

CPSC 340: Machine Learning and Data Mining

Data Exploration

Fall 2017

This lecture roughly follow:

http://www-users.cs.umn.edu/~kumar/dmbook/dmslides/chap2_data.pdf

Admin

- **Assignment 0** is due next Friday: start early.
- Waiting list people: you should be registered next week.
- Bookmark the course webpage:
 - www.cs.ubc.ca/~schmidtm/Courses/340-F17
- Sign up for the course **Piazza** group:
 - www.piazza.com/ubc.ca/winterterm12017/cpsc340/home
- Sign up for a CS undergrad account:
 - <https://www.cs.ubc.ca/getacct>
- Tutorials start next week.
 - Optional, but you need to be registered in a tutorial.
- Office hours start next week:
 - Watch the website for details.
- Auditing: message me on Piazza if you want to audit.

Data Mining: Bird's Eye View

- 1) Collect data.
- 2) Data mining!
- 3) Profit?

Unfortunately, it's often more complicated...

Data Mining: Some Typical Steps

- 1) Learn about the application.
 - 2) Identify data mining task.
 - 3) Collect data.
 - 4) Clean and preprocess the data.
 - 5) Transform data or select useful subsets.
 - 6) Choose data mining algorithm.
 - 7) Data mining!
 - 8) Evaluate, visualize, and interpret results.
 - 9) Use results for profit or other goals.
- (often, you'll go through cycles of the above)

Data Mining: Some Typical Steps

- 1) Learn about the application.
- 2) Identify data mining task.
- 3) Collect data.
- 4) Clean and preprocess the data.
- 5) Transform data or select useful subsets.
- 6) Choose data mining algorithm.
- 7) Data mining!
- 8) Evaluate, visualize, and interpret results.
- 9) Use results for profit or other goals.

(often, you'll go through cycles of the above)

What is Data?

- We'll define data as a collection of **objects**, and their **features**.

Age	Job?	City	Rating	Income
23	Yes	Van	A	22,000.00
23	Yes	Bur	BBB	21,000.00
22	No	Van	CC	0.00
25	Yes	Sur	AAA	57,000.00
19	No	Bur	BB	13,500.00
22	Yes	Van	A	20,000.00
21	Yes	Ric	A	18,000.00

- Each row is an object, each column is a feature.

Object
(Training example, sample)

Feature
(Variables, covariates)

Types of Data

- **Categorical features** come from an unordered set:
 - Binary: job?
 - Nominal: city.
- **Numerical features** come from ordered sets:
 - Discrete counts: age.
 - Ordinal: rating.
 - **Continuous**/real-valued: height.

Converting to Continuous Features

- Often want a real-valued object representation:

Age	City	Income
23	Van	22,000.00
23	Bur	21,000.00
22	Van	0.00
25	Sur	57,000.00
19	Bur	13,500.00
22	Van	20,000.00

→

Age	Van	Bur	Sur	Income
23	1	0	0	22,000.00
23	0	1	0	21,000.00
22	1	0	0	0.00
25	0	0	1	57,000.00
19	0	1	0	13,500.00
22	1	0	0	20,000.00

- We can now **interpret objects as points** in space:
 - E.g., first object is at (23,1,0,0,22000).

↓
"l of K"

Approximating Text with Continuous Features

- **Bag of words** replaces document by word counts:

The **International Conference on Machine Learning** (ICML) is the leading international academic conference in machine learning

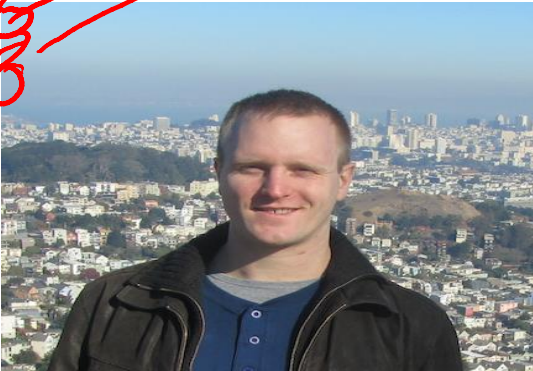
ICML	International	Conference	Machine	Learning	Leading	Academic
1	2	2	2	2	1	1

- Ignores order, but often captures general theme.
- You can compute 'distance' between documents.

Approximating Images and Graphs

- We can think of other data types in this way:

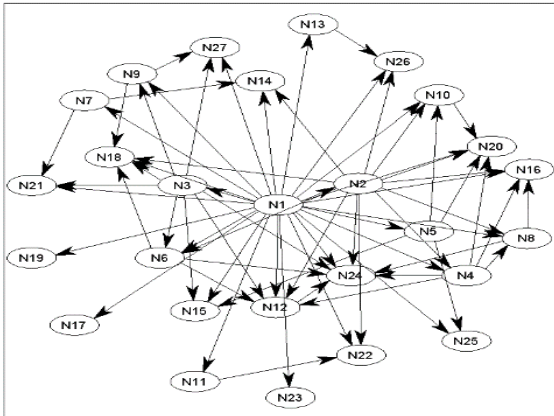
– Images:



graycale
intensity

(1,1)	(2,1)	(3,1)	...	(m,1)	...	(m,n)
45	44	43	...	12	...	35

– Graphs:



adjacency
matrix

	N1	N2	N3	N4	N5	N6	N7
N1	0	1	1	1	1	1	1
N2	0	0	0	1	0	1	0
N3	0	0	0	0	0	1	0
N4	0	0	0	0	0	0	0

N1 is a parent of N5

Data Cleaning

- ML+DM typically assume 'clean' data.
- Ways that data might not be 'clean':
 - Noise (e.g., distortion on phone).
 - Outliers (e.g., data entry or instrument error).
 - Missing values (no value available or not applicable)
 - Duplicated data (repetitions, or different storage formats).
- Any of these can lead to problems in analyses.
 - Want to fix these issues, if possible.
 - Some ML methods are robust to these.
 - Often, [ML is the best way to detect/fix](#) these.

The Question I Hate the Most...

- How much data do we need?
- A difficult if not impossible question to answer.
- My usual answer: “more is better”.
 - With the warning: “as long as the quality doesn’t suffer”.
- Another popular answer: “ten times the number of features”.

A Simple Setting: Coupon Collecting

- Assume we have a categorical variable with 50 possible values:
 - {Alabama, Alaska, Arizona, Arkansas,...}.
- Assume each category has probability of 1/50 of being chosen:
 - How many objects do we need to see before we expect to see them all?
- Expected value is ~225.
- Coupon collector problem: $O(n \log n)$ in general.
 - [Gotta Catch'em all!](#)
- Obvious sanity check, is need **more samples than categories**:
 - Situation is worse if they don't have equal probabilities.
 - Typically want to see categories more than once to learn anything.

Feature Aggregation

- Feature aggregation:
 - Combine features to form new features:

Van	Bur	Sur	Edm	Cal		BC	AB
1	0	0	0	0		1	0
0	1	0	0	0		1	0
1	0	0	0	0	→	1	0
0	0	0	1	0		0	1
0	0	0	0	1		0	1
0	0	1	0	0		1	0

- More province information than city information.

Feature Selection

- Feature Selection:
 - Remove features that are not relevant to the task.

SID:	Age	Job?	City	Rating	Income
3457	23	Yes	Van	A	22,000.00
1247	23	Yes	Bur	BBB	21,000.00
6421	22	No	Van	CC	0.00
1235	25	Yes	Sur	AAA	57,000.00
8976	19	No	Bur	BB	13,500.00
2345	22	Yes	Van	A	20,000.00

- Student ID is probably not relevant.

Feature Transformation

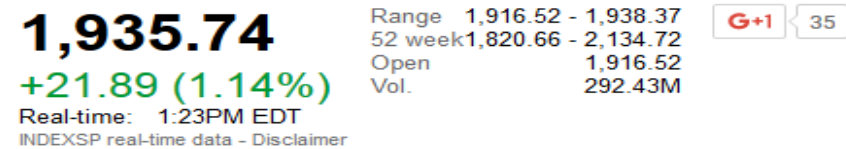
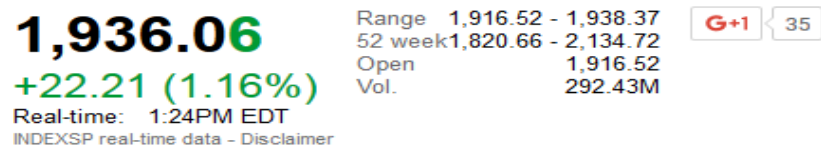
- Mathematical transformations:
 - **Discretization** (binning): turn numerical data into categorical.

Age		< 20	>= 20, < 25	>= 25
23		0	1	0
23	→	0	1	0
22		0	1	0
25		0	0	1
19		1	0	0
22		0	1	0

- Only need consider 3 values.

Feature Transformation

- Mathematical transformations:
 - **Discretization** (binning): turn numerical data into categorical.
 - Square, exponentiation, or take logarithm.

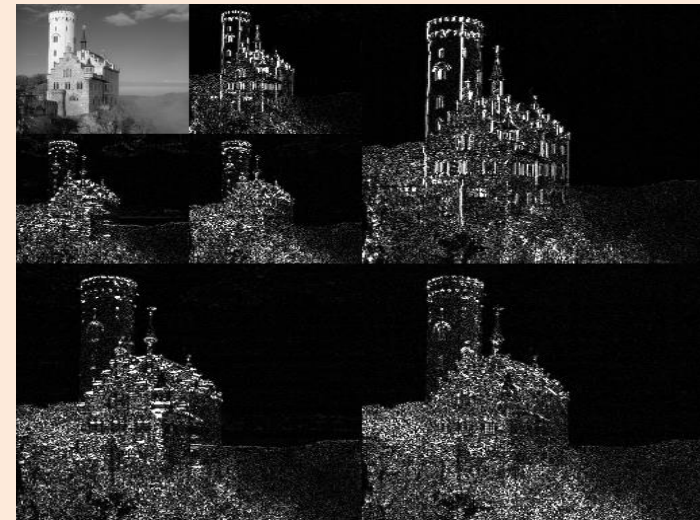
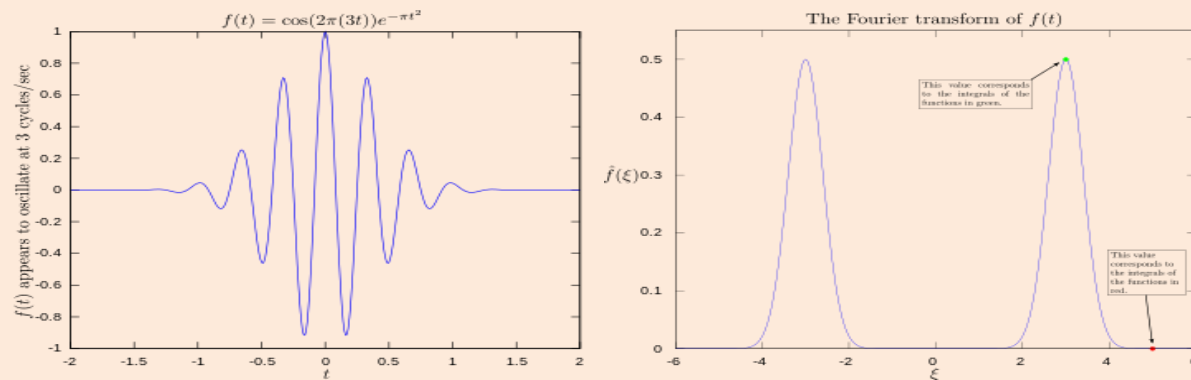


Feature Transformation

- Mathematical transformations:
 - **Discretization** (binning): turn numerical data into categorical.
 - Square, exponentiation, or take logarithm.
 - Scaling: convert variables to comparable scales (E.g., convert kilograms to grams.)

Feature Transformation

- Mathematical transformations:
 - **Discretization** (binning): turn numerical data into categorical.
 - Square, exponentiation, or take logarithm.
 - Scaling: convert variables to comparable scales.
 - Fourier coefficients, spectrograms, and wavelets (signal data).



https://en.wikipedia.org/wiki/Fourier_transform
<https://en.wikipedia.org/wiki/Spectrogram>
https://en.wikipedia.org/wiki/Discrete_wavelet_transform

Links to figure sources will be here.

