



A Data Story: Cholera in 1854

School of Information Studies
Syracuse University

Jeff Hemsley, PhD
Assistant Professor
School of Information Studies

A Story of Patterns

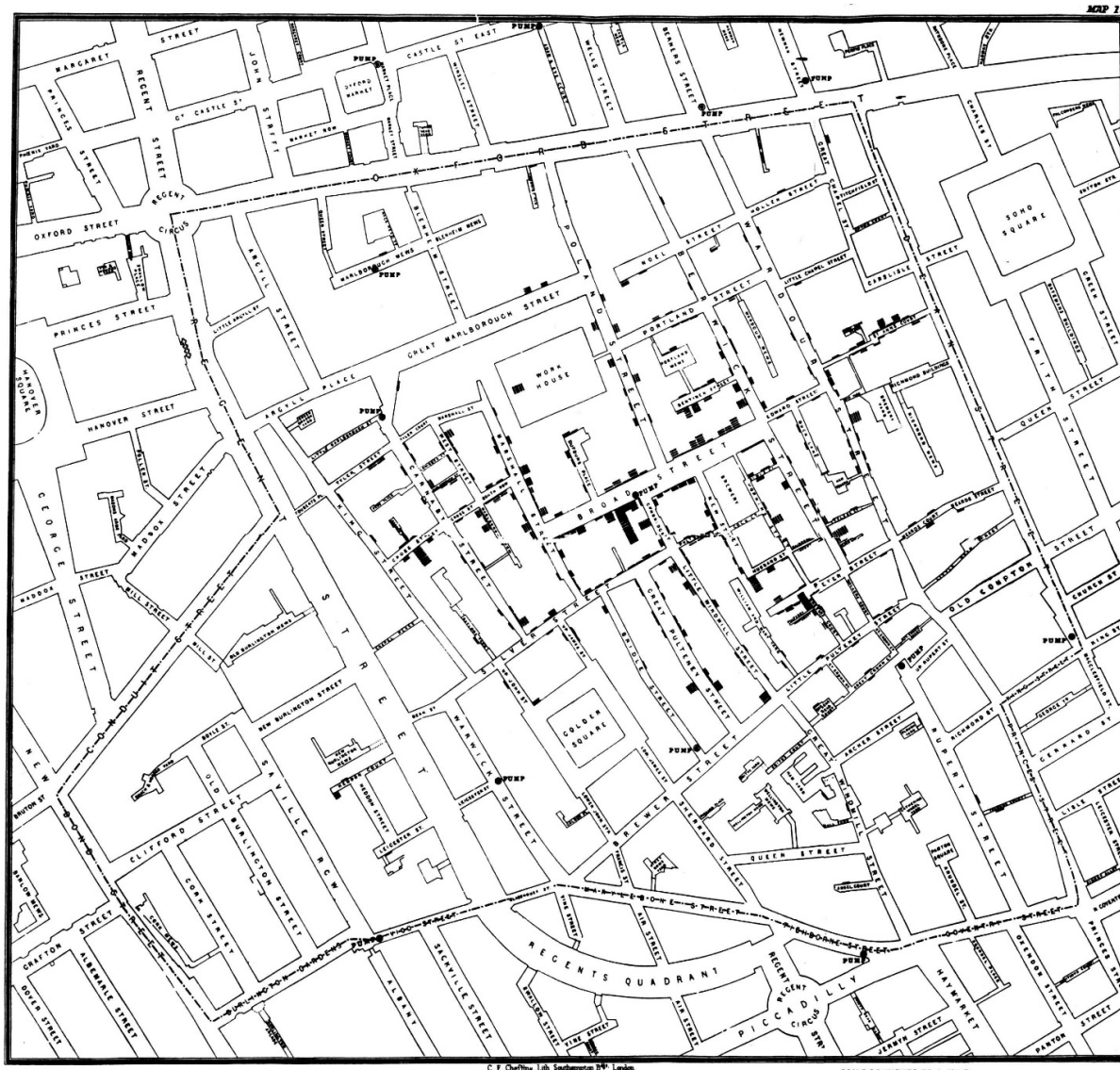
John Snow

Cholera outbreak on August 31, 1854

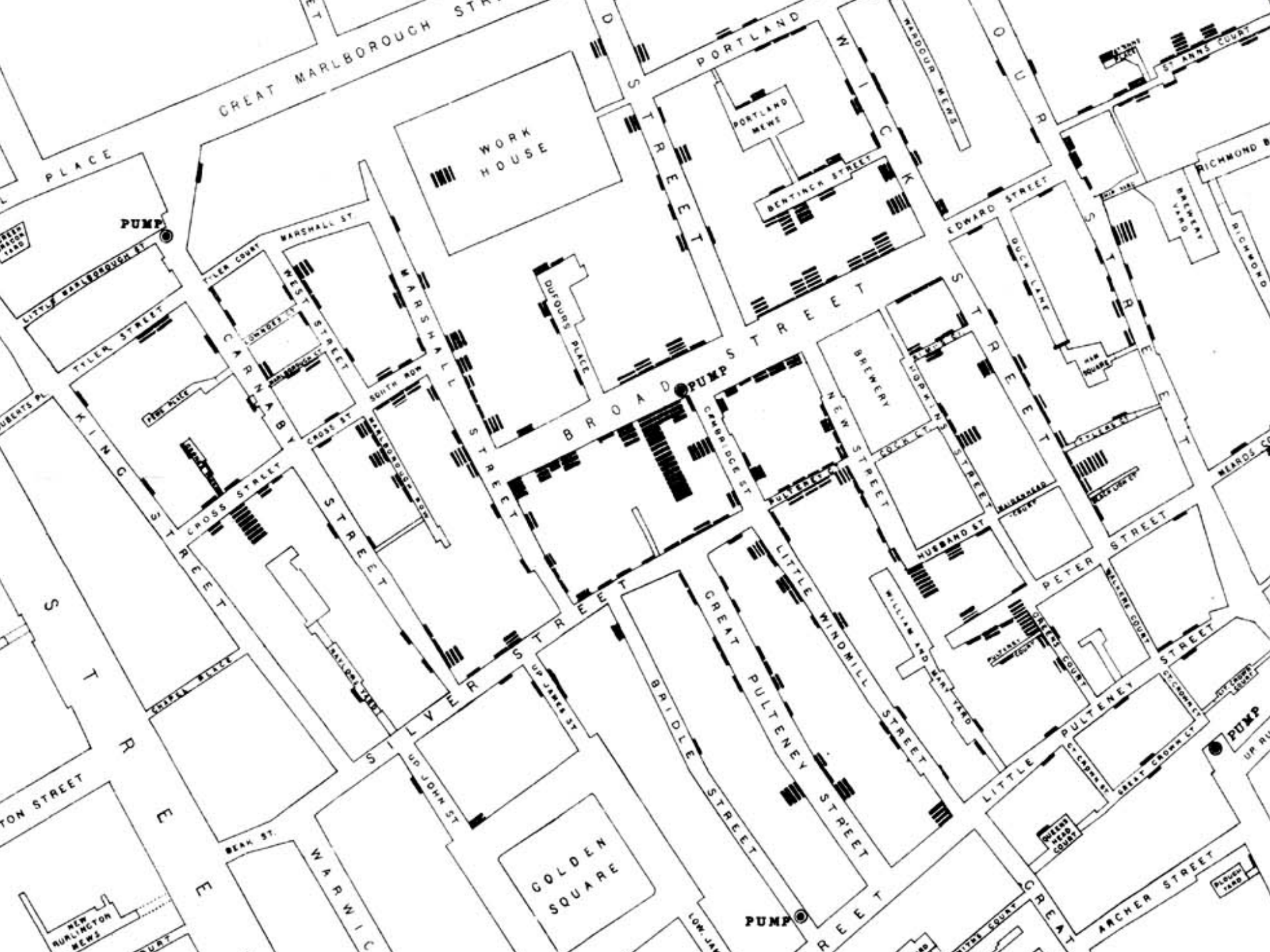
Data: 83 deaths from cholera

Motivating questions:

- Was there a pattern to the outbreak?
Reduce uncertainty
- Could the pattern tell us anything about the cause of the outbreak?
Look for exceptions
- Could impurities in water be to blame?
Hypothesis testing



Commons











Finding Data

Jeff Hemsley, PhD
Assistant Professor
School of Information Studies

School of Information Studies
Syracuse University

Finding Data

Universities

- [Syracuse University data resources](#)

