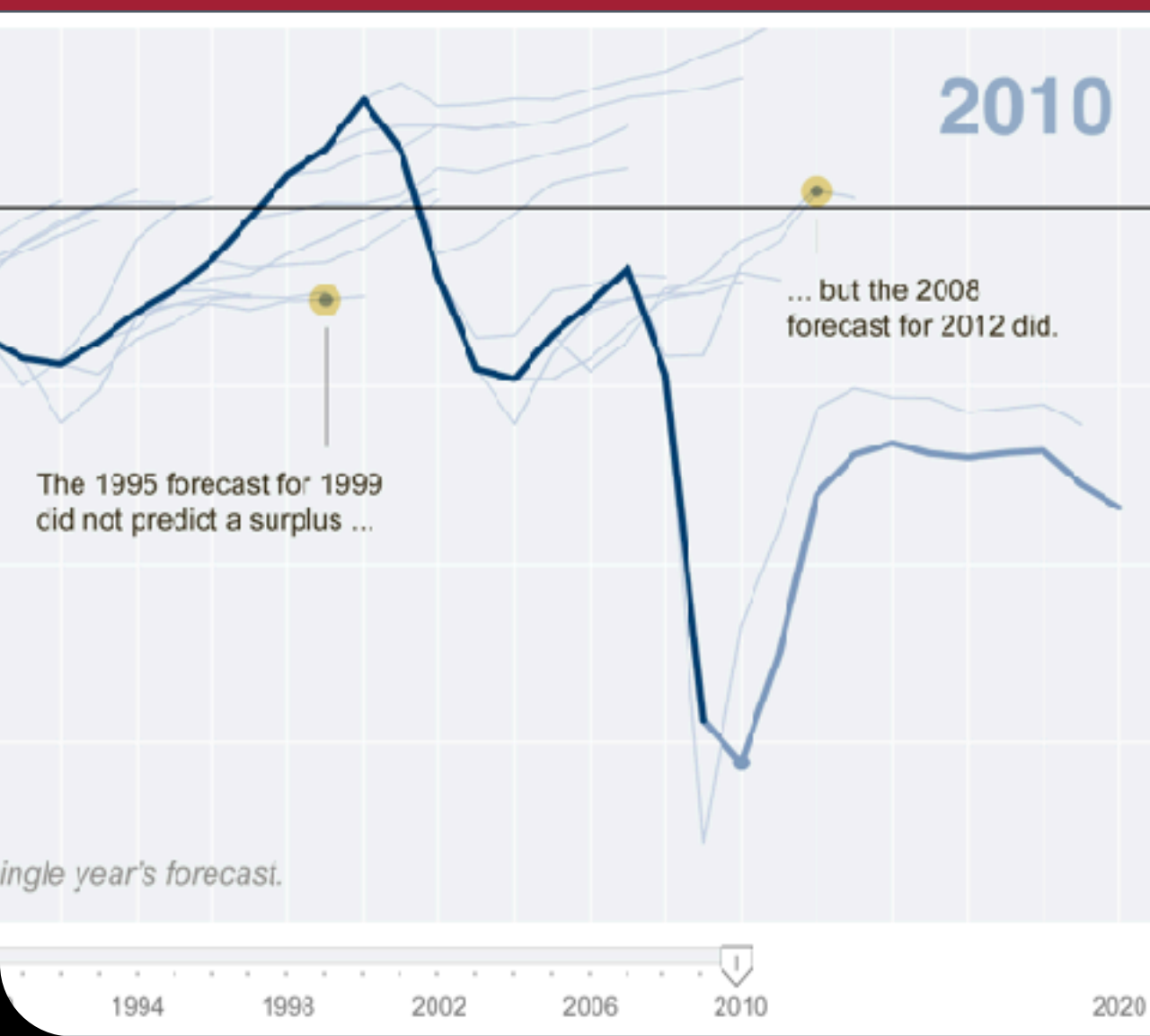


6.859: Interactive Data Visualization

Exploratory Data Analysis

Arvind Satyanarayan

Download data for today's activity:
www.yellkey.com/free



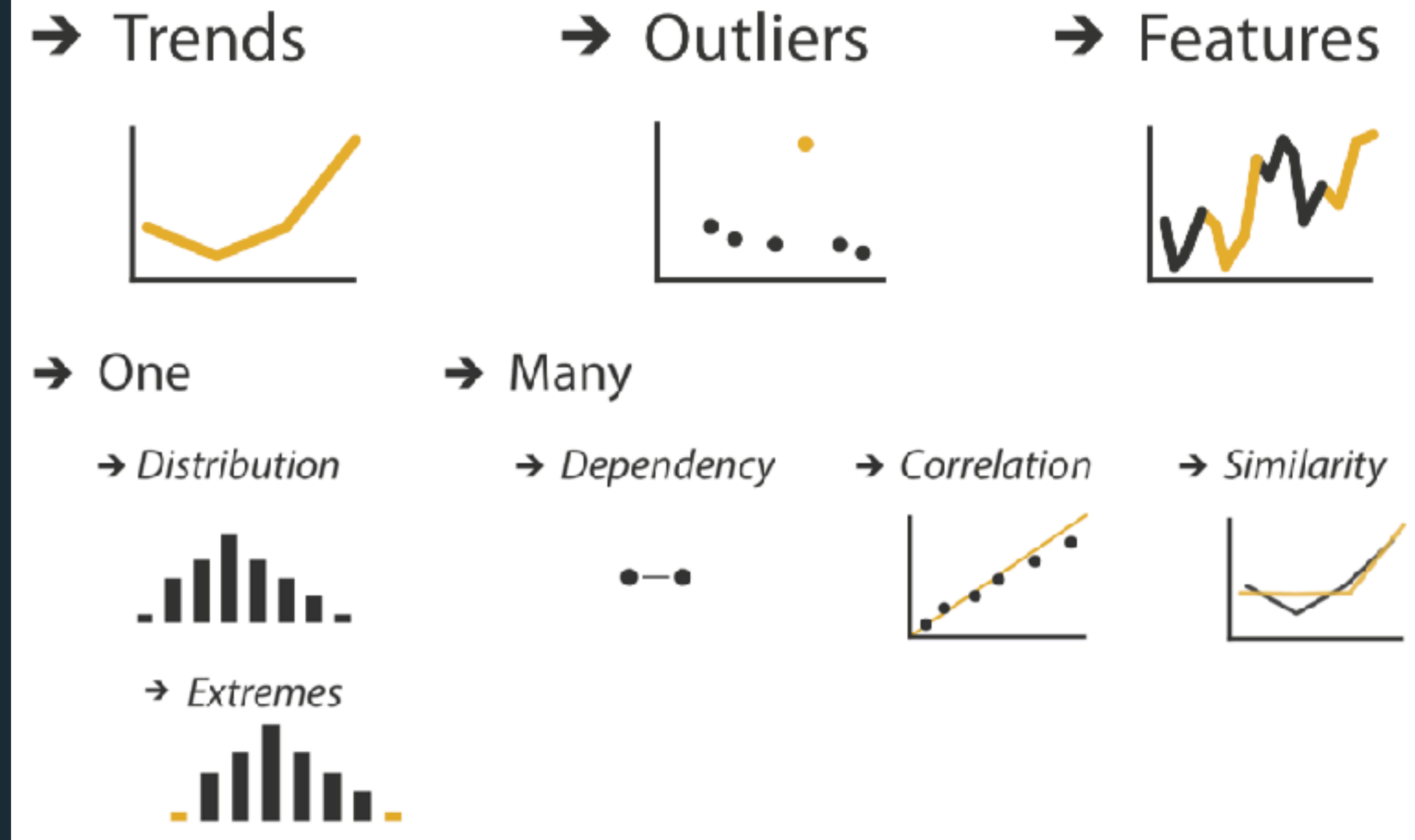
Exploratory Visual Analysis

Process

1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...

Lessons

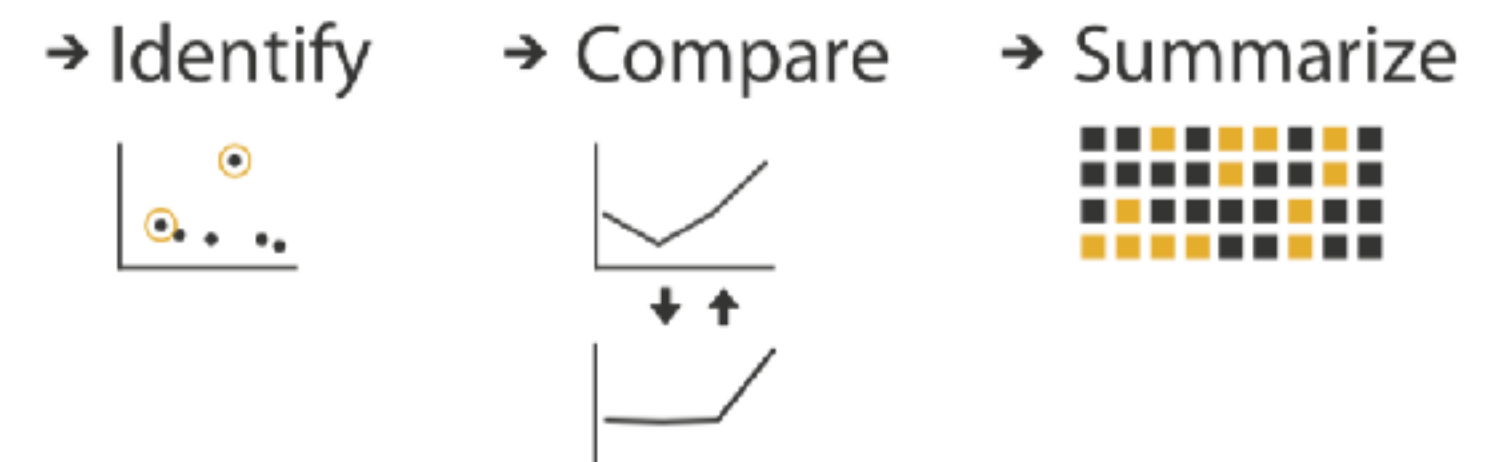
- ✓ Check **data quality** and your **assumptions**.
- ✓ Start with **univariate summaries**, then consider **relationships between variables**.



Search

	Target known	Target unknown
Location known	Lookup	Browse
Location unknown	Locate	Explore

Query



Analysis Example: Antibiotic Effectiveness

Analysis Example: Antibiotic Effectiveness

Collected prior to 1951

Genus of Bacteria String (N)

Species of Bacteria String (N)

Antibiotic Applied String (N)

Gram-Staining? Pos / Neg (N)

Min. Inhibitory Con. (g) Number (Q)

Analysis Example: Antibiotic Effectiveness

Collected prior to 1951

Genus of Bacteria

String (N)

Species of Bacteria

String (N)

Antibiotic Applied

String (N)

Gram-Staining?

Pos / Neg (N)

Min. Inhibitory Con. (g)

Number (Q)

Table 1—Burtin's Data

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus fecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

What questions might we ask?

Collected prior to 1951

Genus of Bacteria

String (N)

Species of Bacteria

String (N)

Antibiotic Applied

String (N)

Gram-Staining?