(pause)

Exploratory Data Analysis

- You should always ‘look’ at the data first.
- But how do you ‘look’ at features and high-dimensional objects?
 - Summary statistics.
 - Visualization.
 - ML + DM (later in course).

Categorical Summary Statistics

- Summary statistics for a categorical variable:
 - **Frequencies** of different classes.
 - **Mode**: category that occurs most often.
 - **Quantiles**: categories that occur more than t times:

**Population by year, by province and territory
(Number)**

	2014
Canada	35,540.4
Newfoundland and Labrador	527.0
Prince Edward Island	146.3
Nova Scotia	942.7
New Brunswick	753.9
Quebec	8,214.7
Ontario	13,678.7
Manitoba	1,282.0
Saskatchewan	1,125.4
Alberta	4,121.7
British Columbia	4,631.3
Yukon	36.5
Northwest Territories	43.6
Nunavut	36.6

Frequency: **13.3%** of Canadian residents live in BC.
Mode: **Ontario** has largest number of residents (38.5%)
Quantile: **6** provinces have **more than 1 million** people.

Continuous Summary Statistics

- Measures of location:
 - **Mean**: average value.
 - **Median**: value such that half points are larger/smaller.
 - **Quantiles**: value such that 't' points are larger.
- Measures of spread:
 - **Range**: minimum and maximum values.
 - **Variance**: measures how far values are from mean.
 - Square root of variance is “standard deviation”.
 - **Intequantile ranges**: difference between quantiles.

Continuous Summary Statistics

- Data: [0 1 2 3 3 5 7 8 9 10 14 15 17 200]
 - Measures of location:
 - Mean(Data) = 21
 - Mode(Data) = 3
 - Median(Data) = 7.5
 - Quantile(Data,0.5) = 7.5
 - Quantile(Data,0.25) = 3
 - Quantile(Data,0.75) = 14
 - Measures of spread:
 - Range(Data) = [0 200].
 - Std(Data) = 51.79
 - IQR(Data,.25,.75) = 11
 - Notice that mean and std are more sensitive to extreme values (“outliers”).
- Handwritten red text: "outlier" with an arrow pointing to the value 200 in the data list.*

Distances and Similarities

- There are also summary **statistics between features 'x' and 'y'**.
 - **Hamming distance:**
 - Number of elements in the vectors that aren't equal.
 - **Euclidean distance:**
 - How far apart are the vectors?
 - **Correlation:**
 - Does one increase/decrease linearly as the other increases?

x	y
0	0
0	0
1	0
0	1
0	1
1	1
0	0
0	1
0	1

Distances and Similarities

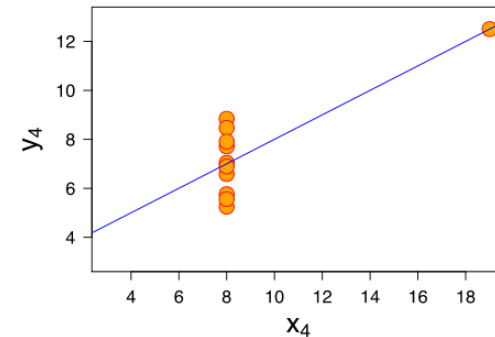
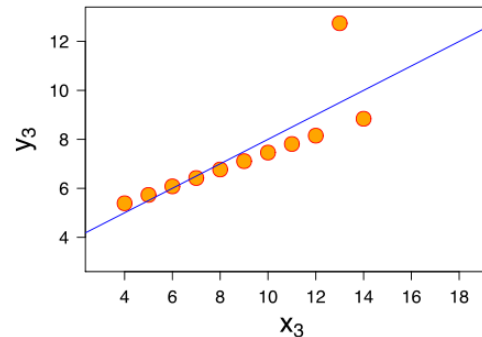
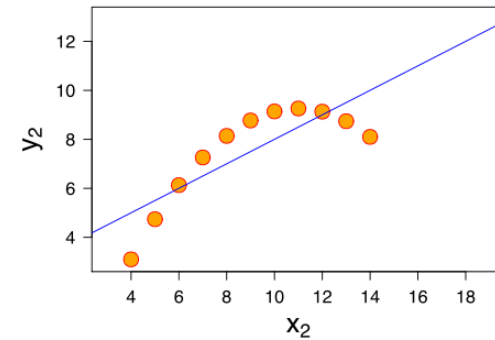
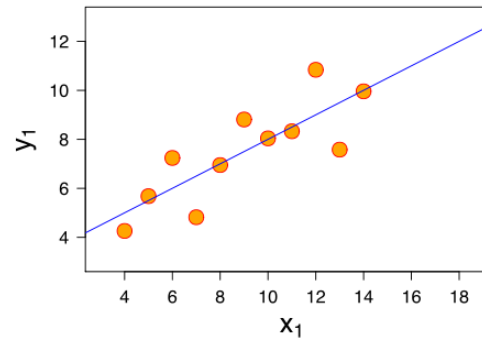
- There are also summary **statistics between features** 'x' and 'y'.
 - **Rank correlation**:
 - Does one increase/decrease non-linearly as the other increases?
- Distances/similarities between other objects:
 - **Jaccard coefficient** (distance between sets):
 - $(\text{size of intersection of sets}) / (\text{size of union of sets})$
 - **Edit distance** (distance between strings):
 - How many characters do we need to change to go from x to y?
 - Computed using dynamic programming (CPSC 320).

x	y
0	0
0	0
1	0
0	1
0	1
1	1
0	0
0	1
0	1

Limitations of Summary Statistics

- On their own **summary statistic can be misleading.**
- [Why not to trust statistics](#)

- Amcomb's quartet:
 - Almost same means.
 - Almost same variances.
 - Almost same correlations.
 - Look completely different.
- [Datasaurus dozen.](#)



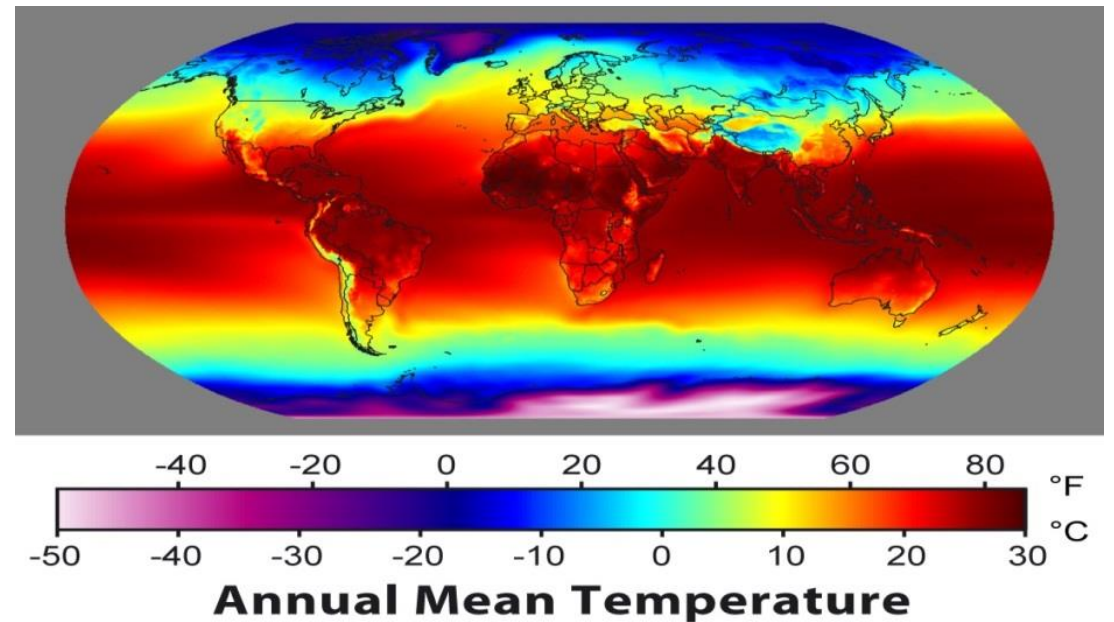
(pause)

Visualization

- You can learn a lot from **2D plots** of the data:
 - Patterns, trends, outliers, unusual patterns.

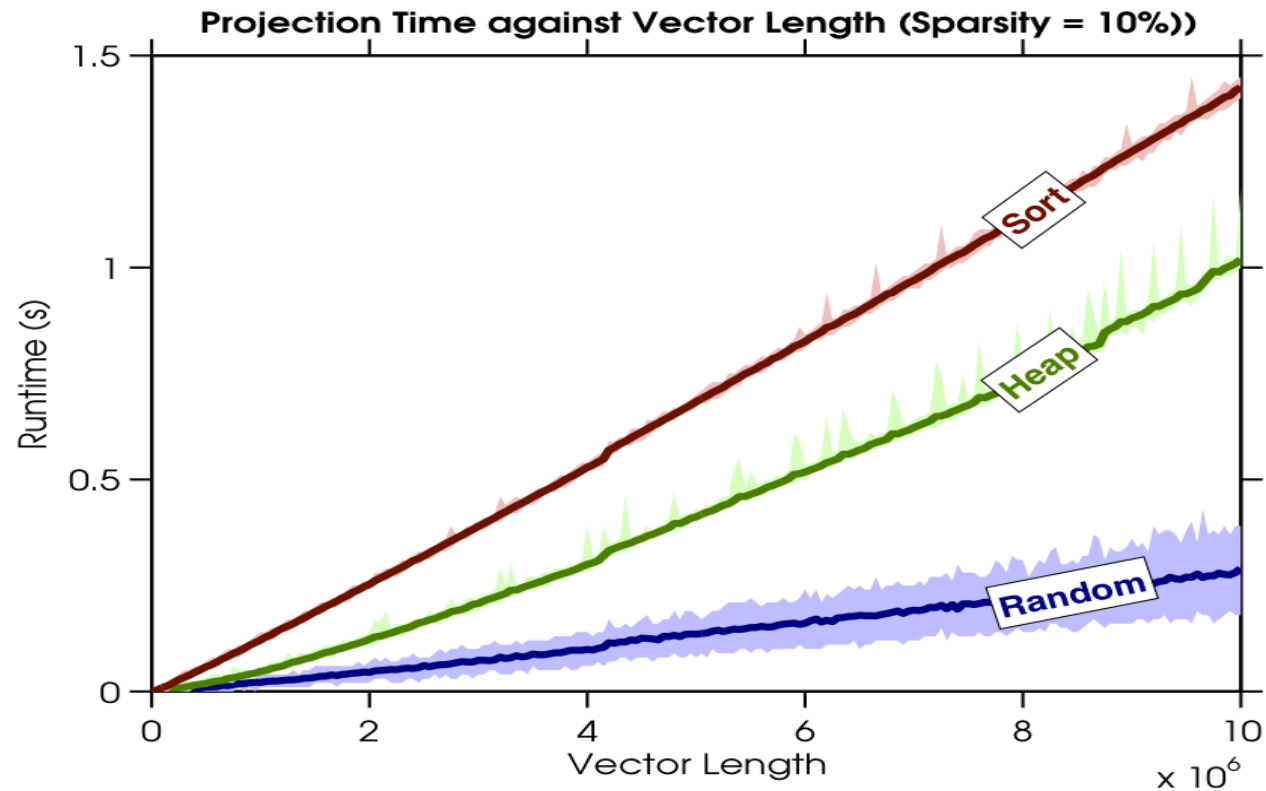
Lat	Long	Temp
0	0	30.1
0	1	29.8
0	2	29.9
0	3	30.1
0	4	29.9
...

vs.