Government

- [U.S. Government's open data](#)
- [U.S. climate data](#)
 - Education/data (tab)
- [NASA data portal](#)

Others

- [Kaggle](#)
- [Reddit datasets](#)

Data clean-up

- [U.S. Census Bureau \(Construction Spending\)](#)



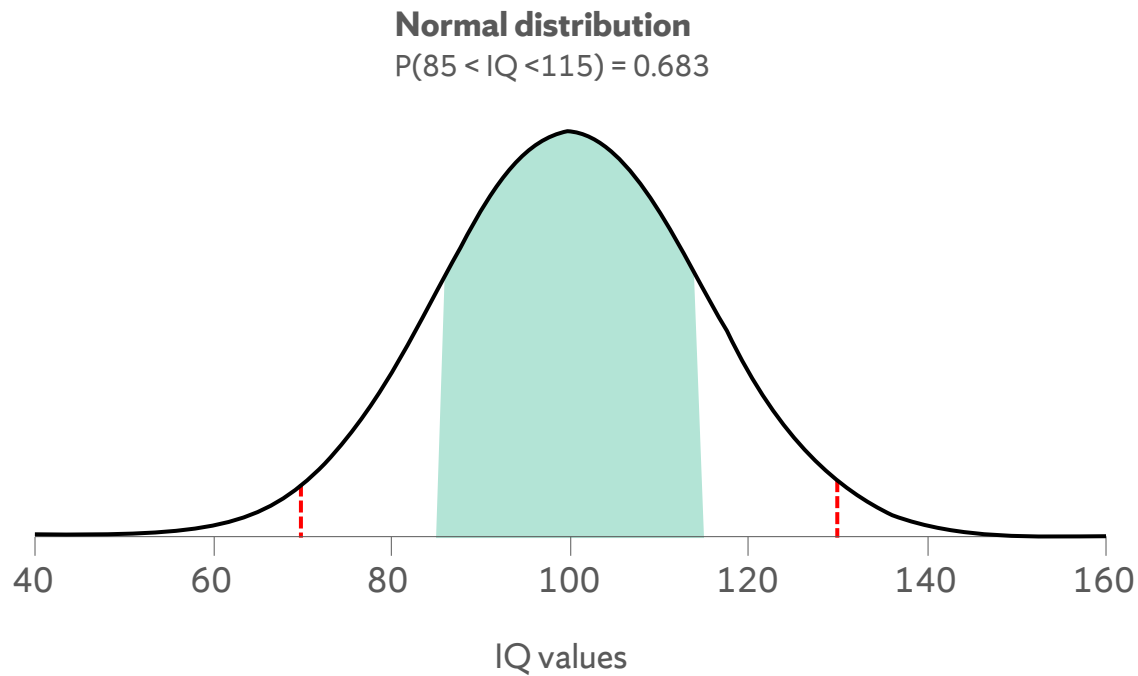


Distributions and Plot Dimensions

School of Information Studies
Syracuse University

Jeff Hemsley, PhD
Assistant Professor
School of Information Studies

Normal Distribution



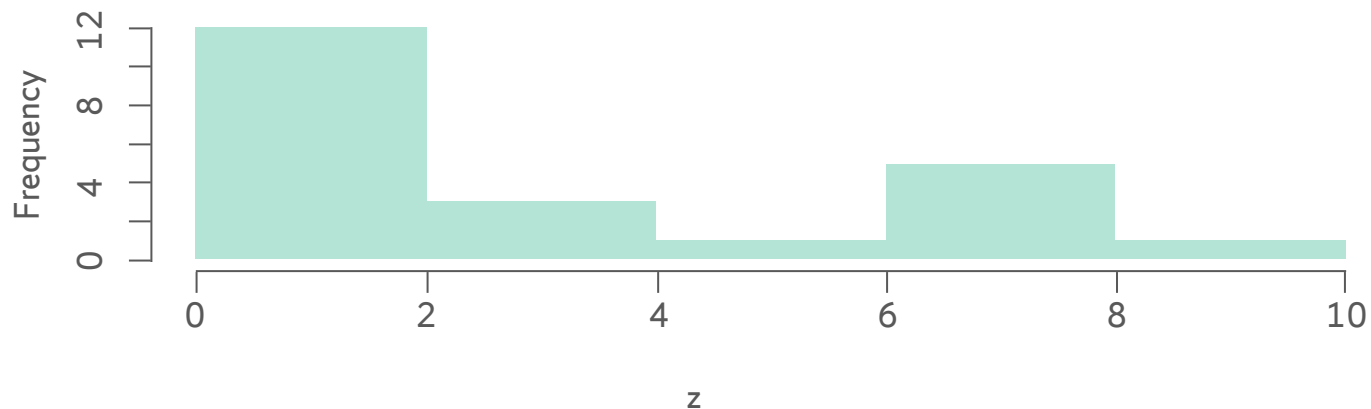
Distribution

