













# Topics

Data Diagnostics Effectiveness of antibiotics Confirmatory analysis Graphical Inference Intro to Tableau



Bureau of Justice Si http://bjs.ojp.usdoj	atistics - Data Online i.gov/				
Reported crime in Alabama					
Year Population 2004 4525375 4029 2005 4548327 3900 2006 4599030 393 2007 4627851 3974 2008 4661900 4082	Property crime rate 9.3 987 2732.4 309.9 9 955.8 2656 289 7 968.9 2645.1 322.9 1.9 980.2 2687 307.7 1.9 1080.7 2712.6 288.6	Burglary rate	Larceny-theft rate	Motor vehicle theft rate	
Reported crime in A	laska				
Year Population 2004 657755 3371 2005 663253 3611 2006 670053 3583 2007 683478 3373 2008 686293 2926	Property crime rate 573.6 2456.7 340.6 6 622.8 2601 391 2 615.2 2588.5 378.3 8.9 538.9 2480 355.1 3.3 470.9 2219.9 237.5	Burglary rate	Larceny-theft rate	Motor vehicle theft rate	
Reported crime in Ar	izona				
Year Population 2004 5739879 5077 2005 5953007 4827 2006 6166318 4747 2007 6338755 4507 2008 6500180 4083	Property crime rate 3118.7 963.5 7 946.2 2958 922 .6 953 2874.1 914.4 2.6 935.4 2780.5 786.7 7.3 894.2 2605.3 587.8	Burglary rate	Larceny-theft rate	Motor vehicle theft rate	
Reported crime in Ar	'kansas				
Year Population 2004 2750000 403 2005 275708 4060 2006 2810872 4022 2007 2834797 394 2008 2855390 384	Property crime rate 1096.4 2699.7 237 3 1085.1 2720 262 .6 1154.4 2596.7 270.4 .5 1124.4 2574.6 246.5 3.7 1182.7 2433.4 227.6	Burglary rate	Larceny-theft rate	Motor vehicle theft rate	
Reported crime in California					
Year Population 2004 35842038 2005 36154147 2006 36457549 2007 36553215 2008 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	
Reported crime in Colorado					
Year Population 2004 4601821 3918	Property crime rate 3.5 717.3 2679.5 521.6	Burglary rate	Larceny-theft rate	Motor vehicle theft rate	

# Data "Wrangling"

One often needs to manipulate data prior to analysis. Tasks include reformatting, cleaning, quality assessment, and integration

#### Some approaches:

Writing custom scripts Manual manipulation in spreadsheets Data Wrangler: <u>http://vis.stanford.edu/wrangler</u> Google Refine: <u>http://code.google.com/p/google-refine</u>

#### How to gauge the quality of a visualization?

"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

200	Age:	95			
200	Sex:	Female		Offensi	
T	Race:	Caucasian			
<u>_</u> -	County (Res):	Prince Georges			
•	Zip Code (Res):	20770			
	Received:	940706		The second second	
	Complaint Sequence:	1		Count	
	Source:	Citizen		County	
	Reason:	Delinquent		1.4.14	
	Alleged Offense:	HARAS			
	Offense Level:	2 - Misdemeanc			
	County (Off):	Prince Georges		A	
	Zip Code (Off):	20770		Area:	
	Area:	V			
	Office:	71610			
	Intake Decision Date:	940729		Contraction of the	
	Intake Decision:	Closed		Office:	
	Days to ID:	23			
	Court Finding:	NONE			
	Disposition Date:	0			
	Disposition:			1. Section	
			l de la companya de l	Intake	
0 10 20 30	40 50 60	70 80	90		
				тс	
Query Result: 4792 out of 4792 (100%)					









## **Visualize Friends by School?**

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

# **Data Quality & Usability Hurdles**

Missing Data Erroneous Values Type Conversion Entity Resolution Data Integration

no measurements, redacted, ...? misspelling, outliers, ...? e.g., zip code to lat-lon diff. values for the same thing? effort/errors when combining data

*LESSON:* Anticipate problems with your data. Many research problems around these issues!

Exploratory Analysis: Effectiveness of Antibiotics

# What questions might we ask?

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

# The Data Set

Genus of Bacteria	
Species of Bacteria	
Antibiotic Applied	
Gram-Staining?	
Min. Inhibitory Concent.	<b>(g</b> )

String String String Pos / Neg Number

Collected prior to 1951

# Will Burtin, 1951



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



### How do the drugs compare?







# Confirmatory Data Analysis

# **Some Uses of Formal Statistics**

What is the probability that the pattern I'm seeing might have arisen by chance?

With what parameters does the data best fit a given function? What is the goodness of fit?

How well do one (or more) data variables predict another?

...and many others.

# Example: Heights by Gender











# Formulating a Hypothesis

	<u> </u>				
Null Hypothesis Alternate Hypotl	(H <sub>0</sub> ): nesis (H	l <sub>a</sub> ):	μ <sub>m</sub> = μ <sub>f</sub> μ <sub>m</sub> ≠ μ <sub>f</sub>	(population) (population)	
A statistical hypothesis test assesses the likelihood of the null hypothesis.					
What is the probability of sampling the observed data assuming population means are equal?					
This is called the	e p valu	e			



Compute a test statistic. This is a number that in essence summarizes the difference.



# **Testing Procedure**

Compute a test statistic. This is a number that in essence summarizes the difference.

The possible values of this statistic come from a known probability distribution.

According to this distribution, look up the probability of seeing a value meeting or exceeding the test statistic. This is the *p* value.



# **Statistical Significance**

- The threshold at which we consider it safe (or reasonable?) to reject the null hypothesis.
- If p < 0.05, we typically say that the observed effect or difference is statistically significant.
- This means that there is a less than 5% chance that the observed data is due to chance.
- Note that the choice of 0.05 is a somewhat arbitrary threshold (chosen by R. A. Fisher)

### **Graphical Inference**

Buja, Cook, Hoffman, Wickham et al.







One plot is the real data. The others are generated using an assumption of normally distributed residuals.









# Tableau demo

#### The dataset:

- Federal Elections Commission Receipts
- Every Congressional Candidate from 1996 to 2002
- 4 Election Cycles
- 9216 Candidacies

## **Data Set Schema**

- Year (Qi)
- Candidate Code (N)
- Candidate Name (N)
- Incumbent / Challenger / Open-Seat (N)
- Party Code (N) [1=Dem,2=Rep,3=Other]
- Party Name (N)
- Total Receipts (Qr)
- State (N)
- **District (N)**
- This is a subset of the larger data set available from the FEC, but should be sufficient for the demo

# Hypotheses?

#### What might we learn from this data?

# Hypotheses?

# What might we learn from this data?

- Has spending increased over time?
- Do democrats or republicans spend more money?
- Candidates from which state spend the most money?

## Tableau Demo