Basic Plot

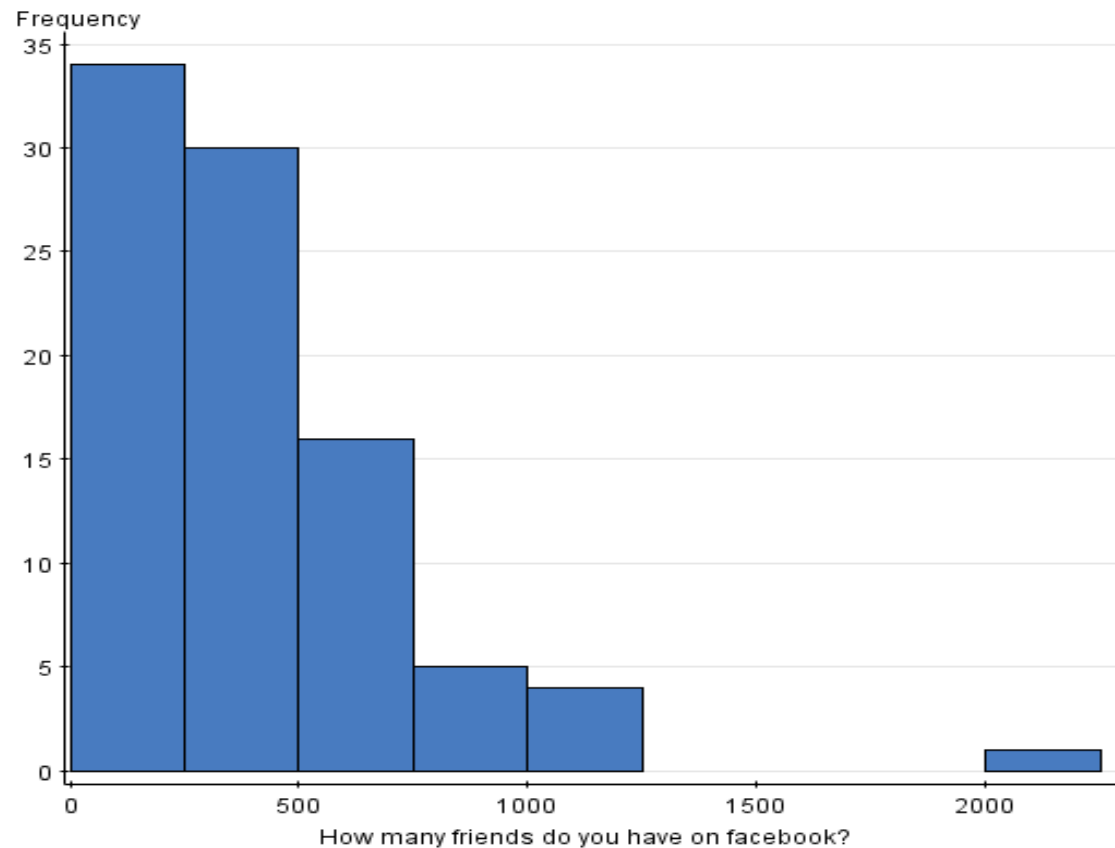
- Visualize one variable as a function of another.



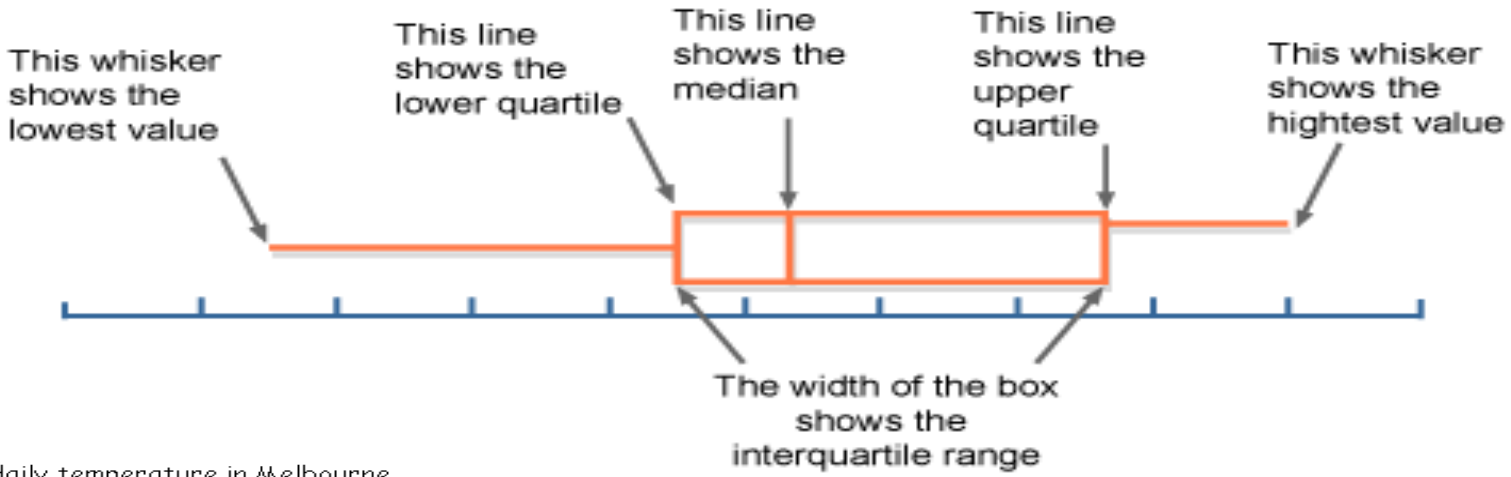
- [Fun with plots.](#)

Histogram

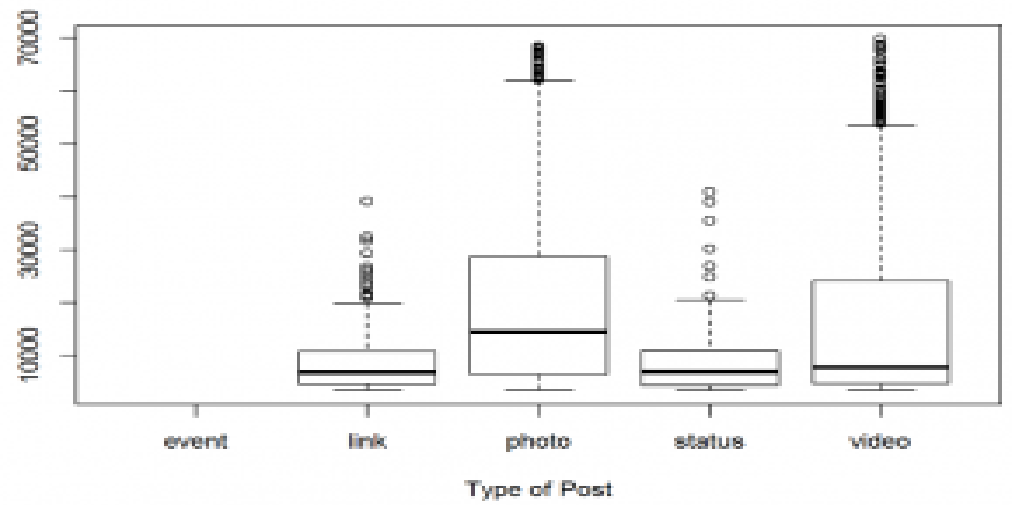
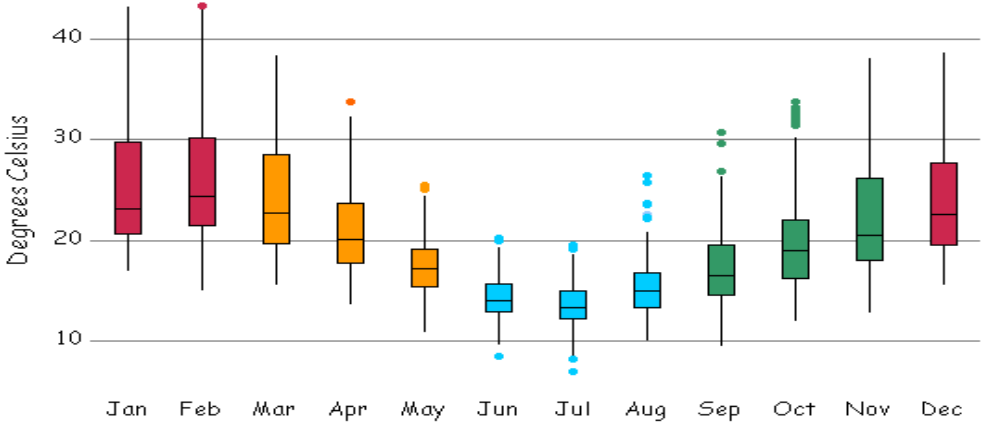
- Histograms display distribution of a variable.



Box Plot



Maximum daily temperature in Melbourne



<http://www.bbc.co.uk/schools/gcsebitesize/maths/statistics/representingdata3hi>
<http://www.scc.ms.unimelb.edu.au/whatisstatistics/weather.html>
<http://r.ramganalytics.com/r/facebook-likes-and-analytics/>

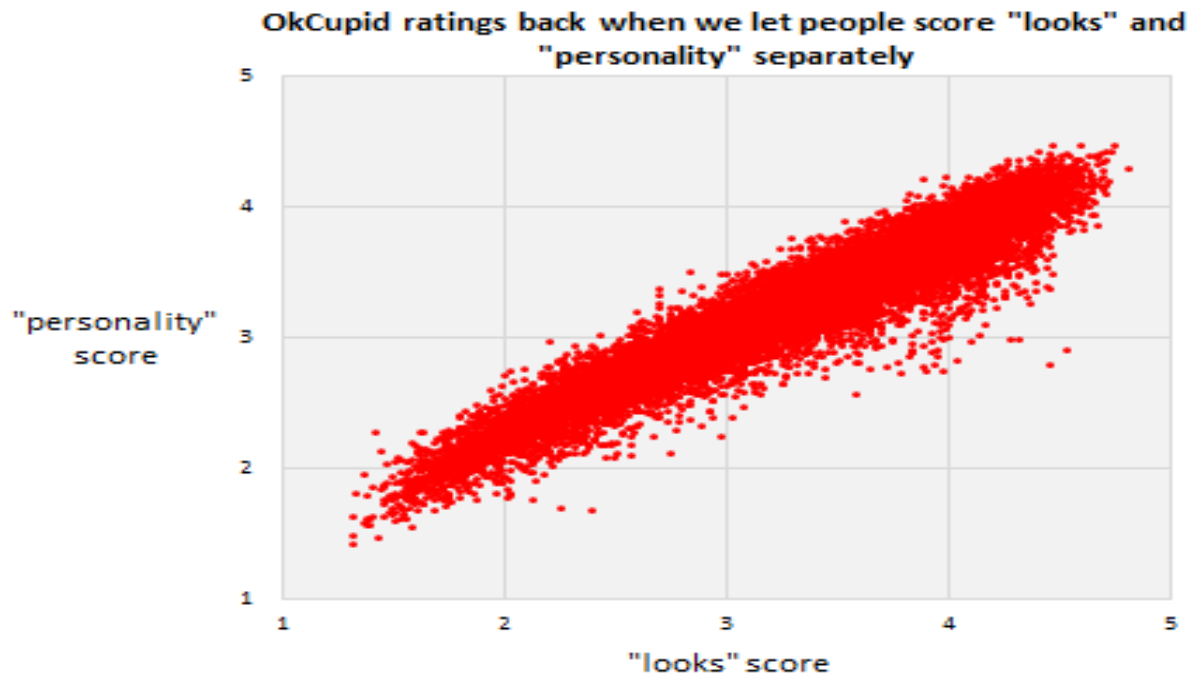
Box Plot

- Photo from CTV Olympic coverage in 2010:



Scatterplot

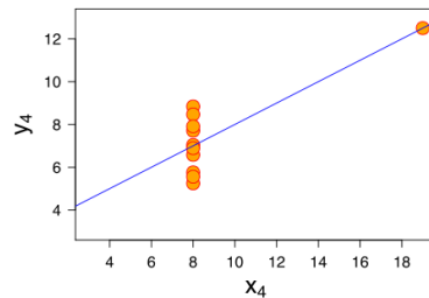
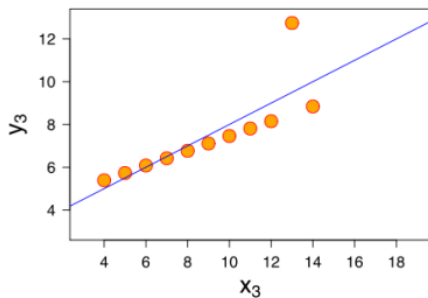
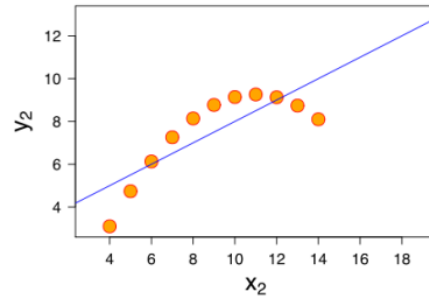
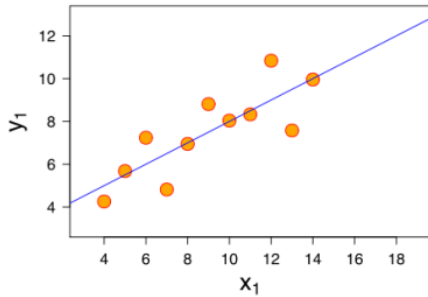
- Look at distribution of two features:
 - Feature 1 on x-axis.
 - Feature 2 on y-axis.
 - Basically a “plot without lines” between the points.