Given a set of numbers, what is the central tendency? What is the shape?

1,5,0,3,8,8,1,2,1,8,7,2,1,2,10,3,2,8,1,0,3,0

0,0,0,1,1,1,1,1,2,2,2,2,3,3,3,5,7,8,8,8,8,10

Histogram of z



Dimensions of Your Plots

An Informal Way to Talk About Plots

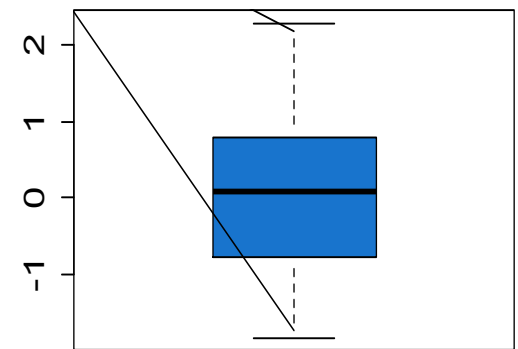
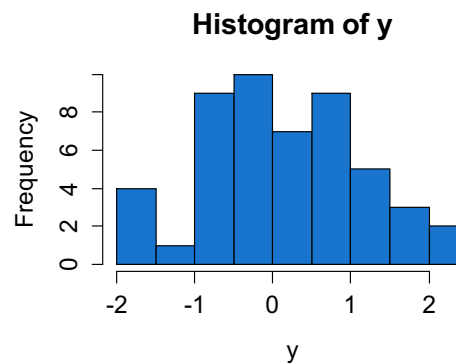
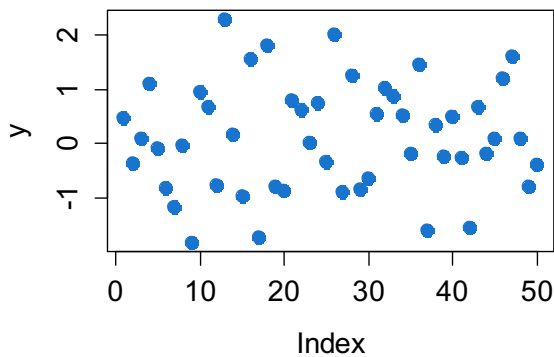
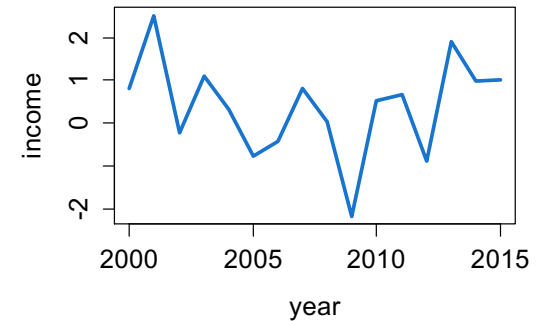
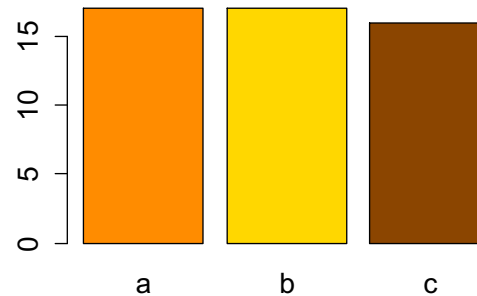
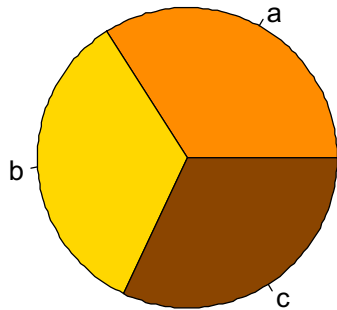
Single