Pos / Neg (N)

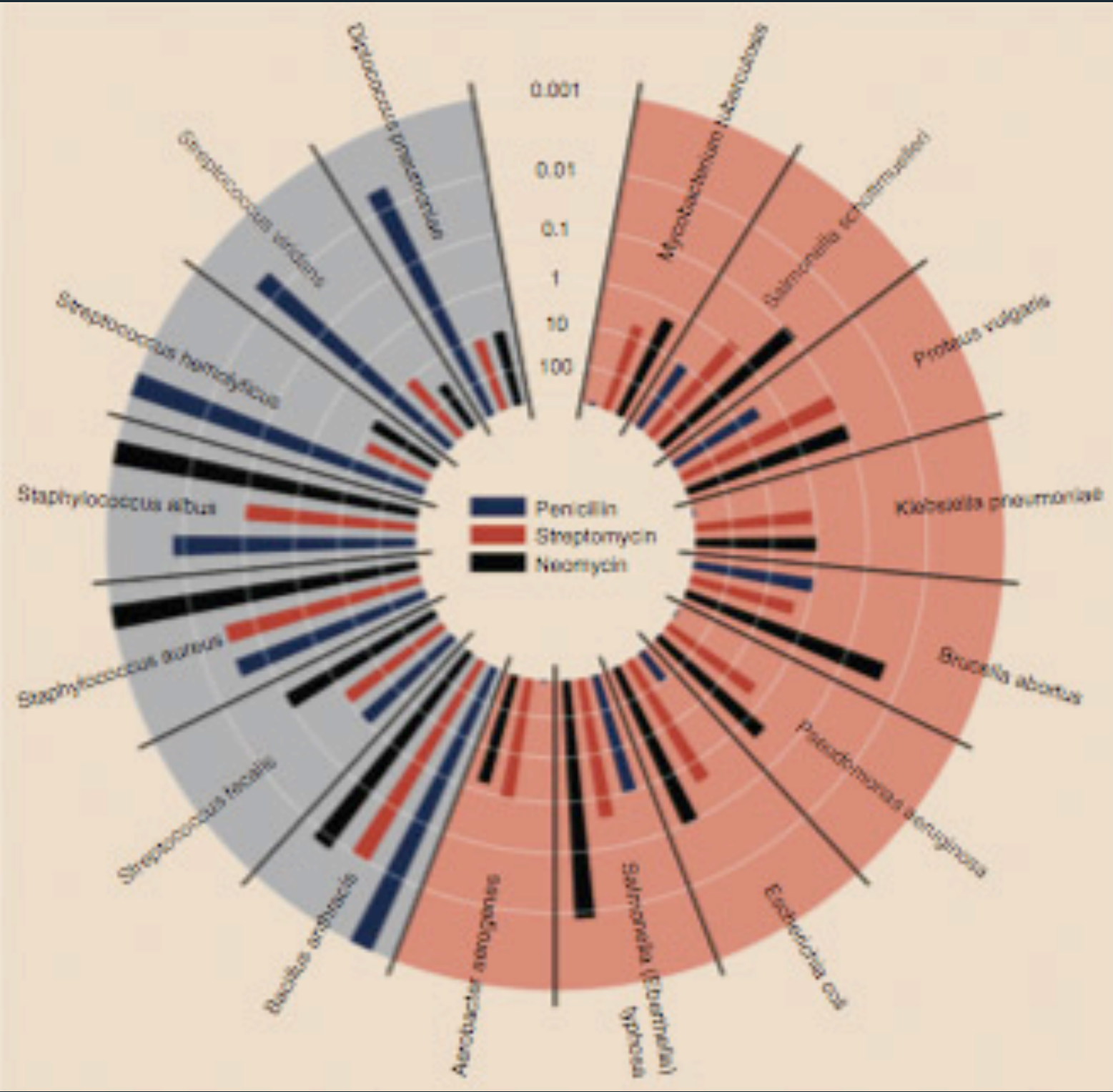
Min. Inhibitory Con. (g)

Number (Q)

Table 1—Burtin's Data

Bacteria	Antibiotic			Gram Staining
	Penicillin	Streptomycin	Neomycin	
<i>Aerobacter aerogenes</i>	870	1	1.6	negative
<i>Brucella abortus</i>	1	2	0.02	negative
<i>Brucella anthracis</i>	0.001	0.01	0.007	positive
<i>Diplococcus pneumoniae</i>	0.005	11	10	positive
<i>Escherichia coli</i>	100	0.4	0.1	negative
<i>Klebsiella pneumoniae</i>	850	1.2	1	negative
<i>Mycobacterium tuberculosis</i>	800	5	2	negative
<i>Proteus vulgaris</i>	3	0.1	0.1	negative
<i>Pseudomonas aeruginosa</i>	850	2	0.4	negative
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	negative
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	negative
<i>Staphylococcus albus</i>	0.007	0.1	0.001	positive
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	positive
<i>Streptococcus fecalis</i>	1	1	0.1	positive
<i>Streptococcus hemolyticus</i>	0.001	14	10	positive
<i>Streptococcus viridans</i>	0.005	10	40	positive

How do the drugs compare?

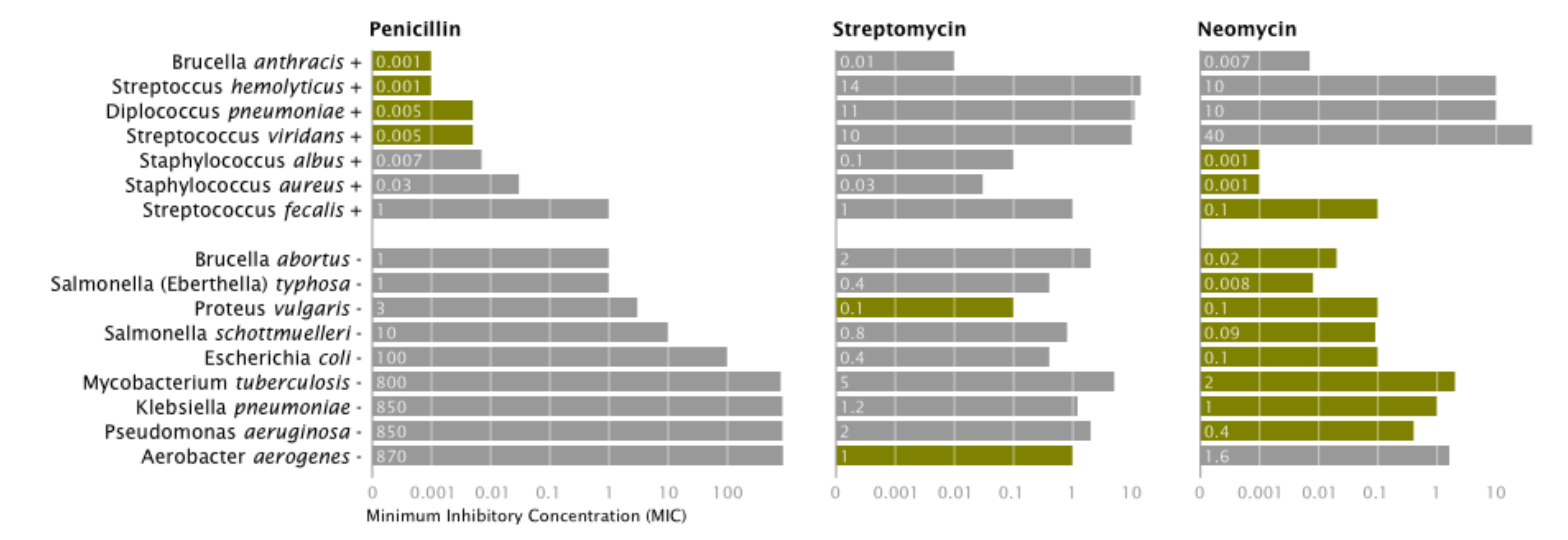


Bacteria	Penicillin	Antibiotic Streptomycin	Neomycin	Gram stain
<i>Aerobacter aerogenes</i>	870	1	1.6	-
<i>Brucella abortus</i>	1	2	0.02	-
<i>Bacillus anthracis</i>	0.001	0.01	0.007	+
<i>Diplococcus pneumoniae</i>	0.005	11	10	+
<i>Escherichia coli</i>	100	0.4	0.1	-
<i>Klebsiella pneumoniae</i>	850	1.2	1	-
<i>Mycobacterium tuberculosis</i>	800	5	2	-
<i>Proteus vulgaris</i>	3	0.1	0.1	-
<i>Pseudomonas aeruginosa</i>	850	2	0.4	-
<i>Salmonella (Eberthella) typhosa</i>	1	0.4	0.008	-
<i>Salmonella schottmuelleri</i>	10	0.8	0.09	-
<i>Staphylococcus albus</i>	0.007	0.1	0.001	+
<i>Staphylococcus aureus</i>	0.03	0.03	0.001	+
<i>Streptococcus fecalis</i>	1	1	0.1	+
<i>Streptococcus hemolyticus</i>	0.001	14	10	+
<i>Streptococcus viridans</i>	0.005	10	40	+

- Encodings**
- Radius:** $1 / \log(\text{MIC})$
- Bar Color:** Antibiotic
- Background Color:** Gram Staining

Original graphic by Will Burtin, 1951.

How do the drugs compare?



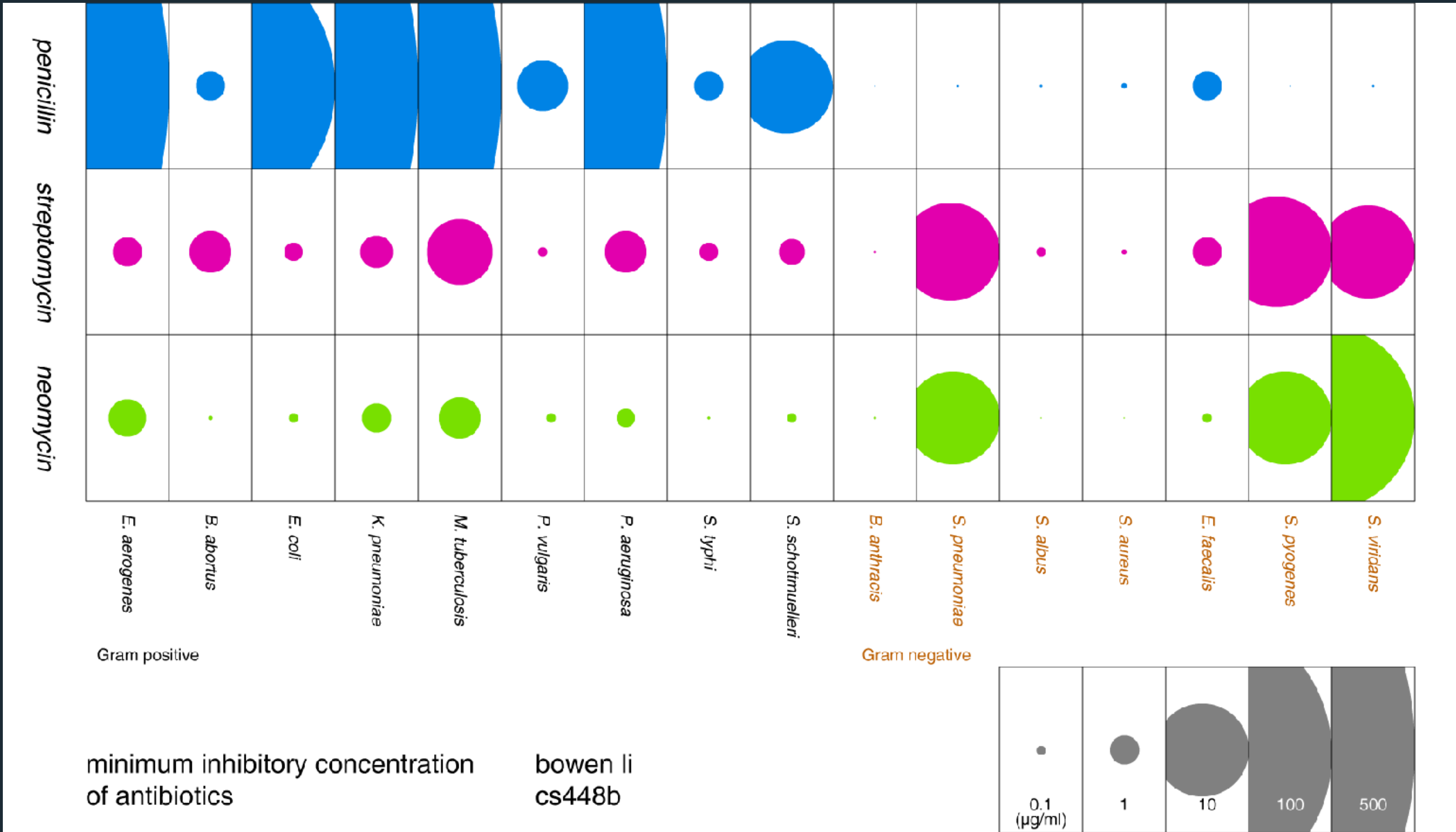
X-Axis: Antibiotic | log(MIC)

Y-Axis: Gram-Staining | Species

Color: Most Effective?