- Shows correlation between “personality” score and “looks” score.

Scatterplot

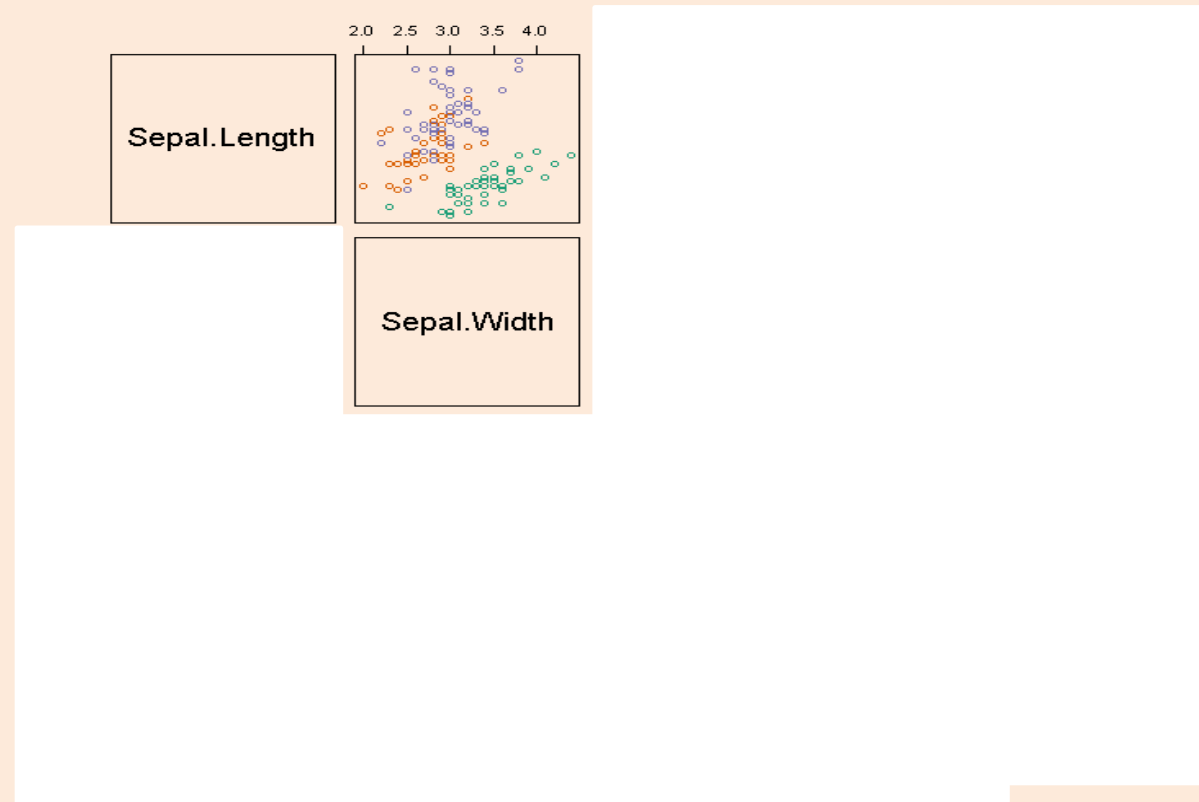
- Look at distribution of two features:
 - Feature 1 on x-axis.
 - Feature 2 on y-axis.
 - Basically a “plot without lines” between the points.



- Shows correlation between “personality” score and “looks” score.
- But scatterplots let you **see more complicated patterns.**

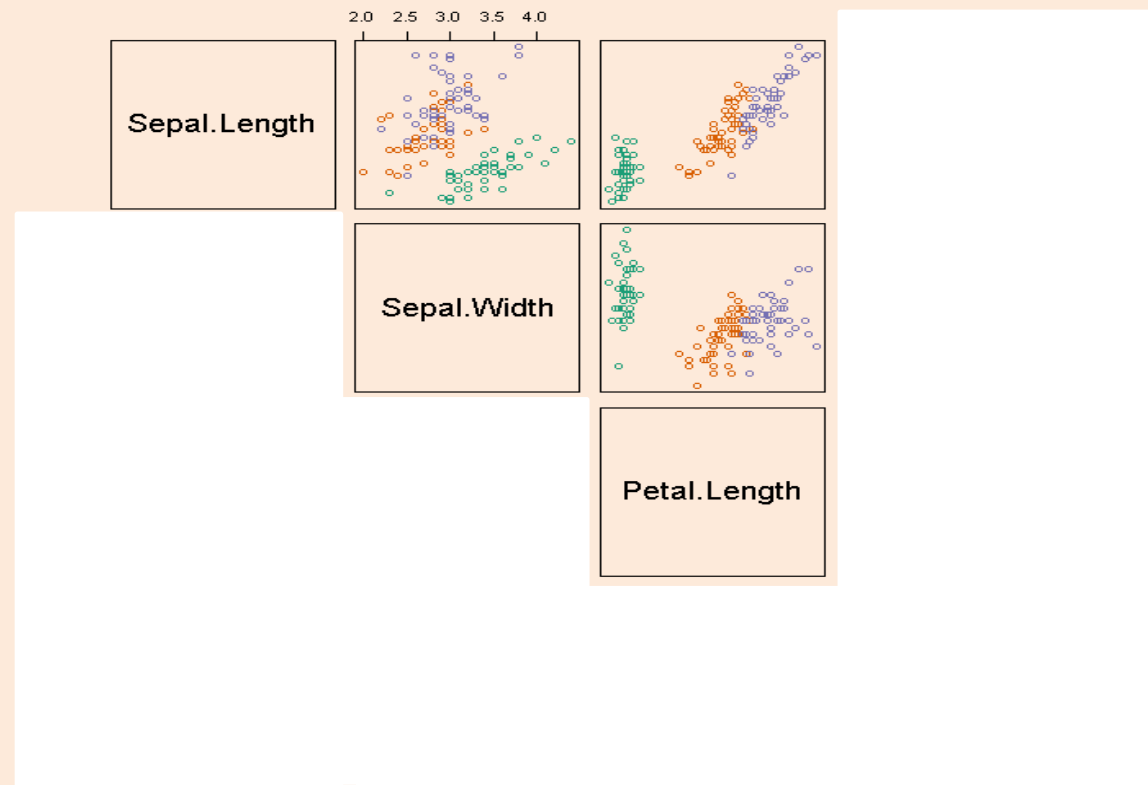
Scatterplot Arrays

- For multiple variables, can use **scatterplot array**.



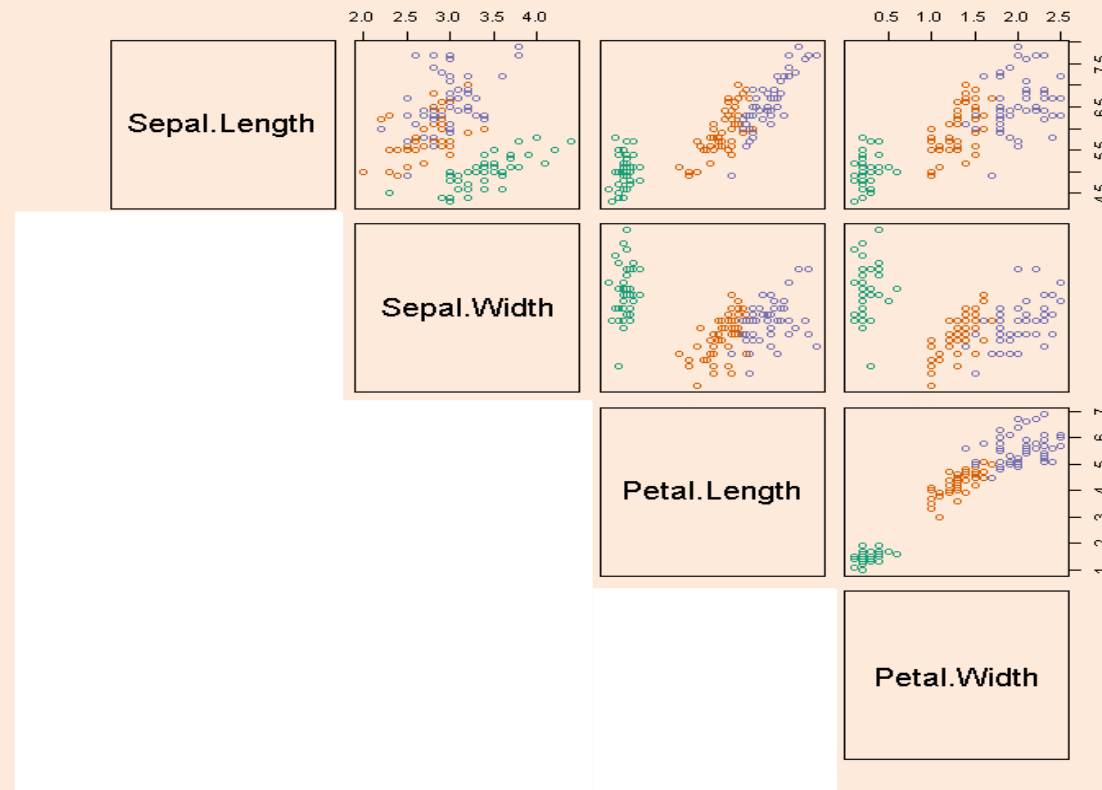
Scatterplot Arrays

- For multiple variables, can use **scatterplot array**.



Scatterplot Arrays

- For multiple variables, can use **scatterplot array**.



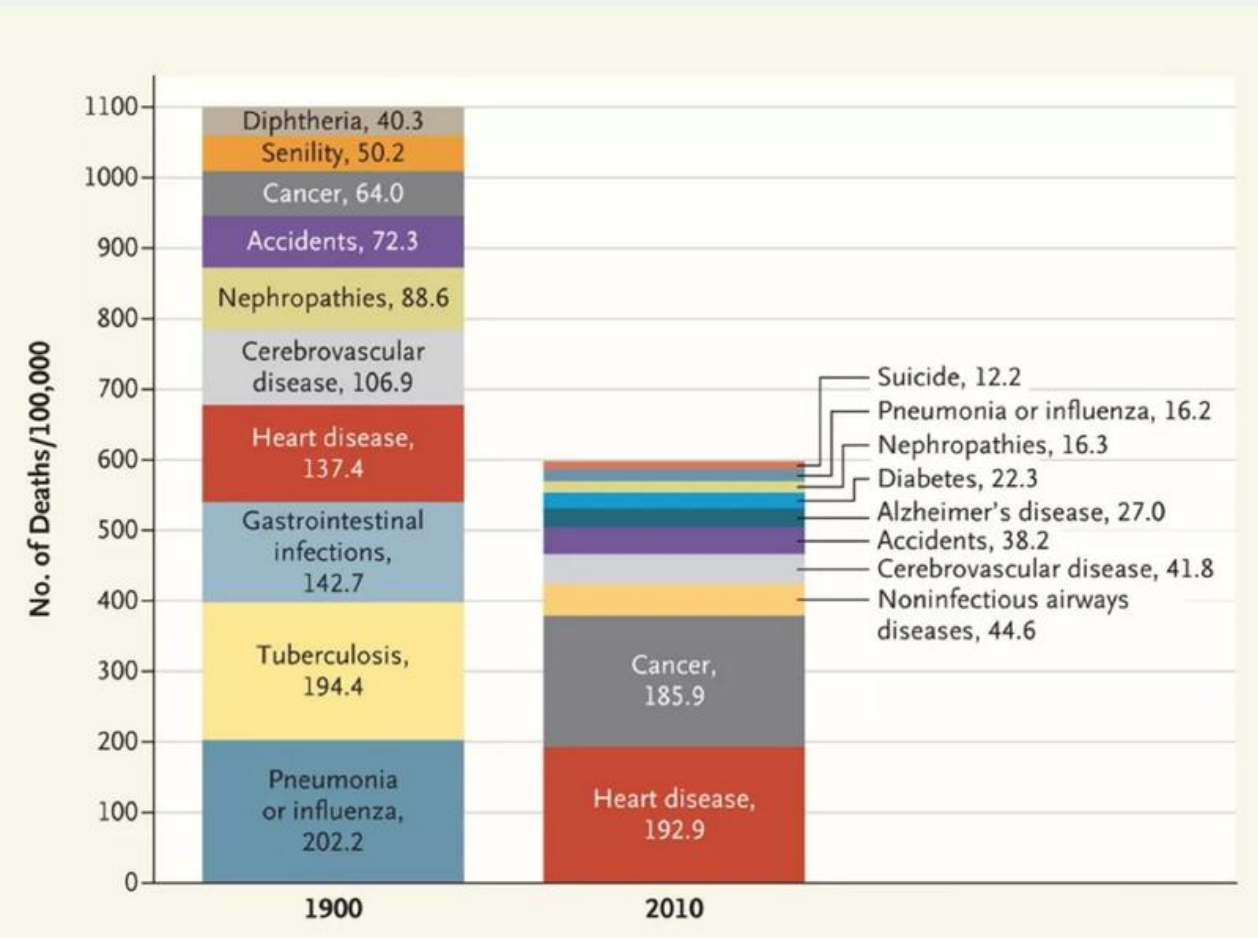
- Colors can indicate a third categorical variable.