- Numeric or categorical
- Distribution within a single variable
- Frequency within a single variable

Multiple

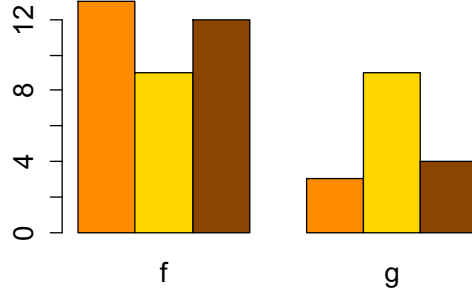
- Relating one variable to another
- Distribution of variable across categories
- Frequency of a categories across other categories

Single Dimension Plots

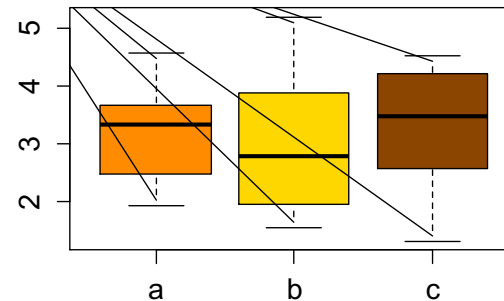


Multi-Dimension Plots

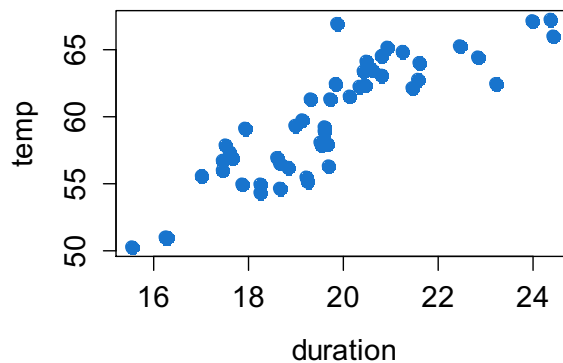
cat by cat



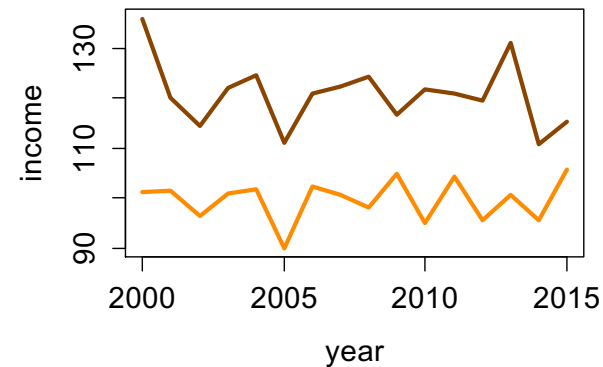
number by cat



num by num

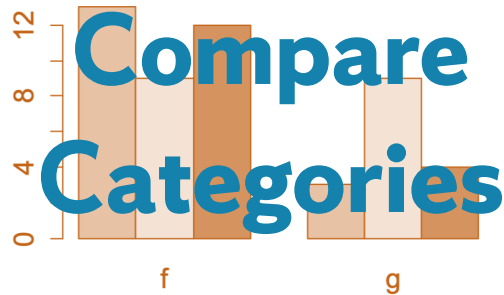


num by time (num/cat)