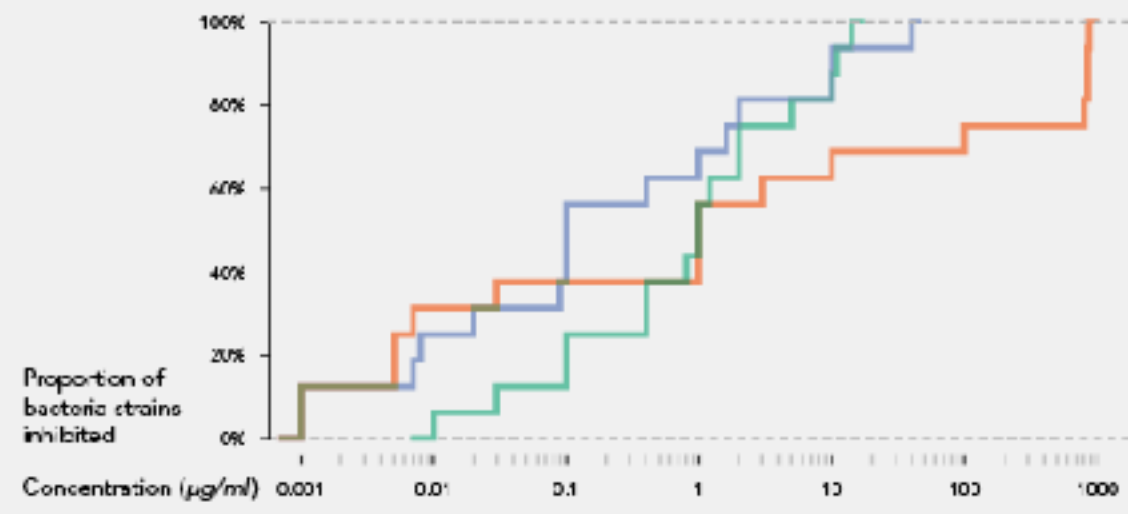
Mike Bostock, *Stanford CS448b* (Winter 2009).

How do the drugs compare?



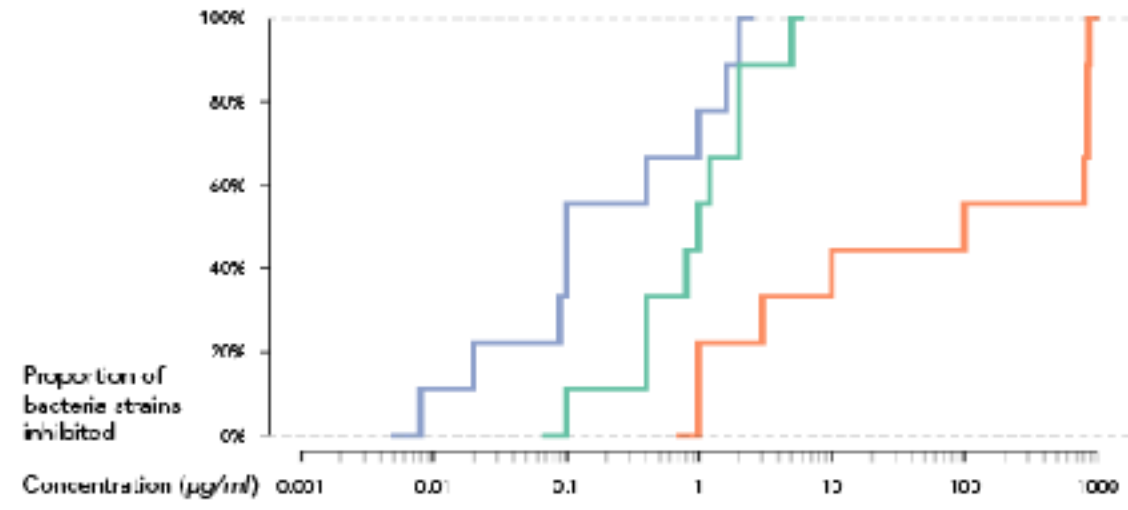
Bowen Li, Stanford CS448b (Fall 2009).

All bacteria



Streptomycin and Neomycin are more efficient broad-spectrum antibiotics than Penicillin.

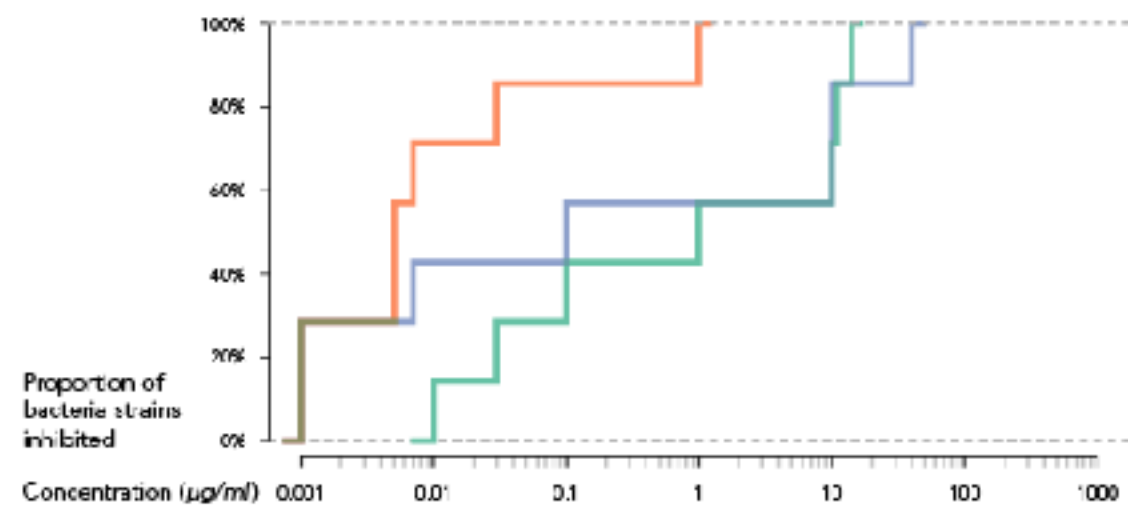
Gram-negative bacteria only



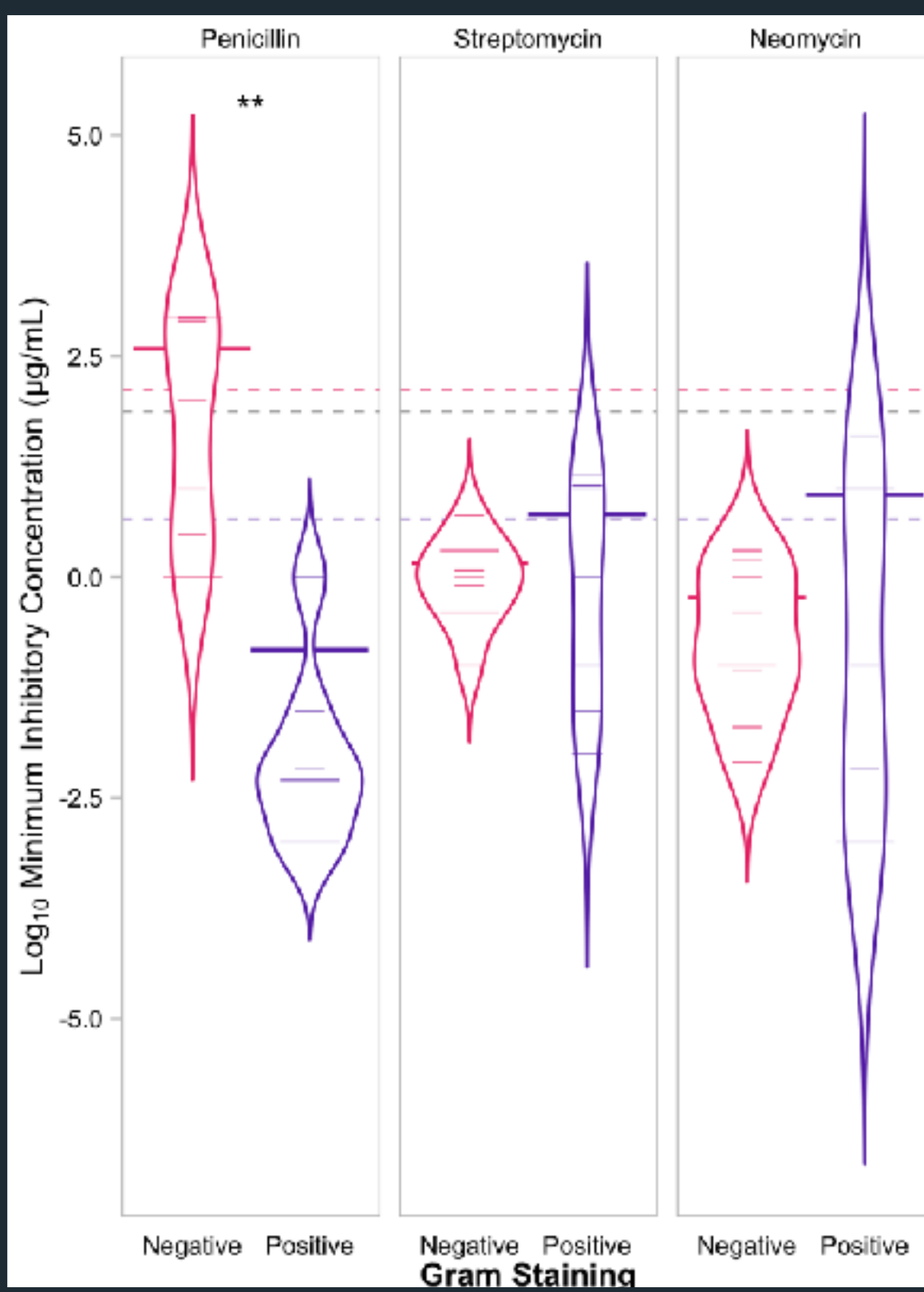
Neomycin and Streptomycin are more efficient against gram-negative bacteria, so can be used at a lower dosage here than above.

Gram staining quickly identifies bacteria as Gram-negative or Gram-positive, which can be used to find a more efficient antibiotic and dosage.