Summary

- **Typical data mining steps:**
 - Involves data collection, preprocessing, analysis, and evaluation.
- **Object-feature representation** and categorical/numerical features.
 - Transforming non-vector objects to vector representations.
- **Feature transformations:**
 - To address coupon collecting or simplify relationships between variables.
- **Exploring data:**
 - Summary statistics and data visualization.
- **Post-lecture bonus slides:** other visualization methods.
- Next week: let's start some machine learning...

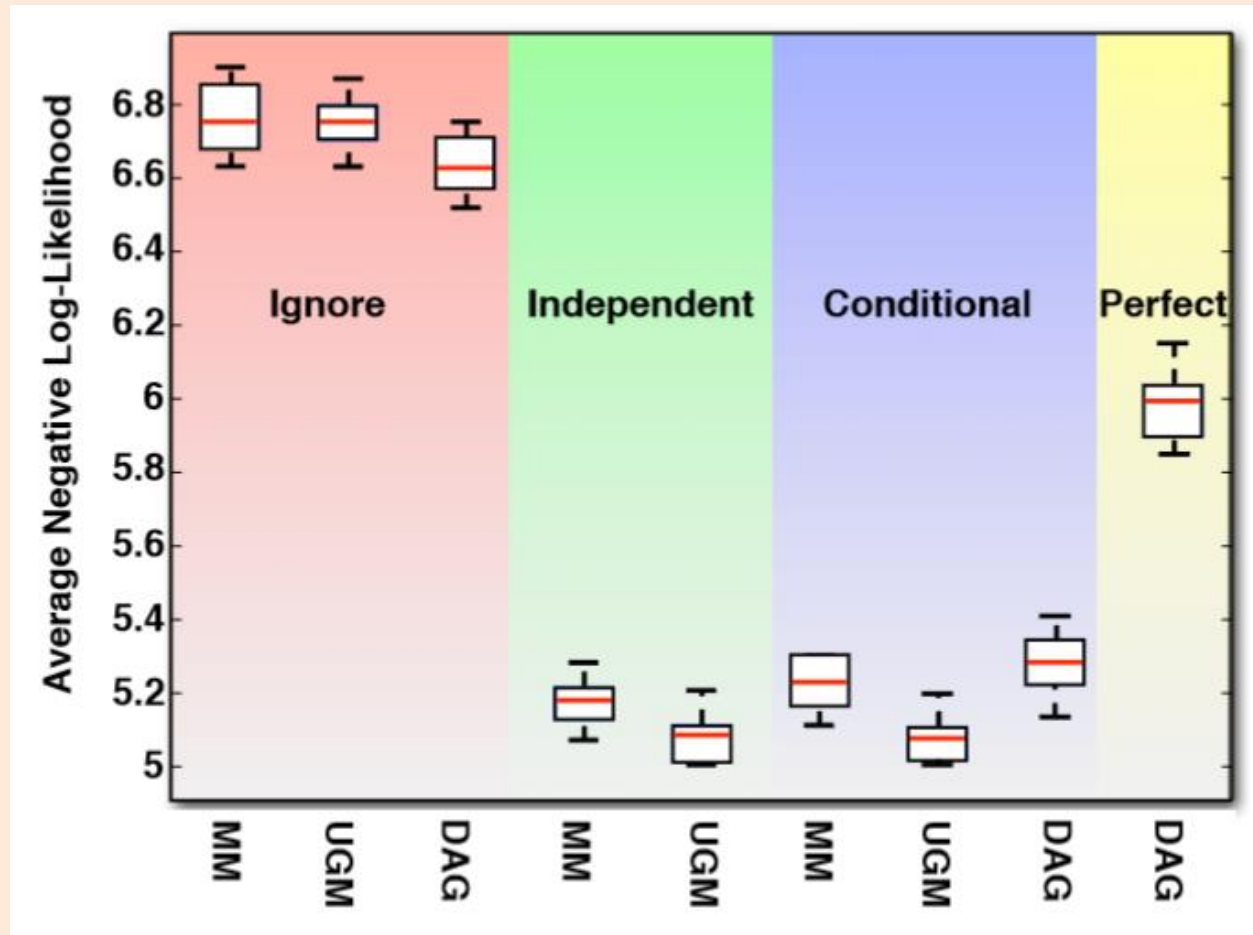
Histogram

- Histogram with grouping:



Box Plots

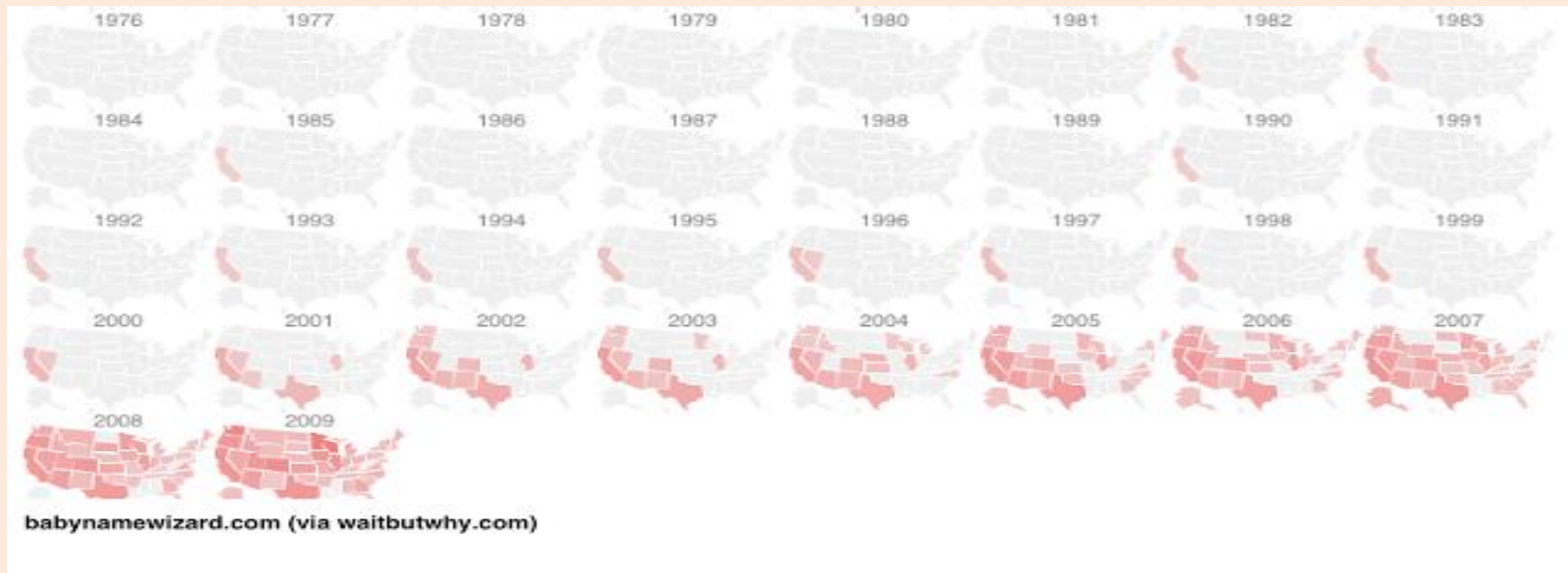
- Box plot with grouping:



Map Coloring

- Color/intensity can represent feature of region.

Popularity of naming baby “Evelyn” over time:

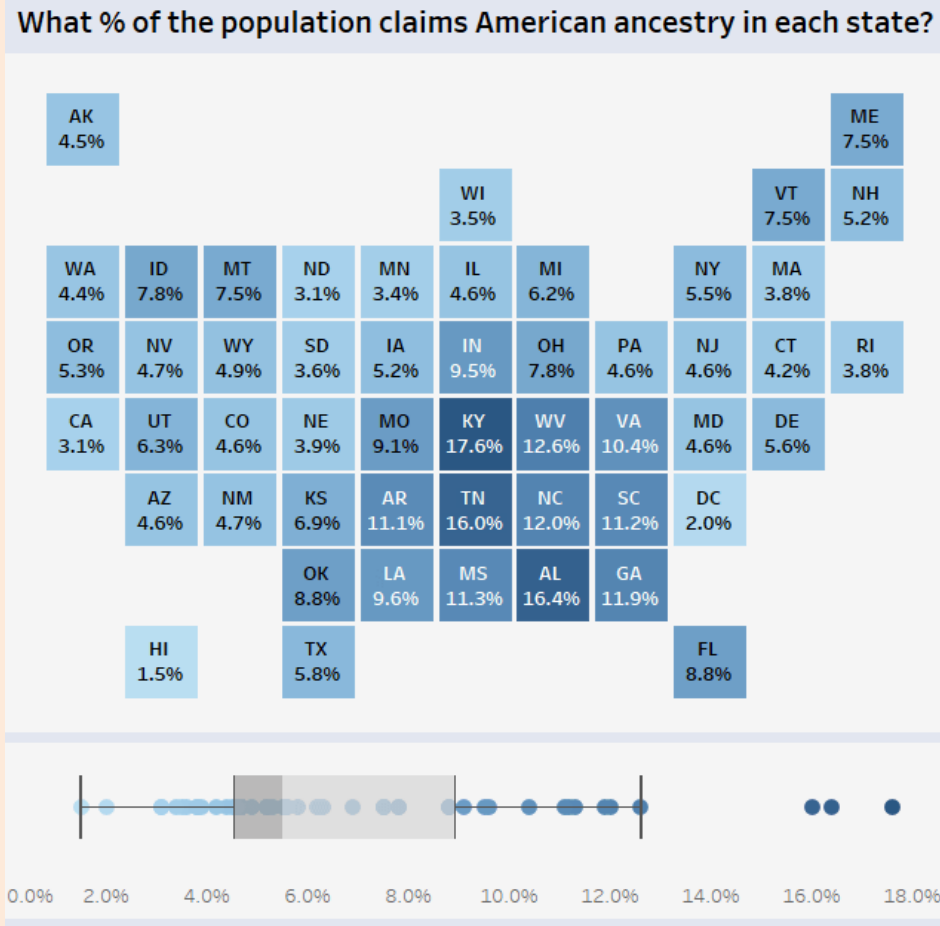


But not very good if some regions are very small.