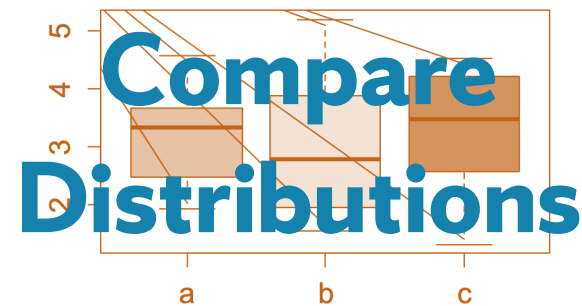


Multi-Dimension Plots

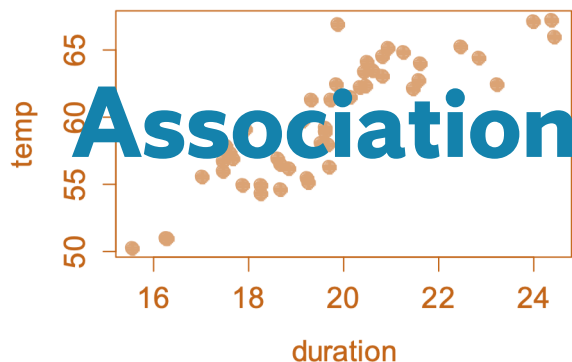
cat by cat



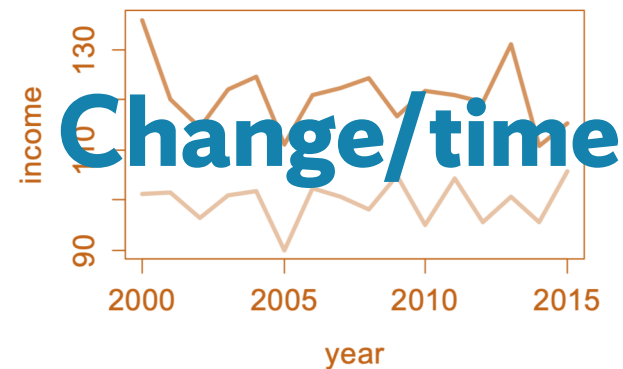
number by cat



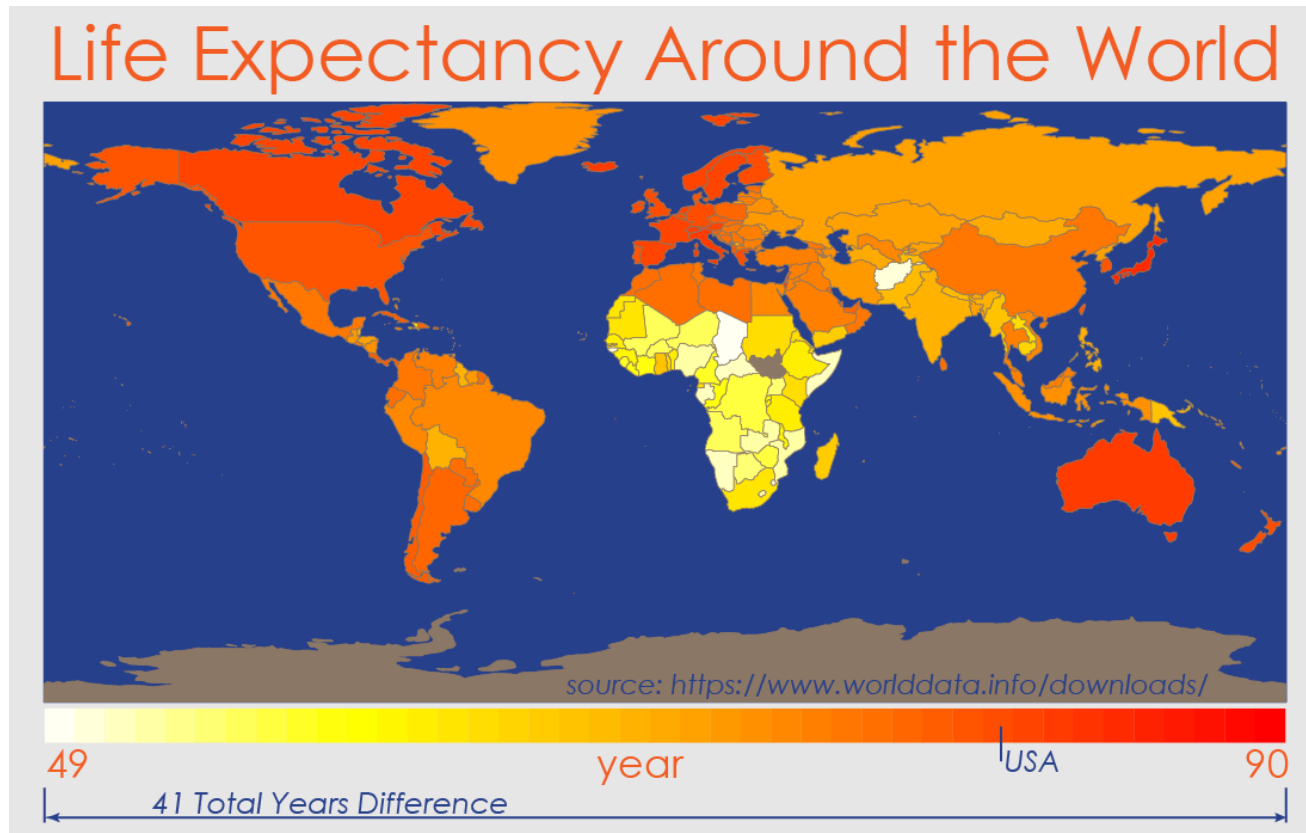
num by num



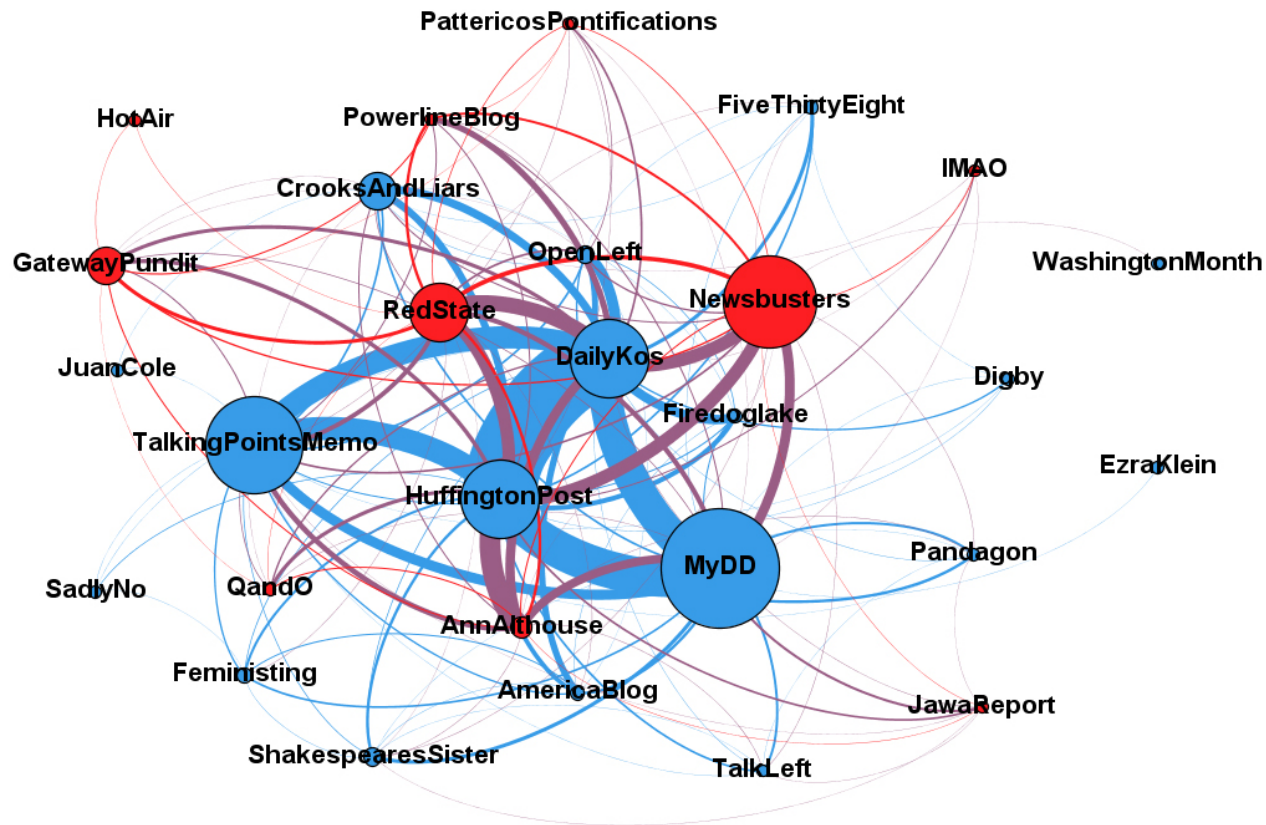
num by time (num/cat)



Multi-Dimension Plots



Multi-Dimension Plots



Comparisons

What kinds of comparisons?

Magnitude

- Nominal: no particular order, frequency
- Ranking: magnitude
- Parts-to-whole: proportions
- Deviation: differences between sets
- Time-series: change over time

Sorting

- Low to high

Position

- Maps

Strength of relationships