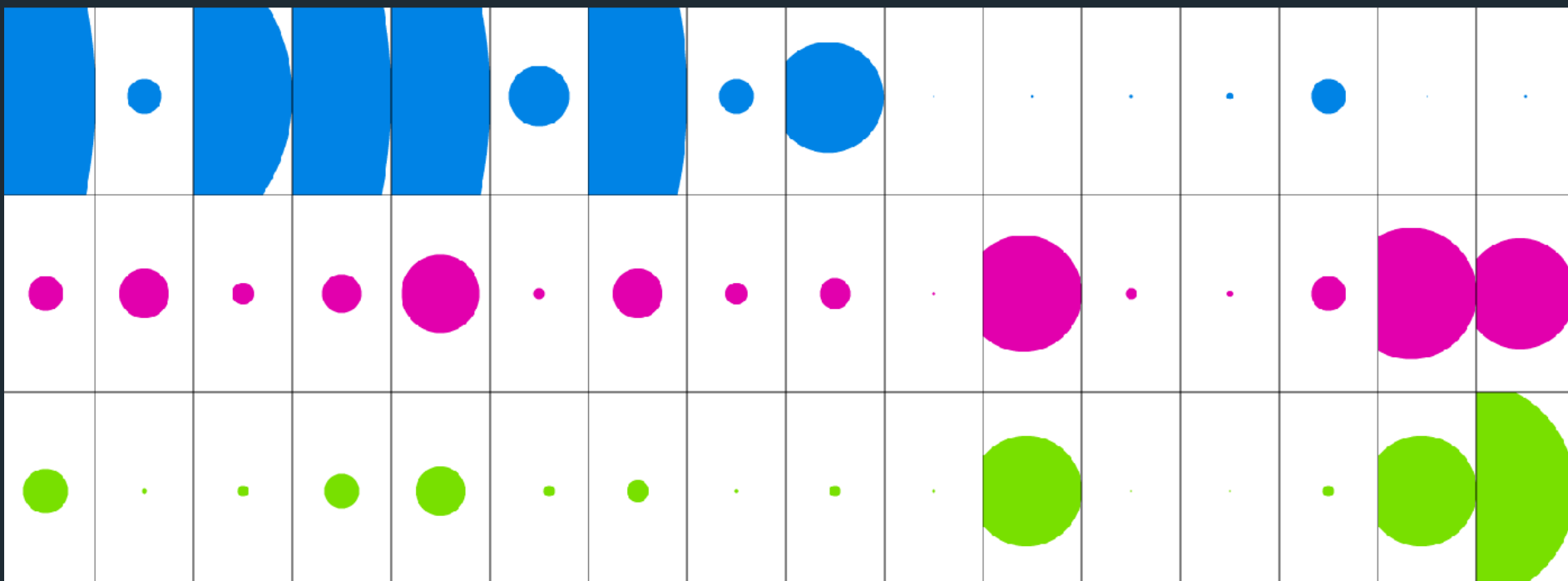
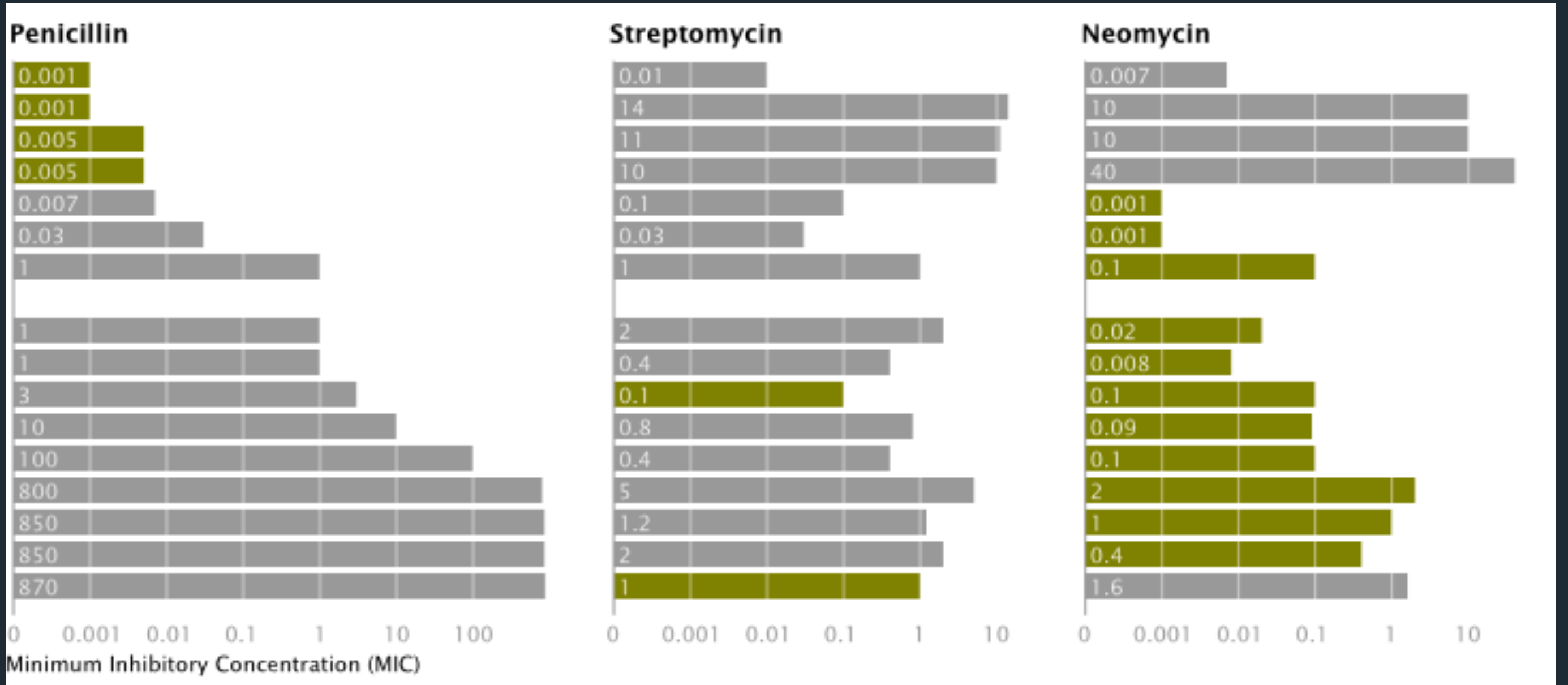
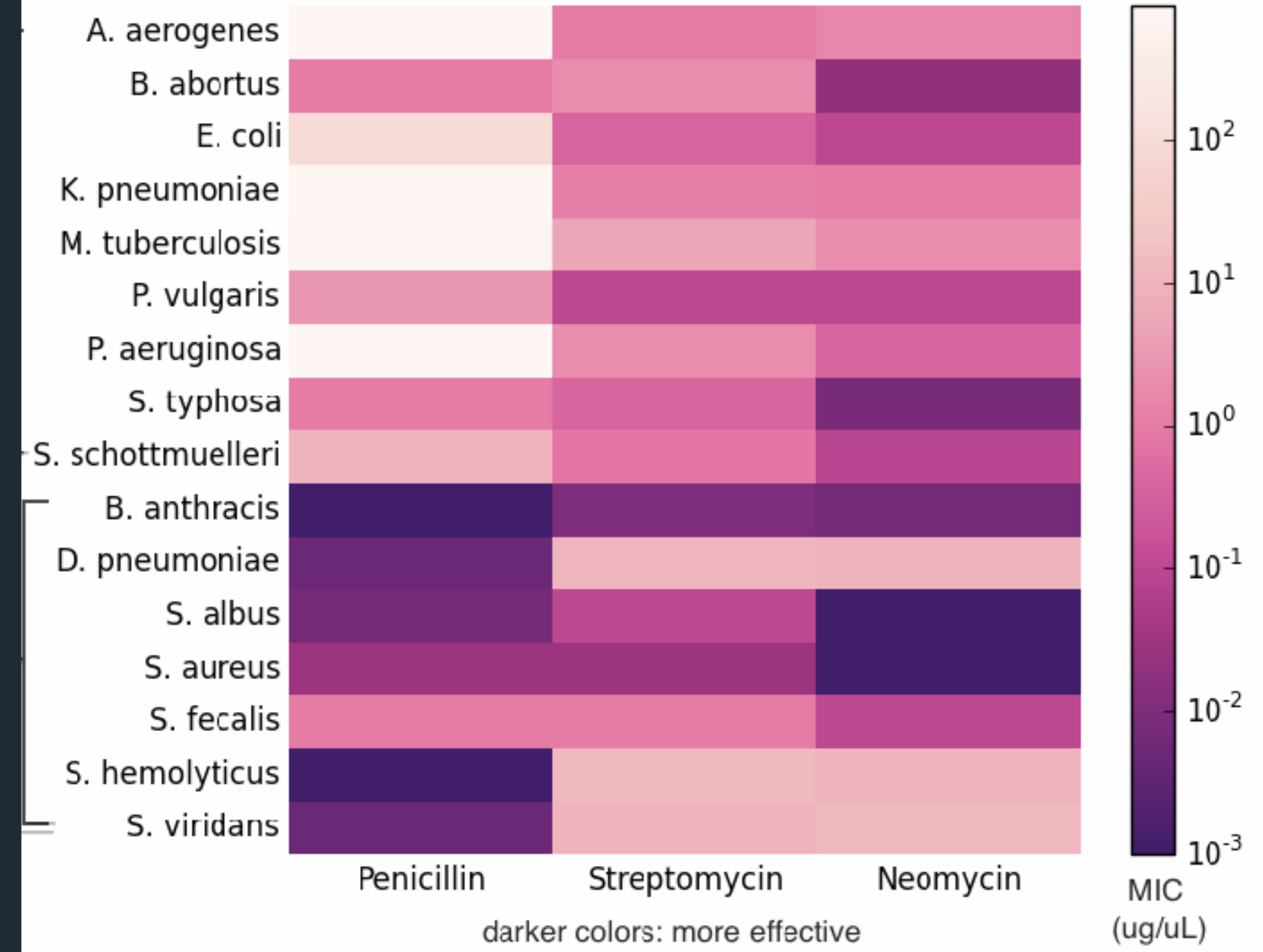
Gram-positive bacteria only