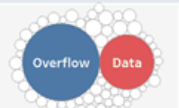
[Canadian Income Mobility](#)

Map Coloring

- Variation just uses fixed-size blocks and tries to arrange geographically:

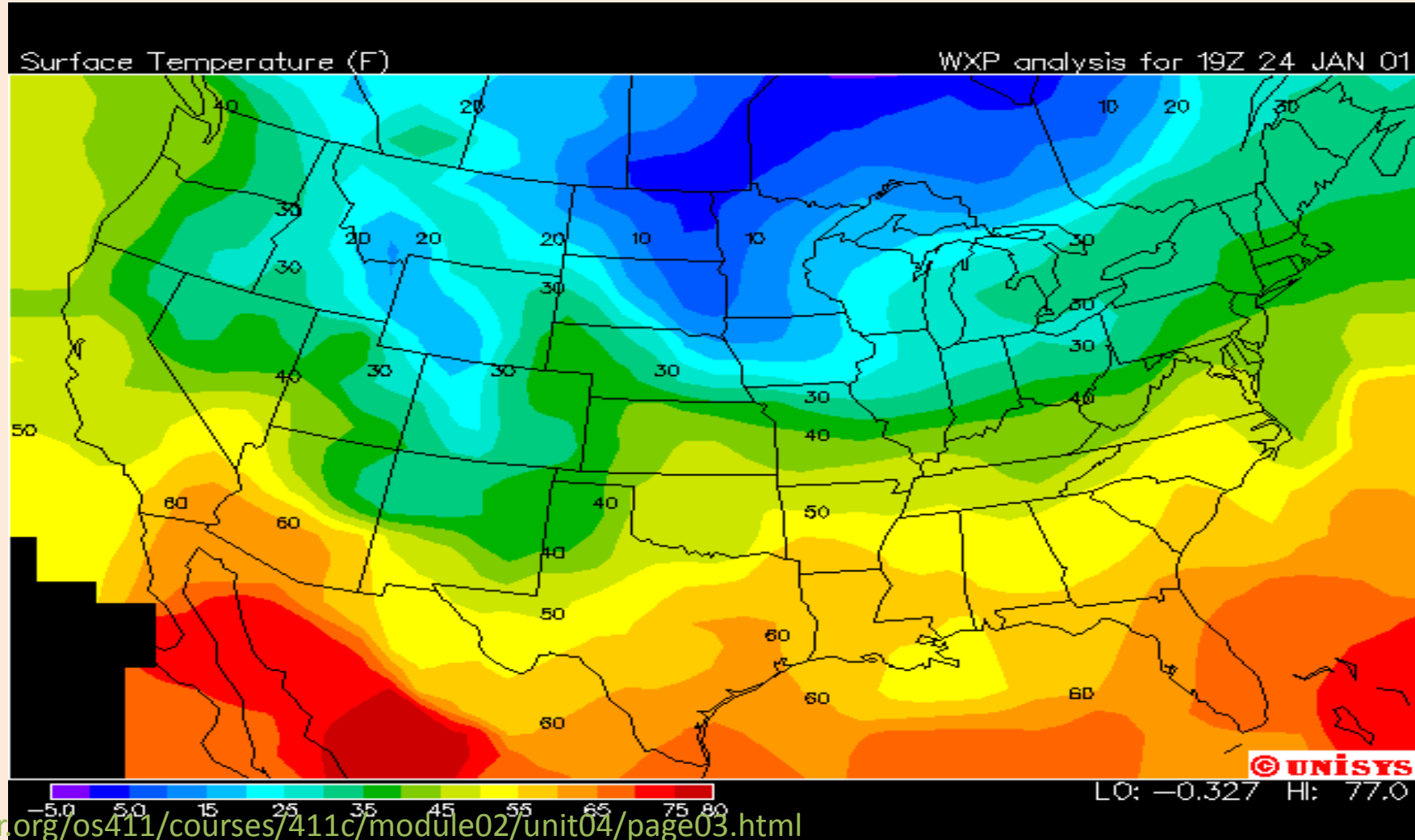


Source: U.S. Census Bureau, 2015 ACS 1 Year Estimates
Click the logo to see more data visualizations.



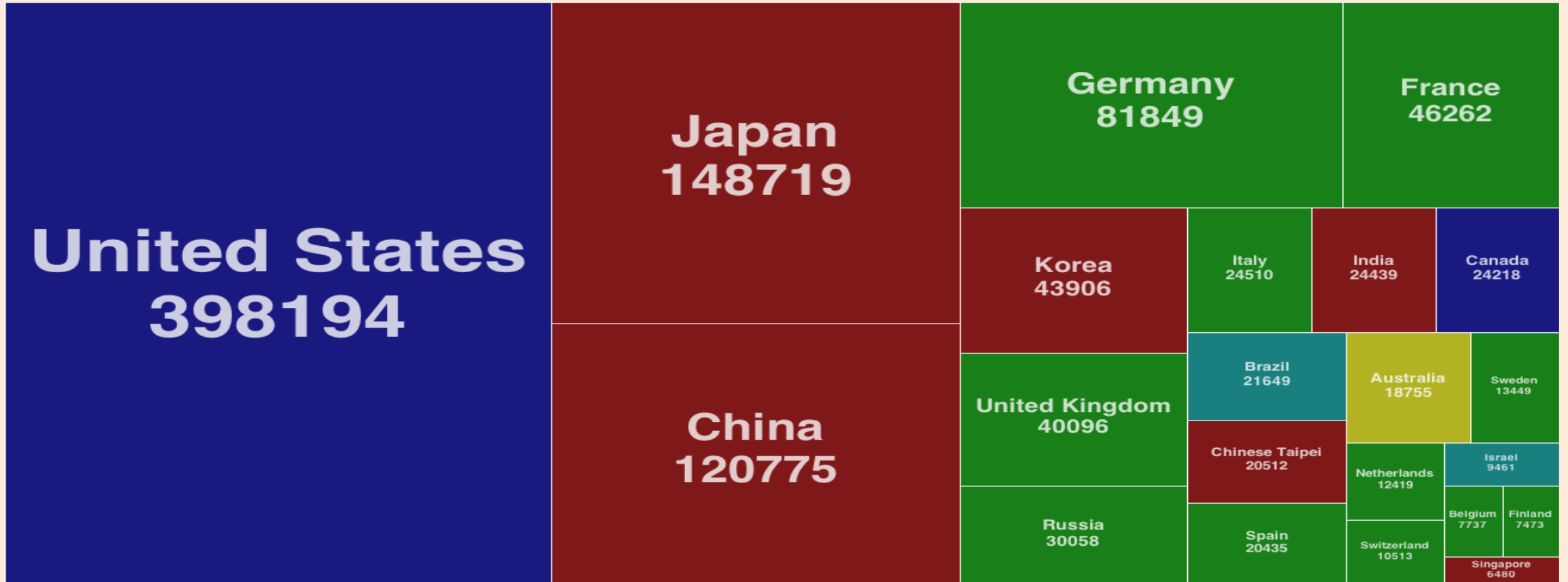
Contour Plot

- Colour visualizes 'z' as we vary 'x' and 'y'.



Treemaps

- Area represents attribute value:

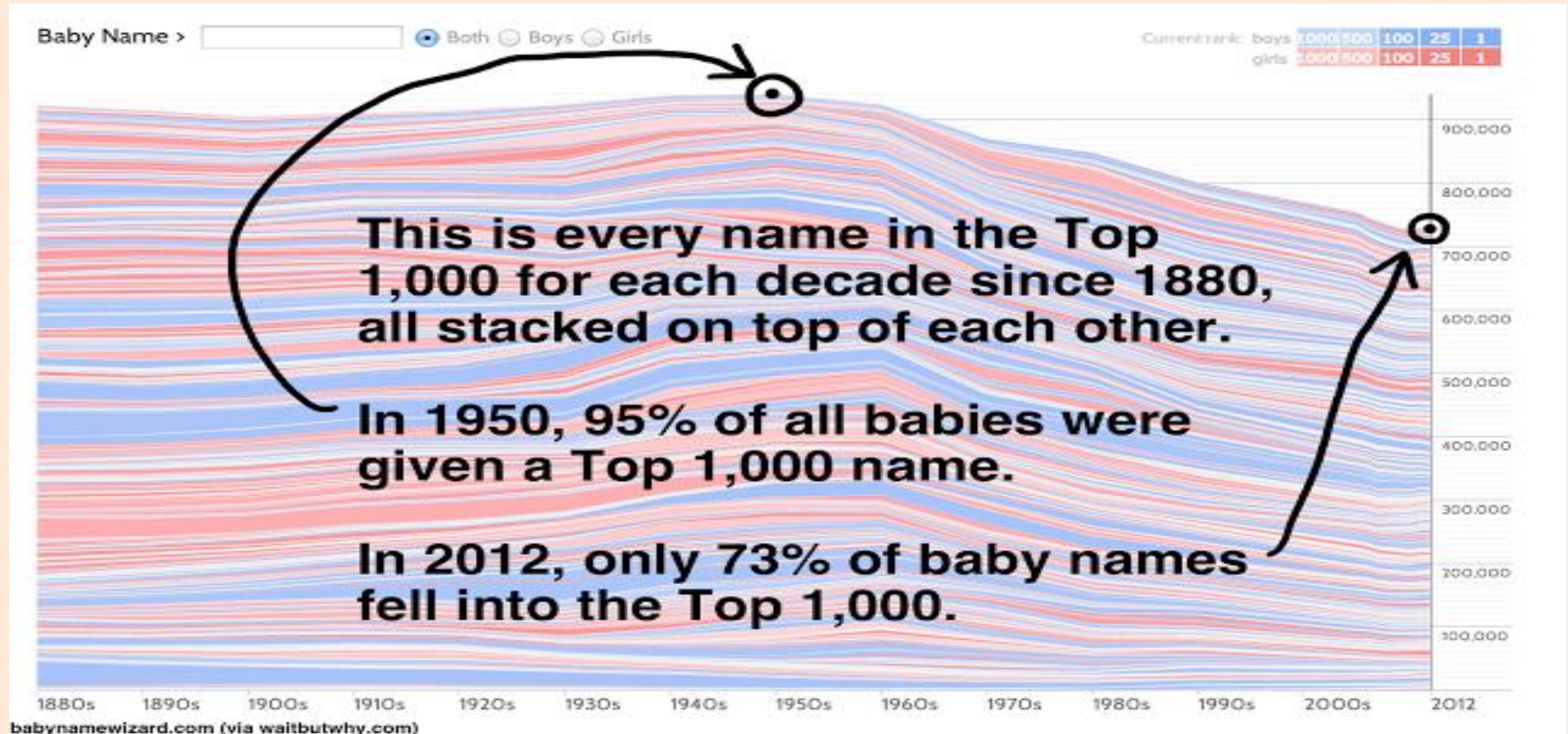


Cartogram

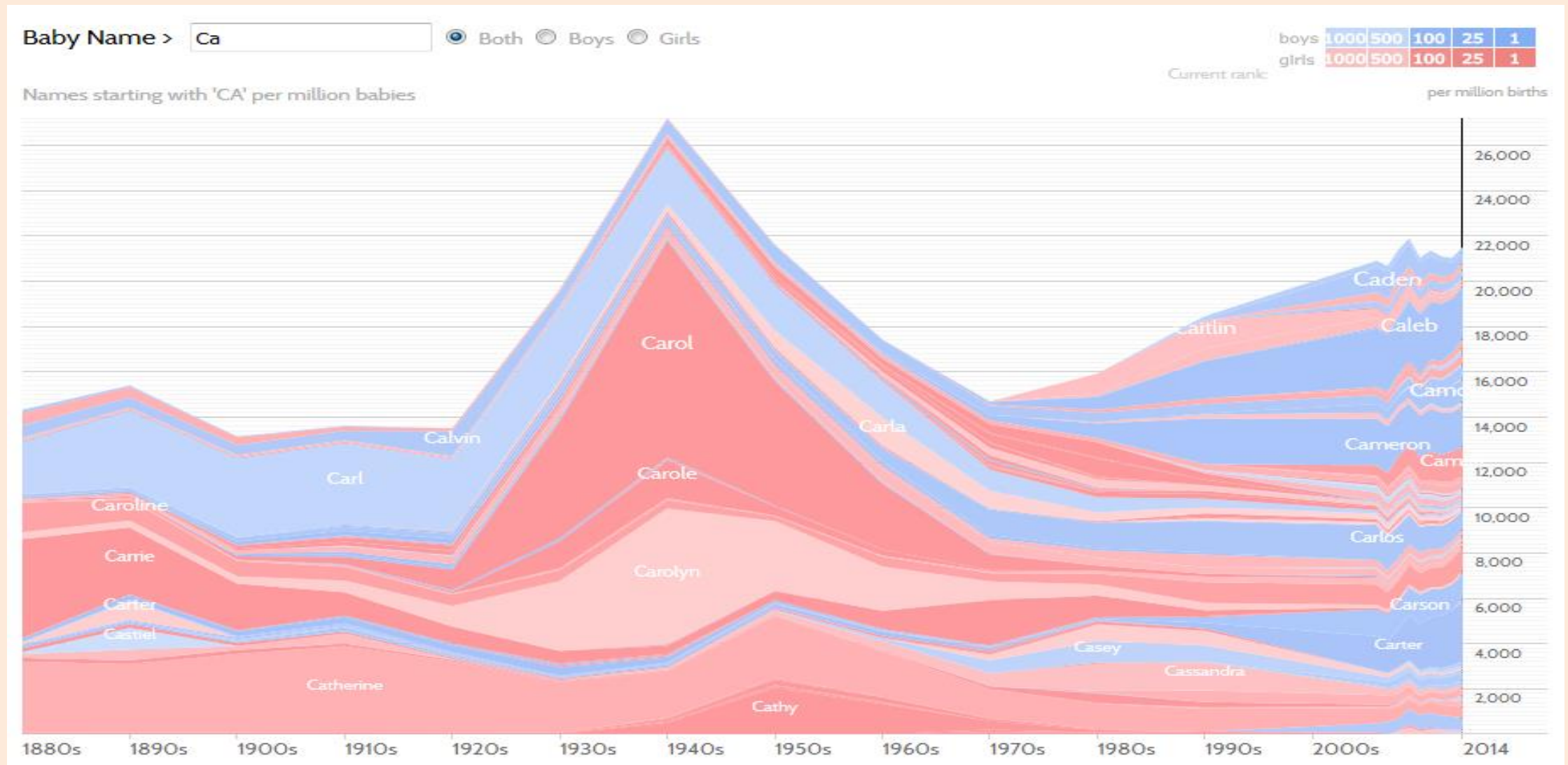
- Fancier version of treemaps:



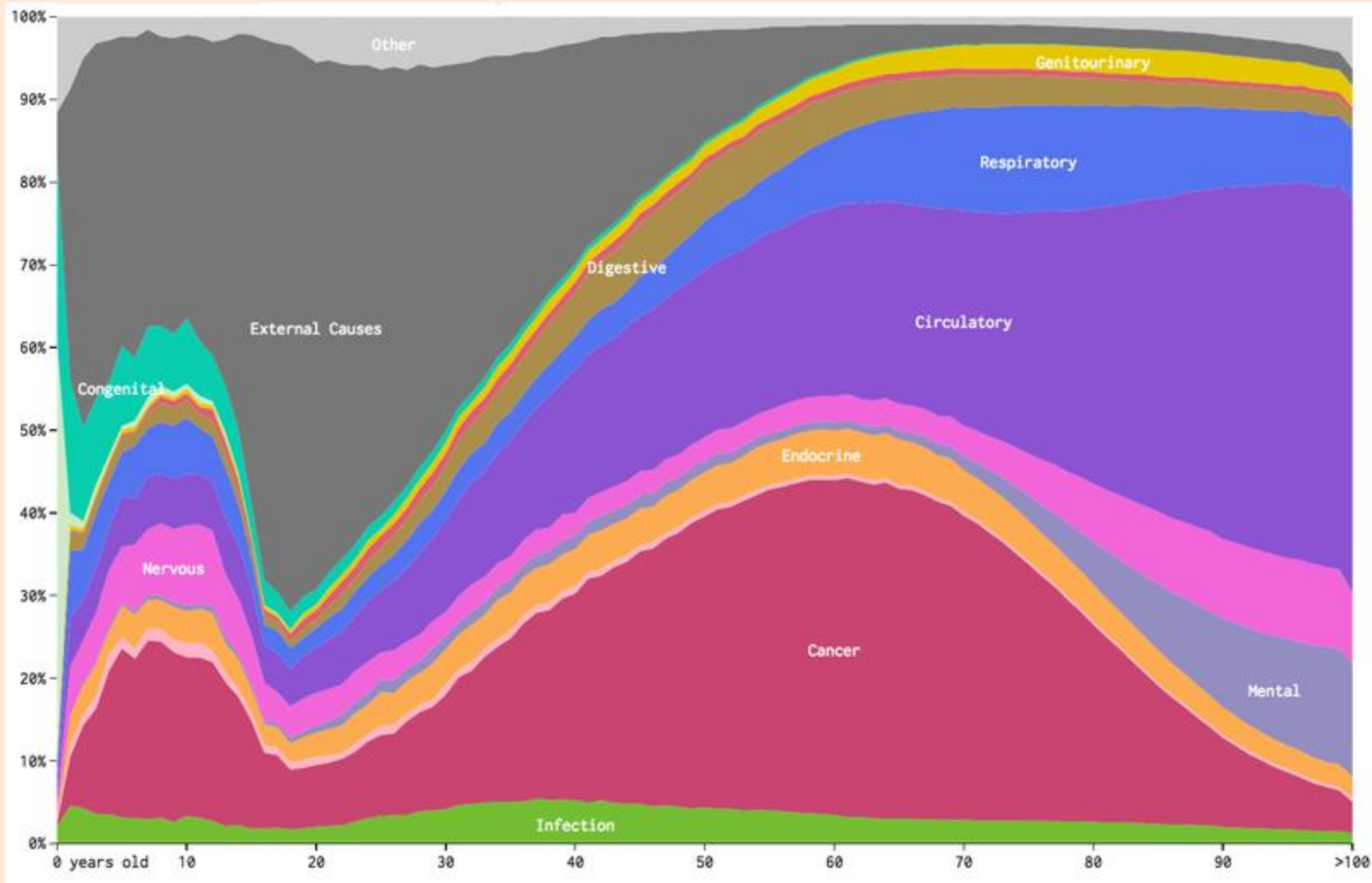
Stream Graph



Stream Graph



Stream Graph



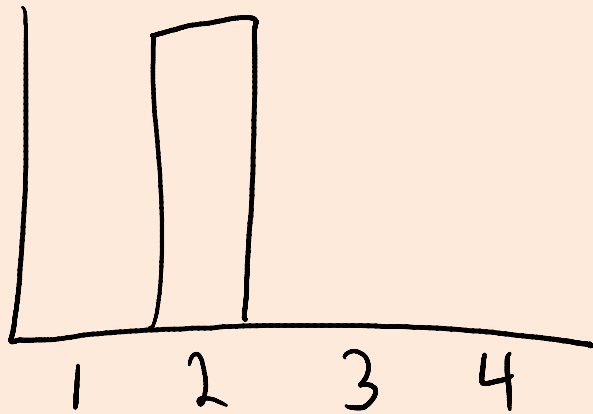
Videos and Interactive Visualizations

- For data recorded over time, videos can be useful:
 - [Map colouring over time.](#)
- There are also lots of neat interactive visualization methods:
 - [Sale date for most expensive paintings.](#)
 - [Global map of wind, weather, and oceans.](#)
 - [Many examples here.](#)

Entropy as Measure of Randomness

- **Entropy** measures “randomness” of a set of variables.
 - See [Wikipedia](#) for definition.

Low entropy means “very predictable”



High entropy means “very random”



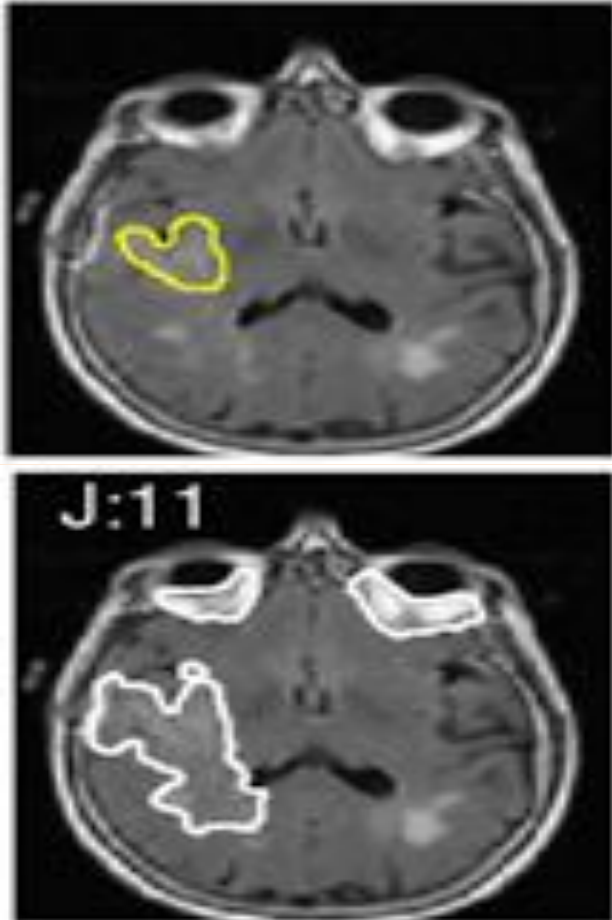
- For categorical features, the uniform distribution has the highest entropy.
- For continuous features with fixed mean and variance:
 - Normal distribution has highest entropy.
- [Entropy and Dr. Seuss.](#)

Hamming Distance vs. Jaccard Coefficient

A	B
1	0
1	0
1	0
0	1
0	1
1	0
0	0
0	0
0	1

- These vectors agree in 2 positions.
 - Normalizing Hamming distance by vector length, similarity is $2/9$.
- If we're really interested in predicting 1s, we could find set of 1s in both and compute Jaccard:
 - A \rightarrow {1,2,3,6}, B \rightarrow {4,5,9}
 - No intersection so Jaccard similarity is actually 0.

Hamming Distance vs. Jaccard Coefficient



- Let's say we want to find the tumour in an MR image.
- We have an expert label (top) and a prediction from our ML system (bottom).
- The normalized Hamming distance between the predictions at each pixel is 0.91. This sounds good, but since there are so many non-tumour pixels this is misleading.
- The ML system predicts a much bigger tumour so hasn't done well. The Jaccard coefficient between the two sets of tumour pixels is only 0.11 so reflects this.

Coupon Collecting

- Consider trying to collect 50 uniformly-distributed states, drawing at random.
- The probability of getting a new state if there 'x' states left: $p=x/50$.
- So expected number of samples before next "success" (getting a new state) is $50/x$.