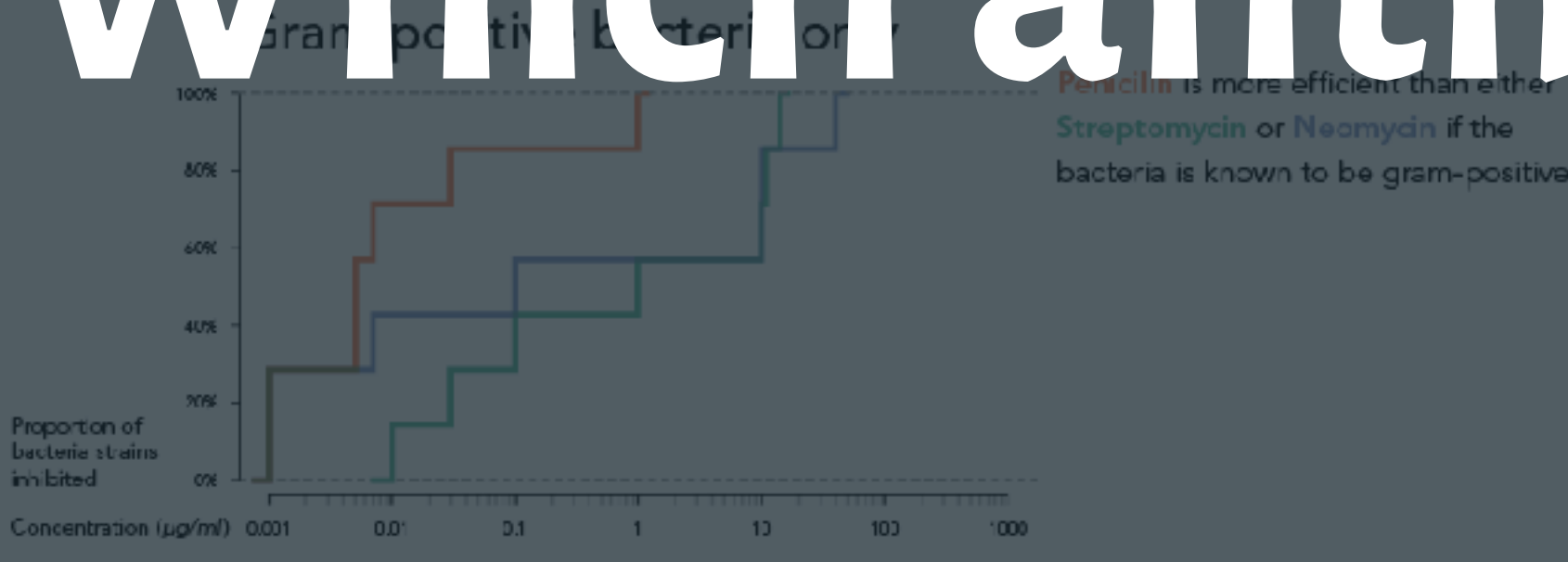
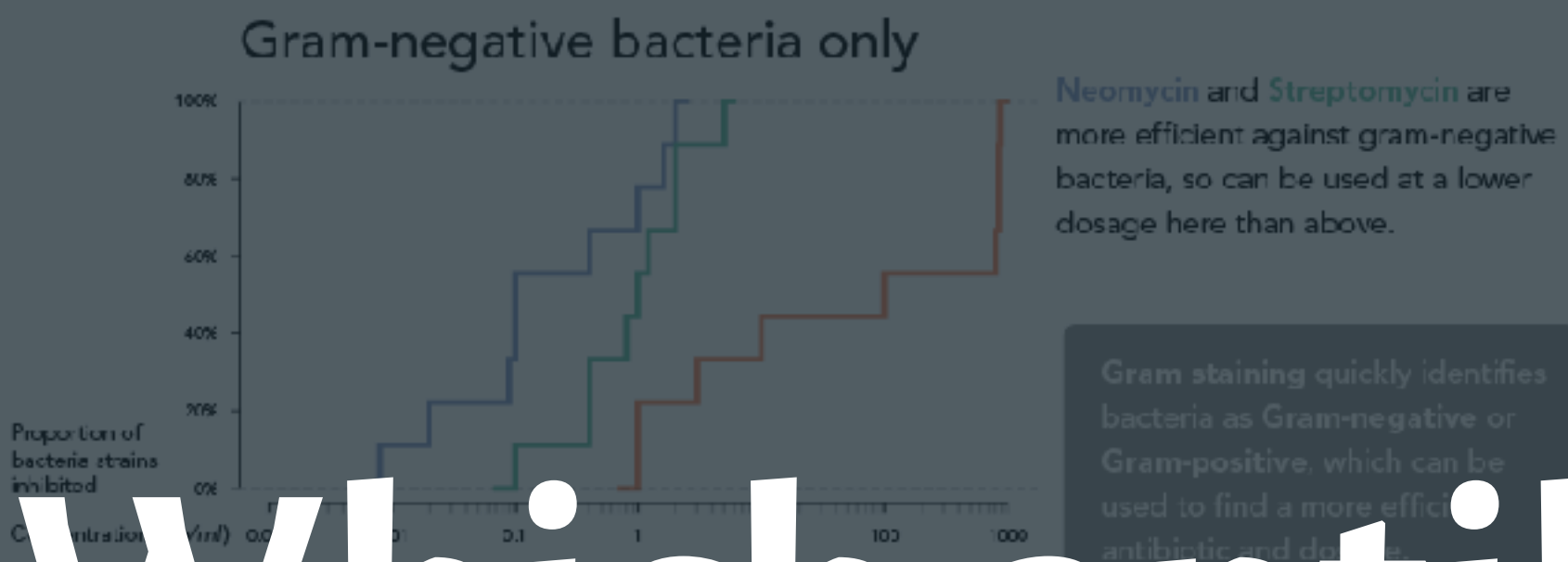
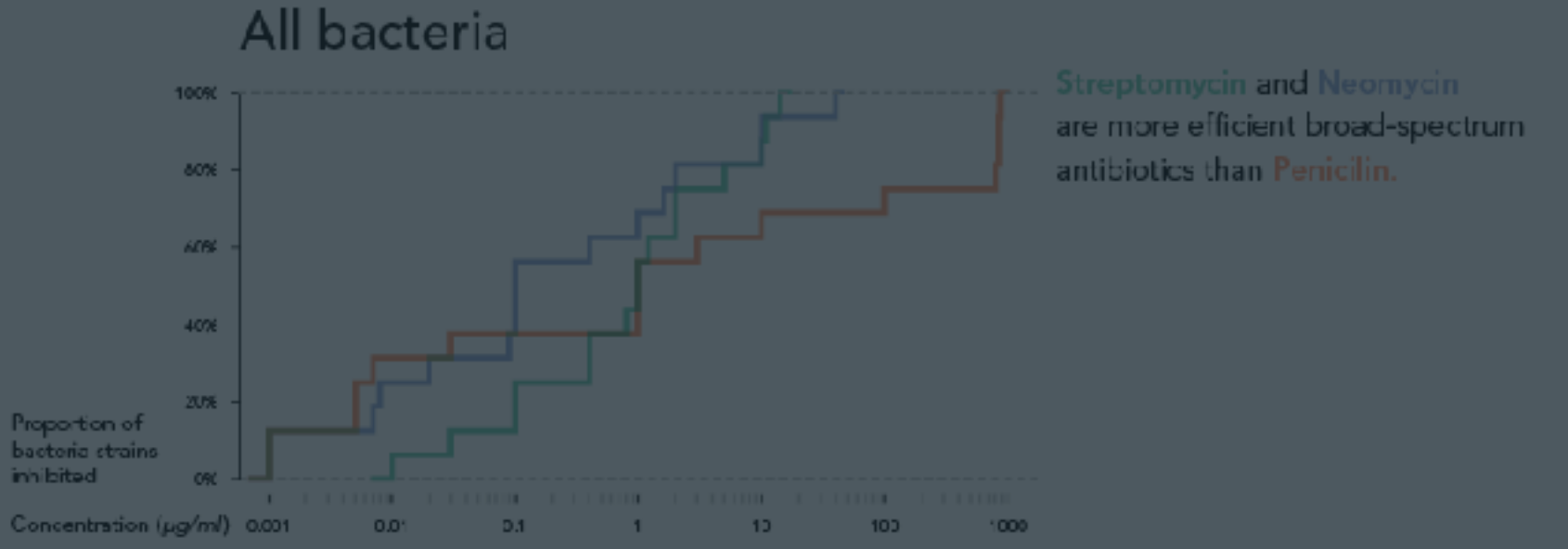


Penicillin is more efficient than either Streptomycin or Neomycin if the bacteria is known to be gram-positive.