(mean of geometric random variable with $p=x/50$)

- So the expected number of draws is the sum of $50/x$ for $x=1:50$.
- For 'n' states instead of 50, summing until you have all 'n' gives:

$$\sum_{i=1}^n \frac{n}{i} = n \sum_{i=1}^n \frac{1}{i} \leq n(1 + \log(n)) = O(n \log n)$$

Huge Datasets and Parallel/Distributed Computation

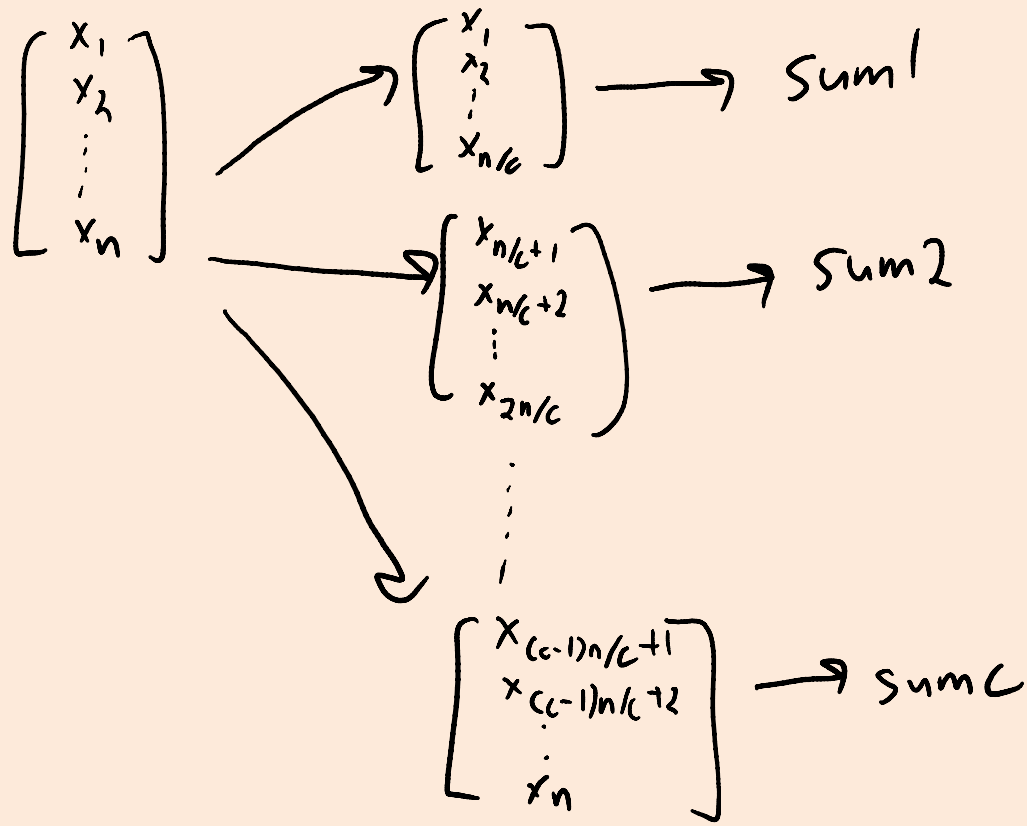
- Most **sufficient statistics can be computed in linear time.**
- For example, the mean of 'n' numbers is computed as:

$$\text{mean}(x_1, x_2, x_3, \dots, x_n) = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

- This costs $O(n)$, which is great.
- But if 'n' is really big, we can go even faster with parallel computing...

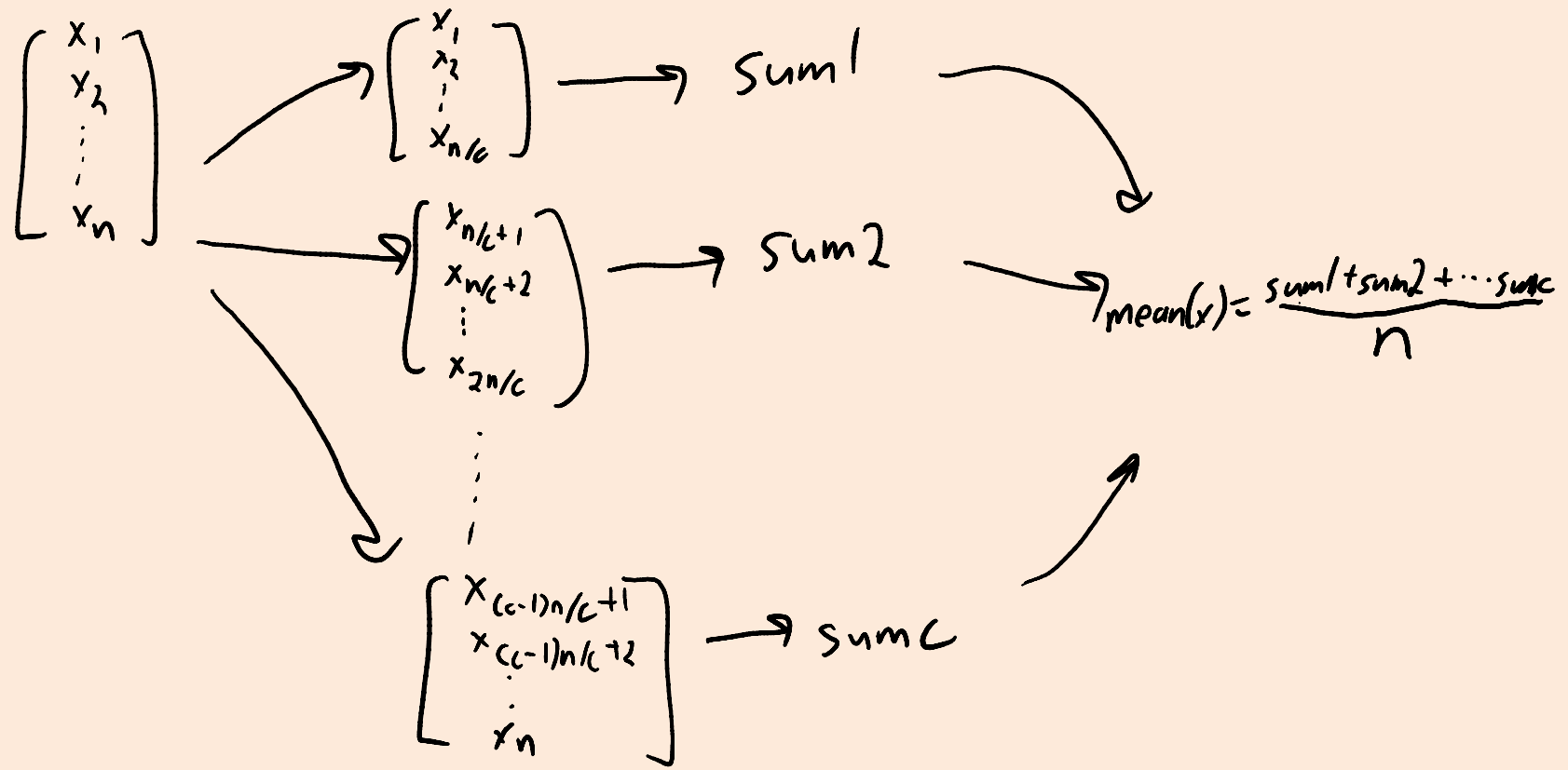
Huge Datasets and Parallel/Distributed Computation

- Computing the mean with **multiple cores**:
 - Each of the 'c' cores computes the sum of $O(n/c)$ of the data:



Huge Datasets and Parallel/Distributed Computation

- Computing the mean with **multiple cores**:
 - Each of the 'c' cores computes the sum of $O(n/c)$ of the data:
 - Add up the 'c' results from each core to get the mean.



Huge Datasets and Parallel/Distributed Computation

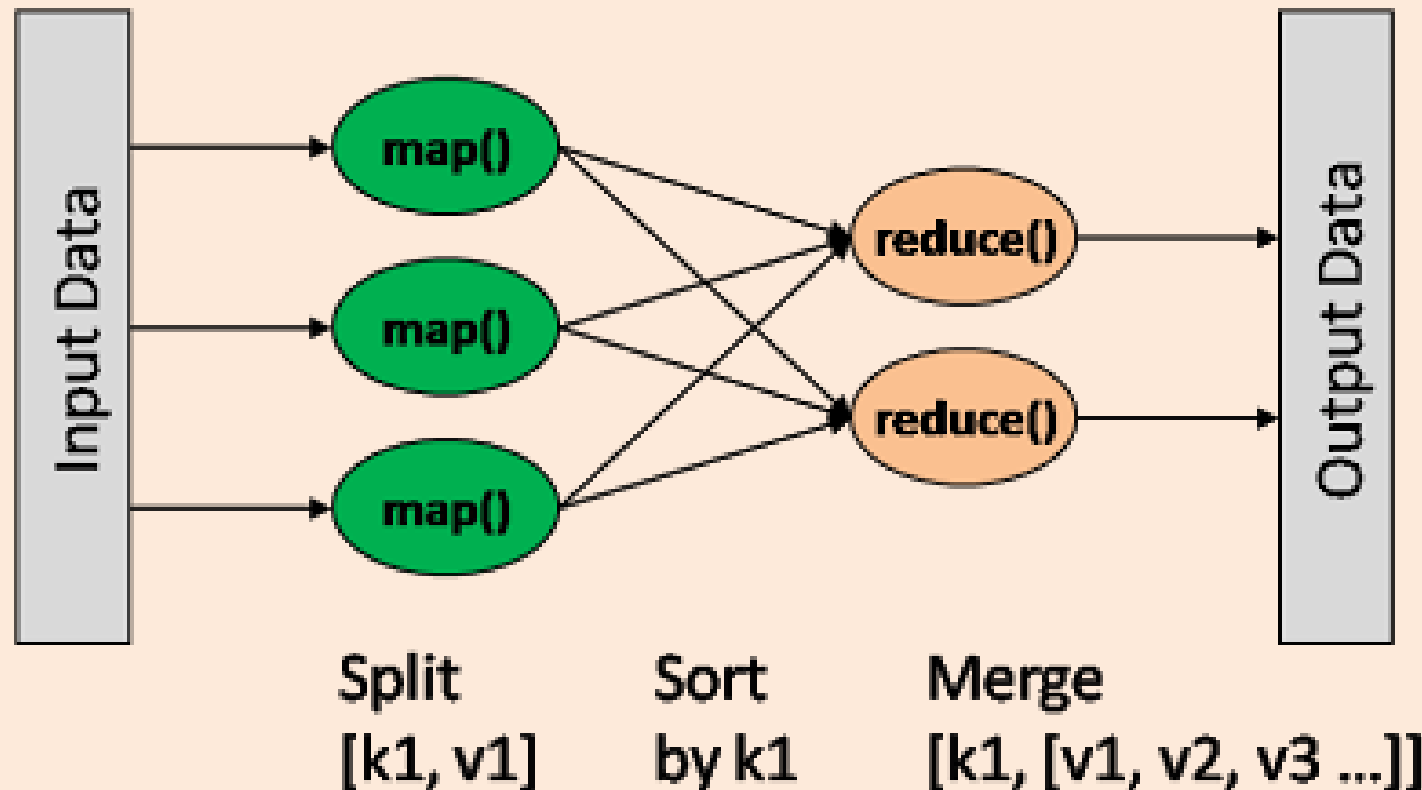
- Computing the mean with **multiple cores**:
 - Each of the 'c' cores computes the sum of $O(n/c)$ of the data.
 - Add up the 'c' results from each core to get the mean.
 - Cost is only $O(n/c + c)$, which can be much faster for large 'n'.
- This assumes cores can access data in parallel (not always true).
- Can reduce cost to $O(n/c)$ by having cores write to same register.
 - But need to “lock” the register and might effectively cost $O(n)$.

Huge Datasets and Parallel/Distributed Computation

- Sometimes 'n' is so big that **data can't fit on one computer**.
- In this case the data might be distributed across 'c' machines:
 - Hopefully, each machine has $O(n/c)$ of the data.
- We can solve the problem similar to the multi-core case:
 - “**Map**” step: each machine computes the sum of its data.
 - “**Reduce**” step: each machine communicates sum to a “master” computer, which adds them together and divides by 'n'.

Huge Datasets and Parallel/Distributed Computation

- Many problems in DM and ML have this flavour:
 - “Map” computes an operation on the data on each machine (in parallel).
 - “Reduce” combines the results across machines.



Huge Datasets and Parallel/Distributed Computation

- Many problems in DM and ML have this flavour:
 - “Map” computes an operation on the data on each machine (in parallel).
 - “Reduce” combines the results across machines.
 - These are standard operations in parallel libraries like [MPI](#).
- Can solve many problems almost ‘c’ times faster with ‘c’ computers.
- To make it up for the **high cost communicating across machines**:
 - Assumes that most of the computation is in the “map” step.
 - Often need to assume data is already on the computers at the start.

Huge Datasets and Parallel/Distributed Computation

- Another challenge with “Google-sized” datasets:
 - You may need so many computers to store the data, that it’s **inevitable that some computers are going to fail**.
- Solution to this is a **distributed file system**.
- Two popular examples are Google’s MapReduce and Hadoop DFS:
 - Store data with redundancy (same data is stored in many places).
 - And assume data isn’t changing too quickly.
 - Have a strategy for restarting “map” operations on computers that fail.
 - Allows fast calculation of more-fancy things than sufficient statistics:
 - Database queries and matrix multiplications.

Data Clean and the Duke Cancer Scandal

- See the Duke cancer scandal:
 - http://www.nytimes.com/2011/07/08/health/research/08genes.html?_r=2&hp
- Basic sanity checks for data cleanliness show problems in these (and many other) studies:
 - E.g., flipped labels, off-by-one mistakes, switched columns etc.
 - <https://arxiv.org/pdf/1010.1092.pdf>