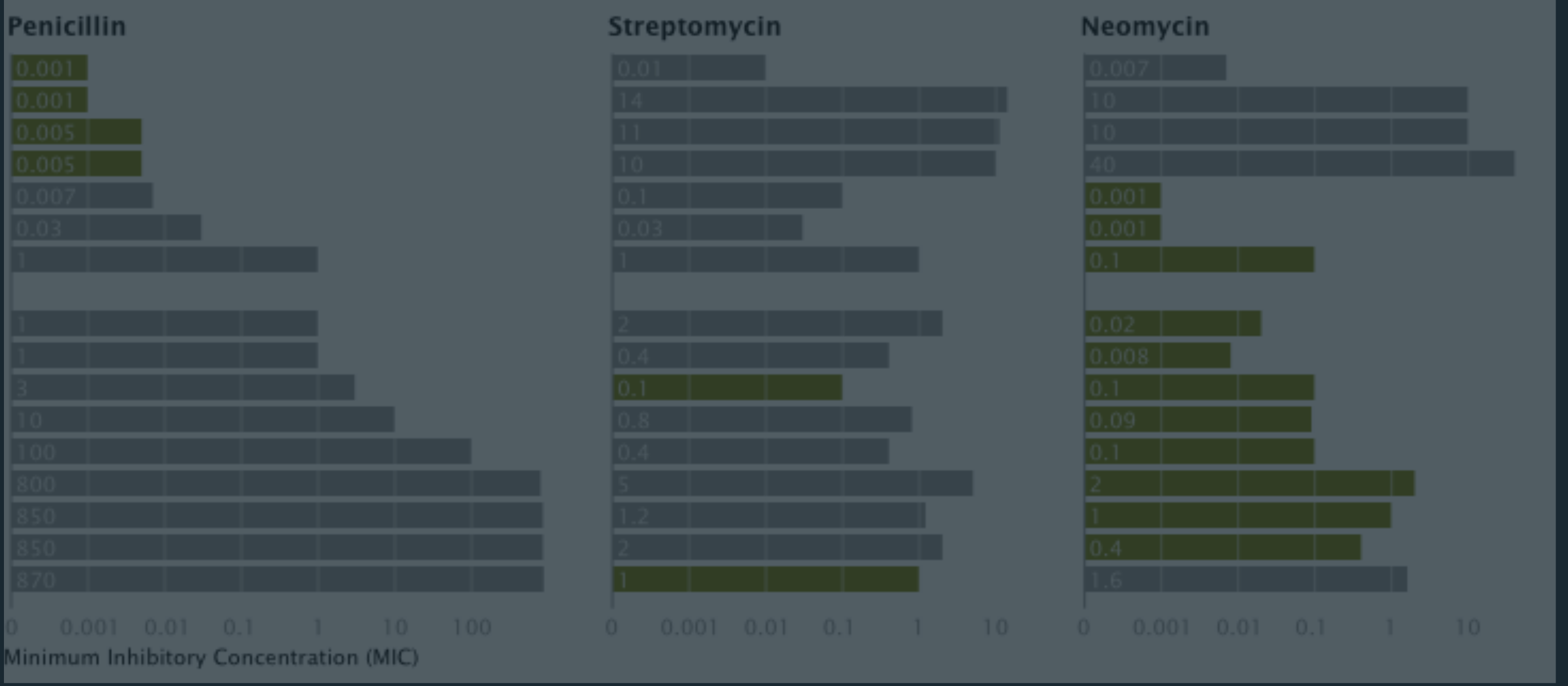
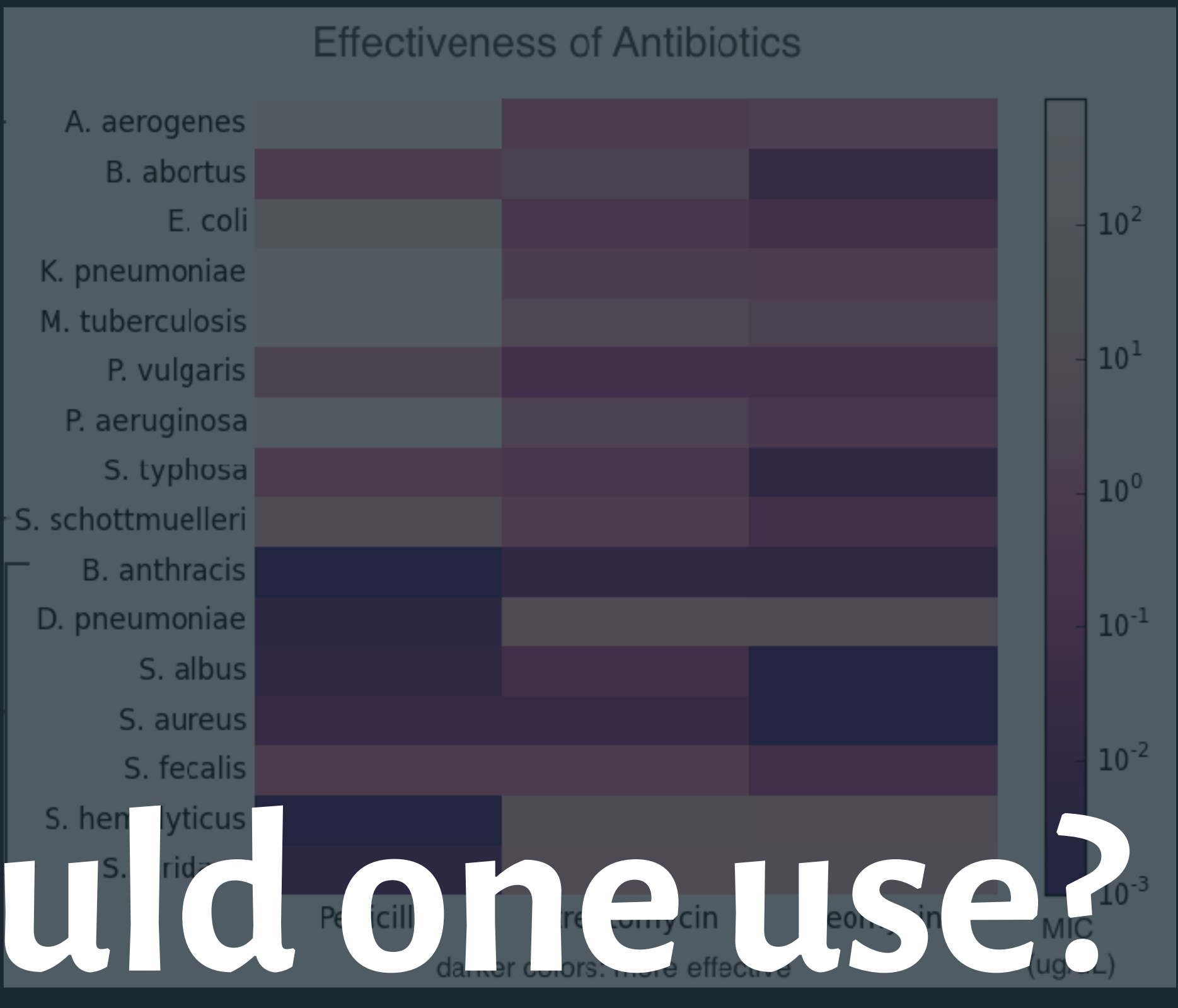
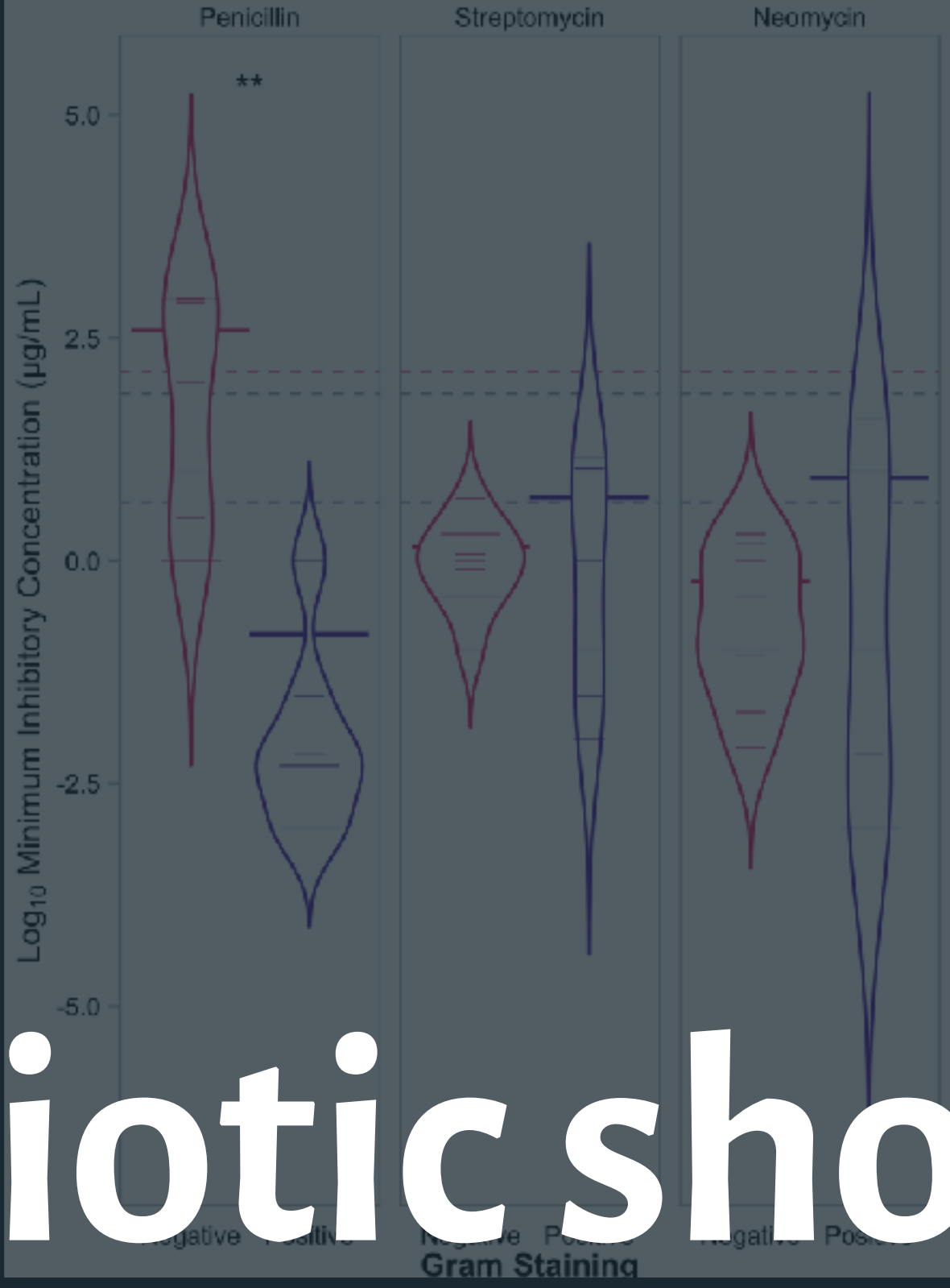
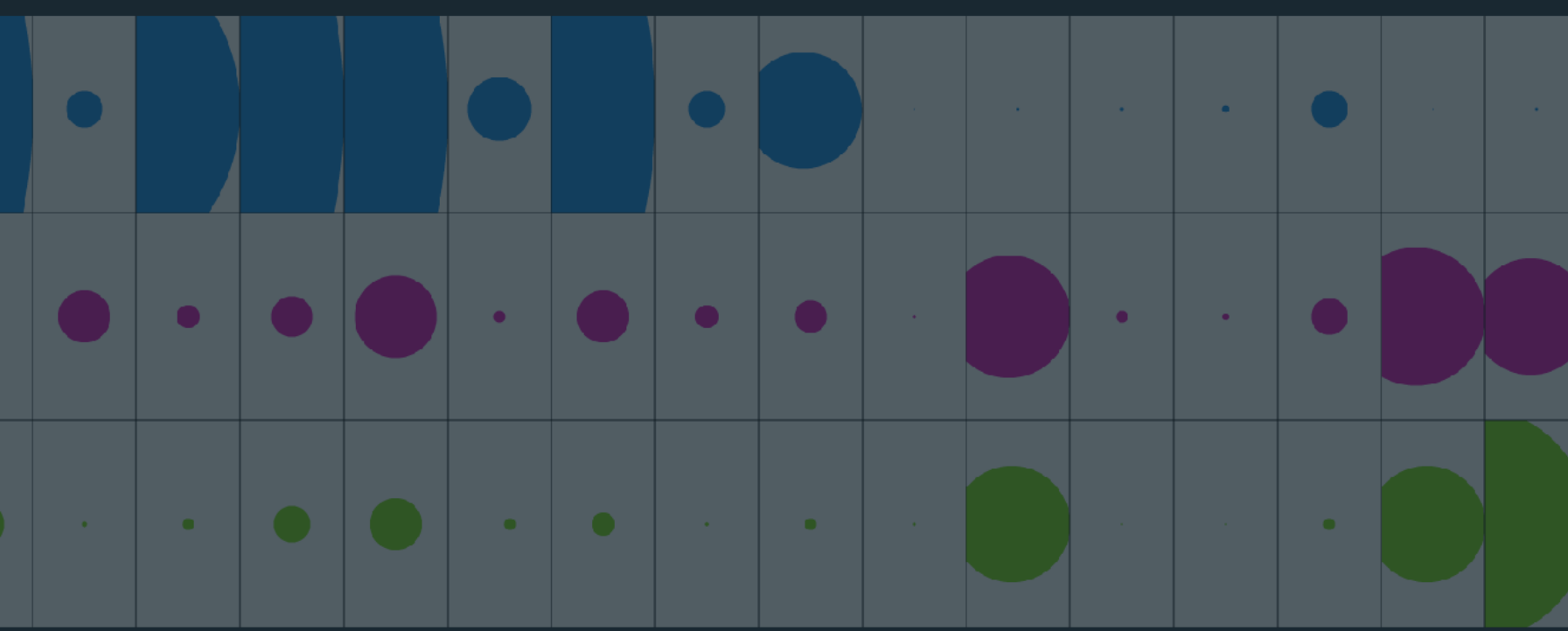


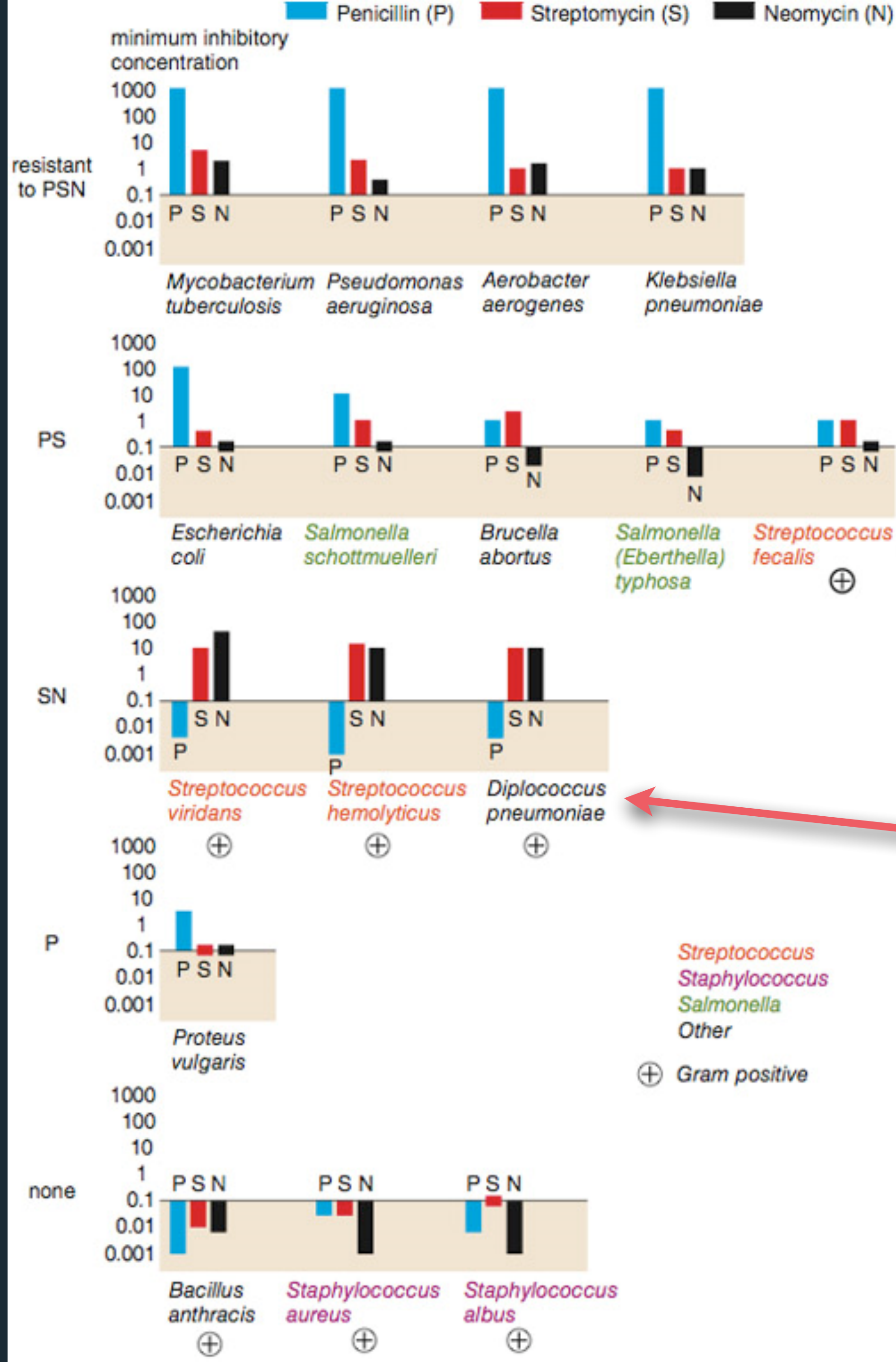
Effectiveness of Antibiotics





Which antibiotic should one use?





Do the bacteria group by antibiotic resistance?

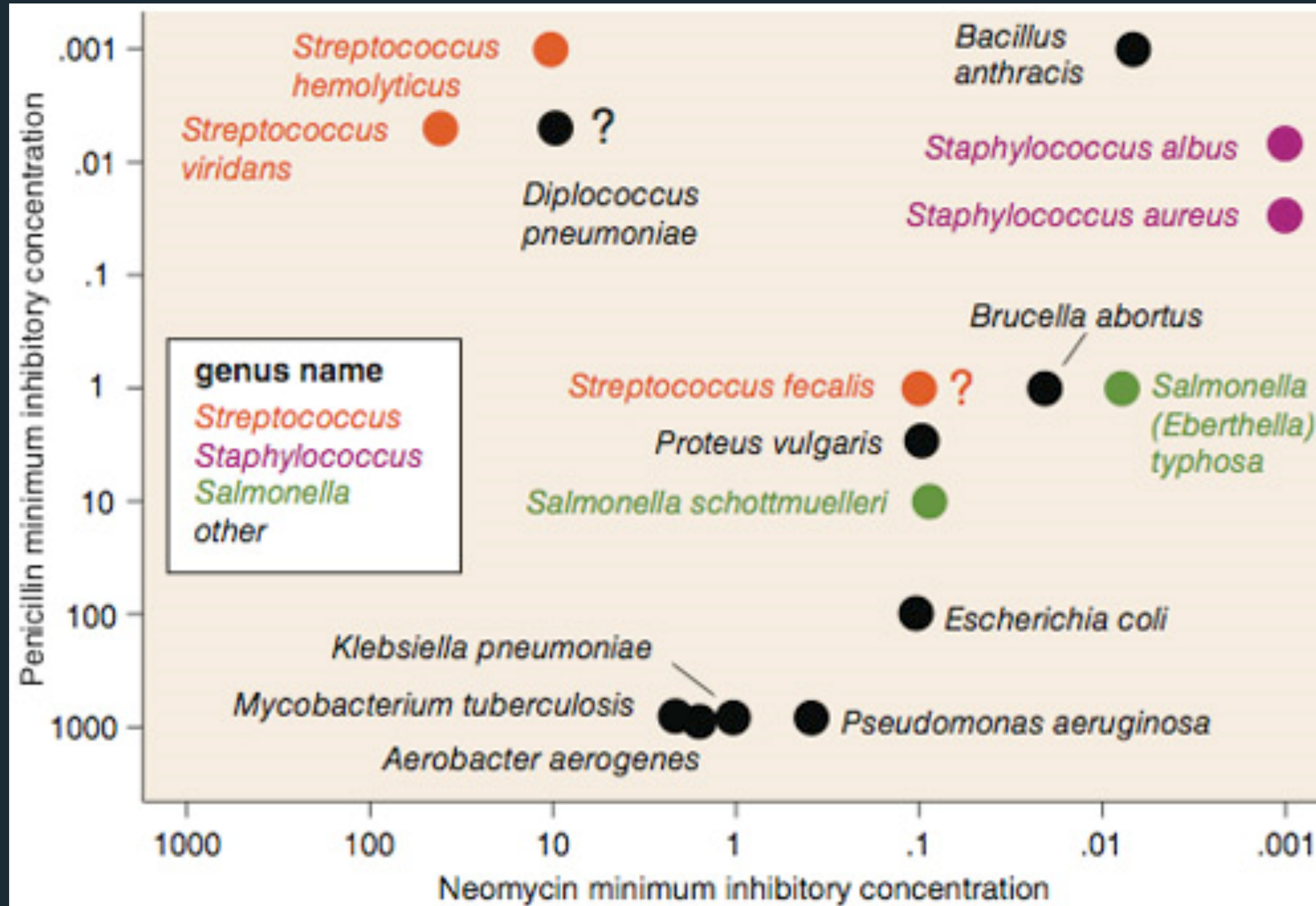
Not a streptococcus!
(realized ~30 yrs later)

Really a streptococcus!
(realized ~20 yrs later)

Do the bacteria group by resistance?

Do different drugs correlate?

Wainer & Lysen. American Scientist, 2009



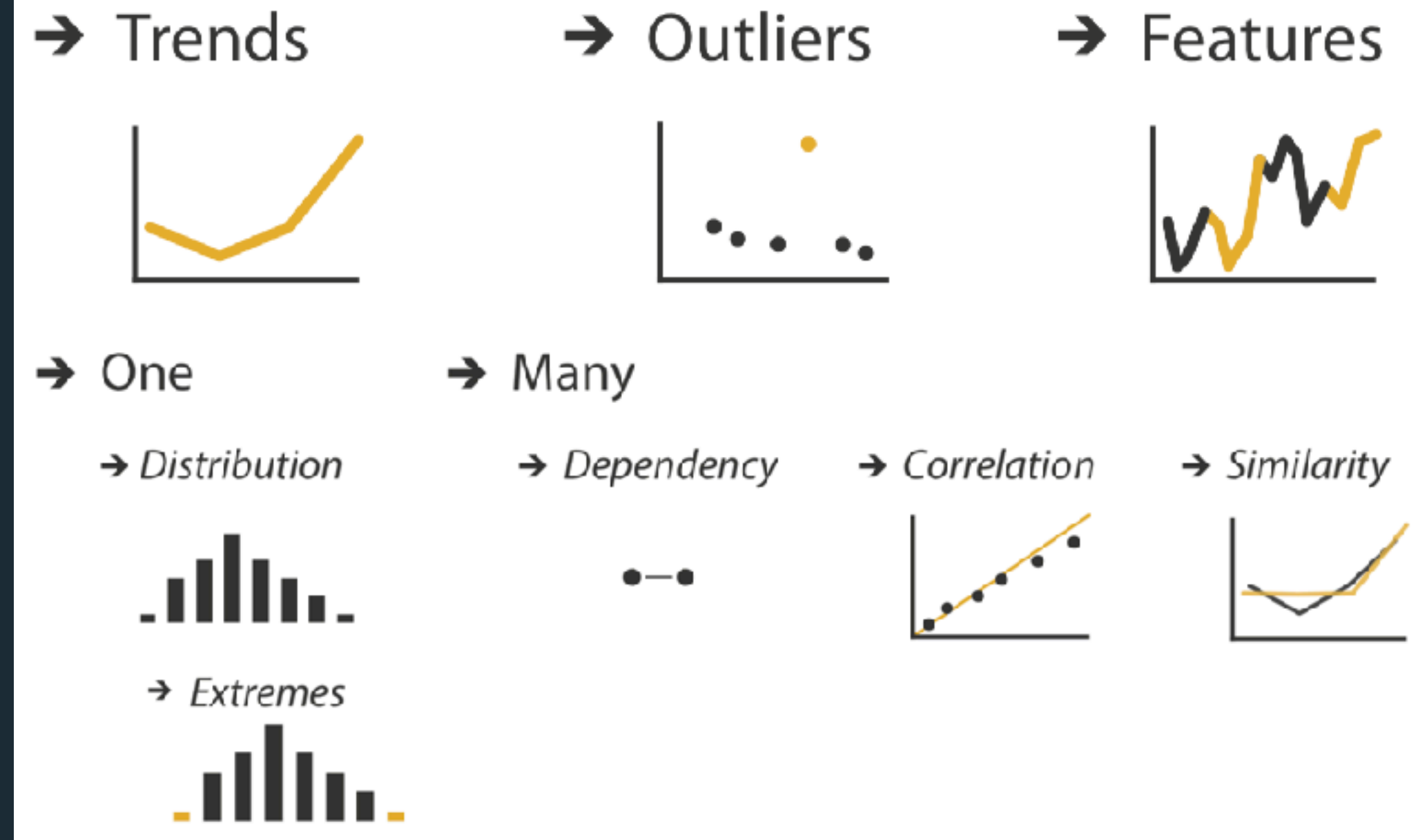
Exploratory Visual Analysis

Process

1. Construct graphics to address questions.
2. Inspect "answer" and ask new questions.
3. Iterate...

Lessons

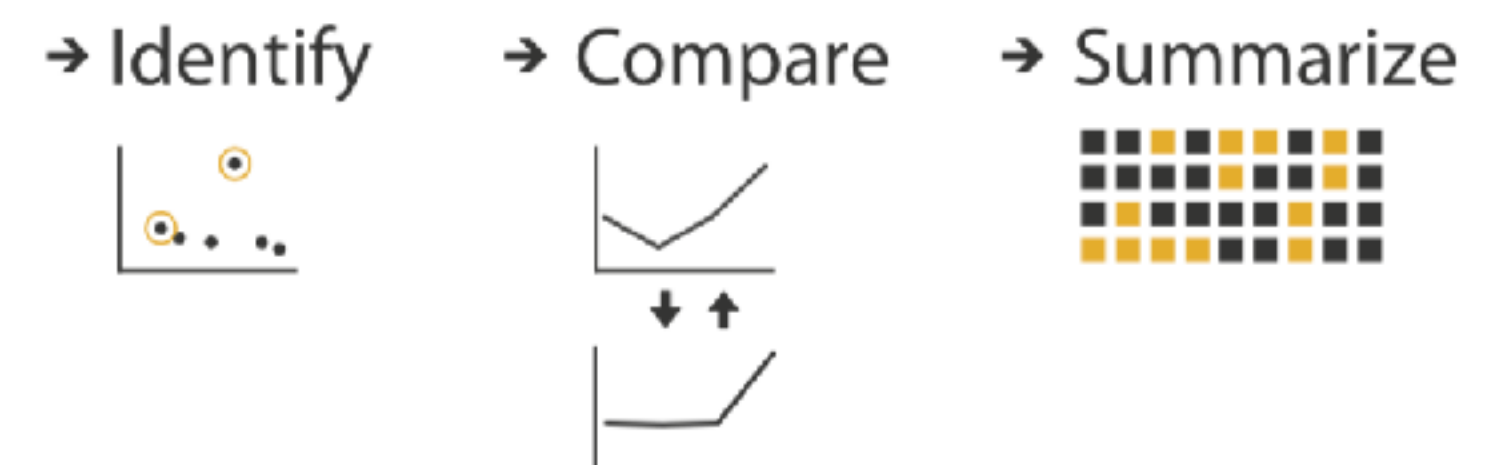
- ✓ Check **data quality** and your **assumptions**.
- ✓ Start with **univariate summaries**, then consider **relationships between variables**.
- ✓ Avoid **premature fixation**: balance **data variation** and **design variation**.



Search

	Target known	Target unknown
Location known	Lookup	Browse
Location unknown	Locate	Explore

Query



Is EDA/EVA fishing?

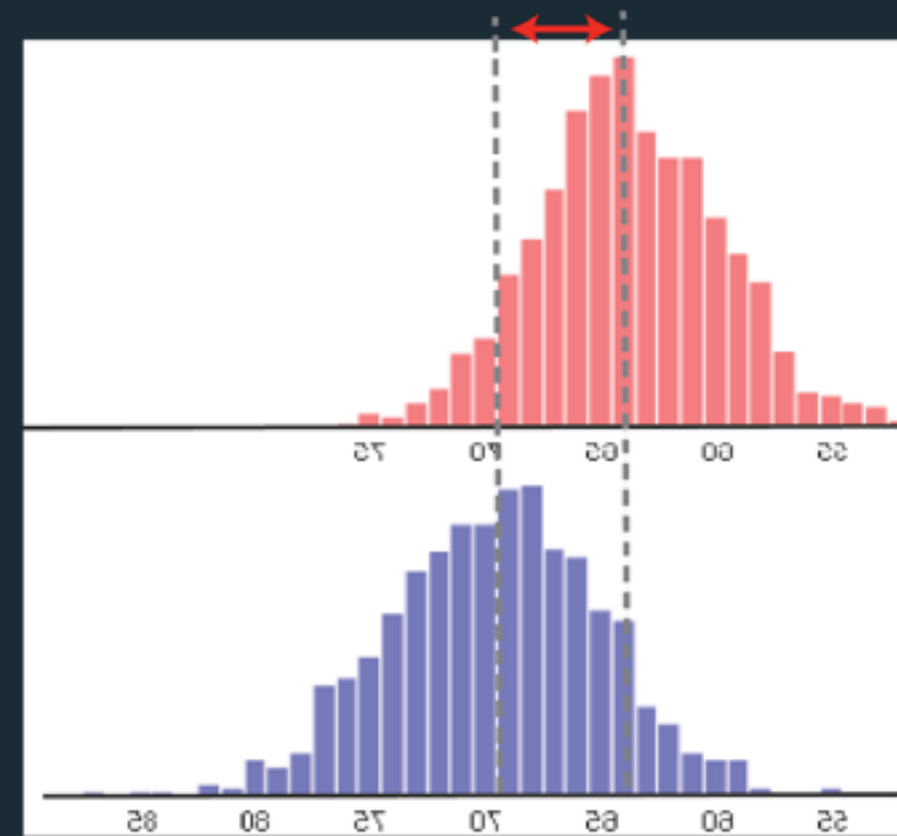


Is EDA/EVA Fishing? 🎣🐟

Some statisticians have proposed that when we look for patterns in visualizations, we're doing a series of *visual hypothesis tests*.

Multiple comparisons problem: the more hypothesis tests, the greater the chance of a curious finding (since the Null Hypothesis Significant Test admits 5% false positives).

No, because there's not a clear separation between *exploratory* and *confirmatory* analysis.

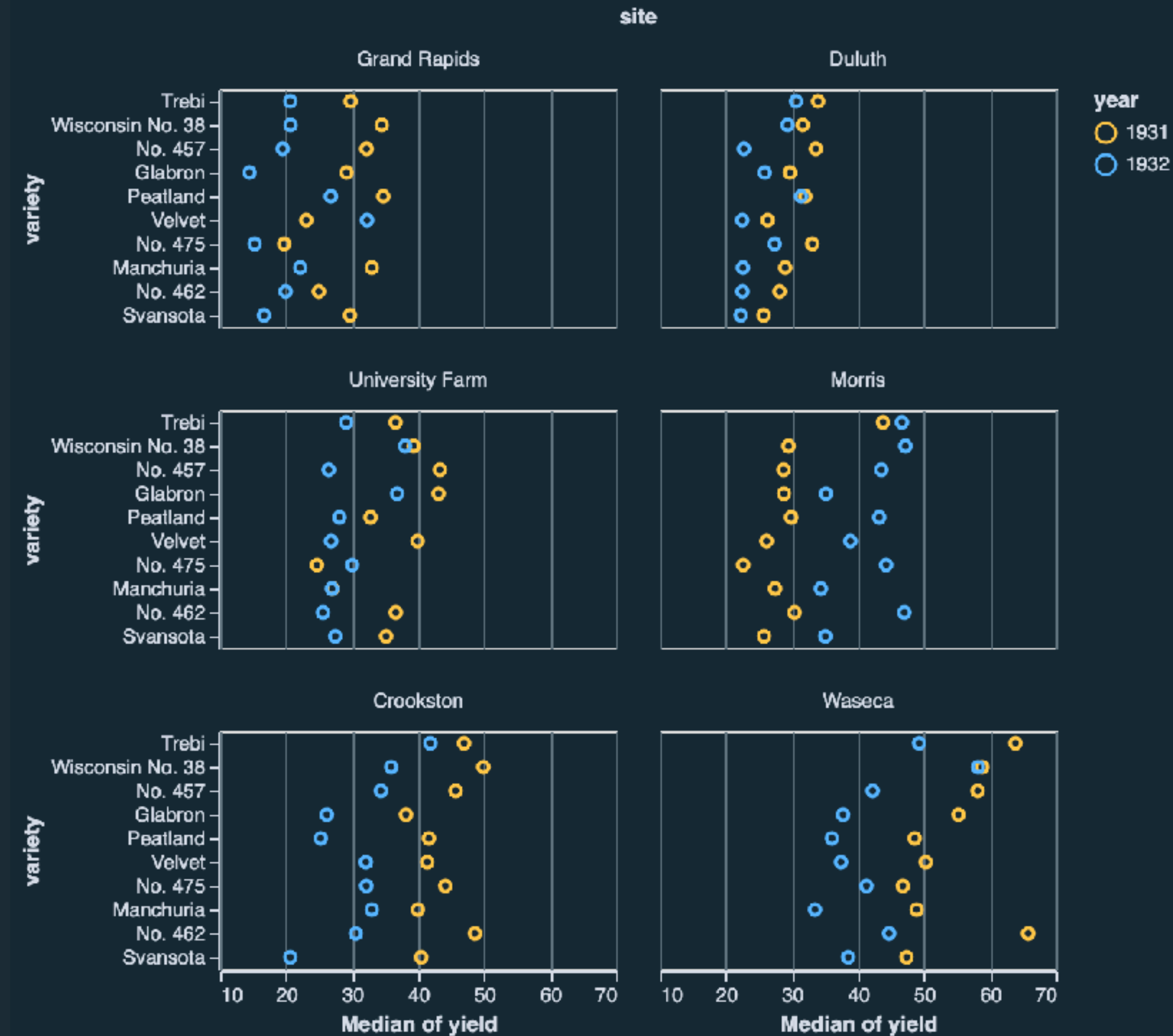


Is EDA/EVA Fishing? 🎣🐟

Some statisticians have proposed that when we look for patterns in visualizations, we're doing a series of *visual hypothesis tests*.

Multiple comparisons problem: the more hypothesis tests, the greater the chance of a spurious finding (since the Null Hypothesis Significant Test admits 5% false positives).

No, because there's not a clear separation between *exploratory* and *confirmatory* analysis.



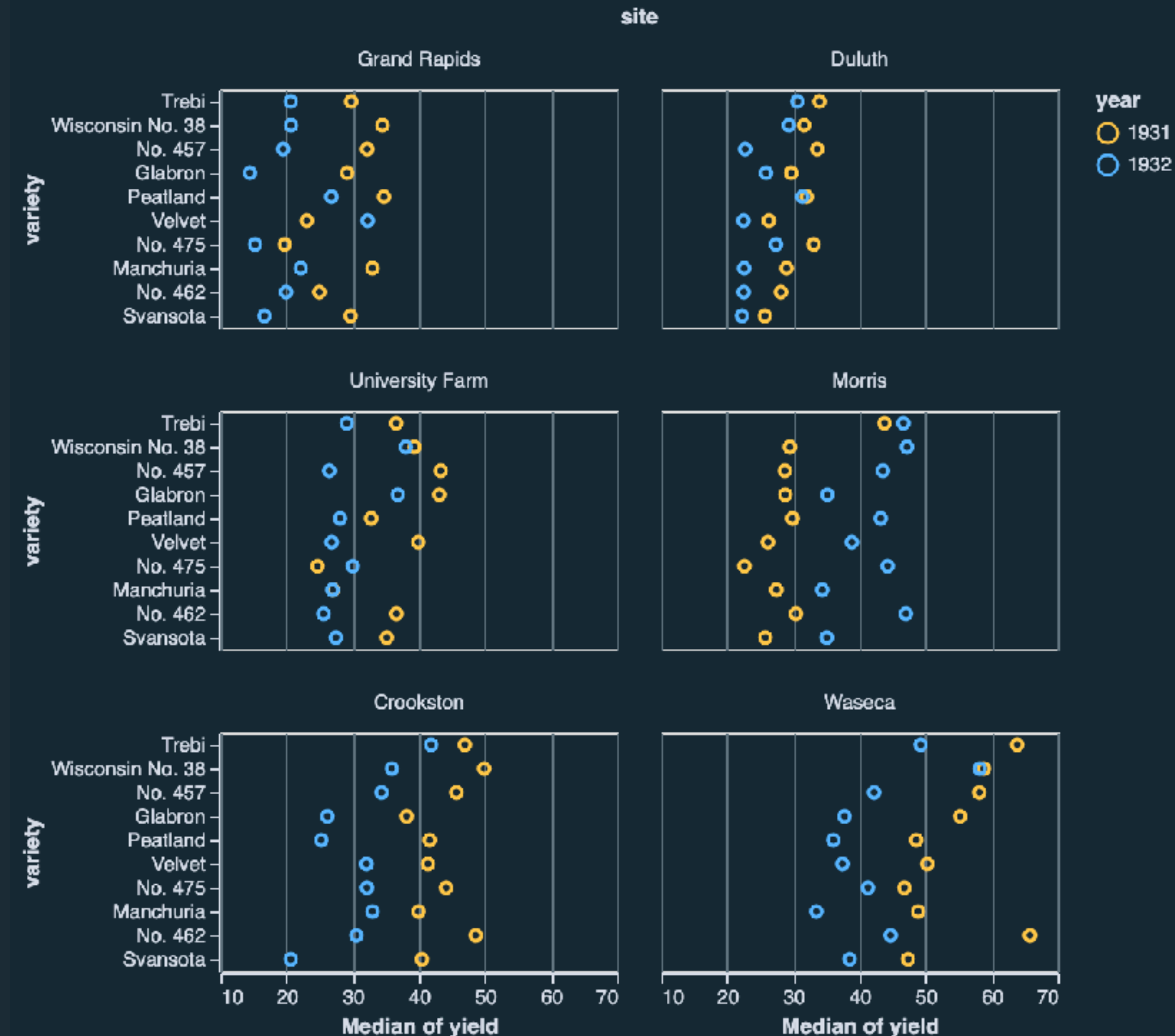
Is EDA/EVA Fishing? 🎣🐟

Some statisticians have proposed that when we look for patterns in visualizations, we're doing a series of *visual hypothesis tests*.

Multiple comparisons problem: the more hypothesis tests, the greater the chance of a spurious finding (since the Null Hypothesis Significant Test admits 5% false positives).

No, because there's not a clear separation between *exploratory* and *confirmatory* analysis.

Sort of if you use the same dataset to make hunches *and* then test them (i.e., you cannot collect more data). **Try a hold out set (e.g., separate training vs. test data).**



Based on slides by Jessica Hullman