

# Admin

- **Assignment 1** is out (due next Friday, start early).
- Tutorials start Monday:
  - 11am, 2pm, and 4pm in DMP 201.
  - **New tutorial section**: 5pm in DMP 101.
  - Make sure you sign up for one.
  - No requirement to attend, but helps with assignment.
- Office hours:
  - Being **moved to Tuesdays at 10 and 4**.
  - Watch the website for details.
- Sign up for the course page on **Piazza**.

**CPSC 340:**  
**Machine Learning and Data Mining**

Data Preprocessing and Exploration

September 11, 2015

# Outline

- 1) Typical steps in knowledge discovery from data.
- 2) Data Representations
- 3) Data Exploration

These notes roughly follow:

[http://www-users.cs.umn.edu/~kumar/dmbook/dmslides/chap2\\_data.pdf](http://www-users.cs.umn.edu/~kumar/dmbook/dmslides/chap2_data.pdf)

# 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)

# Outline

- 1) Typical steps in knowledge discovery from data.
- 2) Data Representations**
- 3) Data Exploration

# 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.



# Types of Data

- **Discrete features** come from an unordered set:
  - Binary: job?
  - Nominal/categorical: 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).

# Bag of Words

- **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.

# Other Data Types

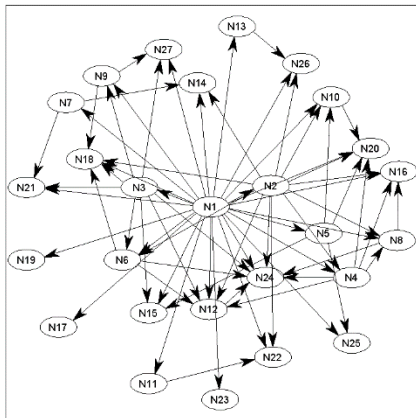
- We can think of other data types in this way:
  - Images:



→  
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
0	1	1	1	1	1	1
0	0	0	1	0	1	0
0	0	0	0	0	1	0
0	0	0	0	0	0	0

# 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 (exact of otherwise).
- 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.

# IID Assumption

- Almost all of ML+DM assumes objects are IID:
  - “Independent and identically distributed”.
  - The order of the objects doesn’t matter.
  - New objects will behave like the existing objects.
- The IID assumption implies that our conclusions will probably also apply to new data.
- This assumption is rarely true:
  - But it is often a good approximation.
  - Some methods relax this assumption.

# How much data do we need?

- Assume we have a categorical variable with 50 values: {Alabama, Alaska, Arizona, Arkansas,...}.
- We can turn this into 50 binary variables.
- If each category has equal probability, **how many IID objects before we see each category once?**
- Expected value is  $\sim 225$ .
- Coupon collector problem:  $O(n \log n)$  in general.
- **Need more data than categories:**
  - Situation is worse if don't have equal probabilities.
  - Typically want to see categories more than once.

# 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

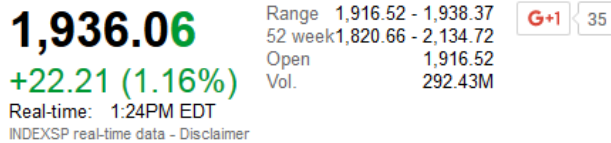
- Feature Selection:
  - Remove features that are not relevant to the task.

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

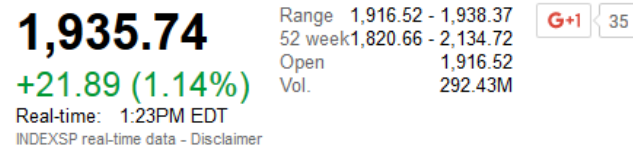
- Student ID is probably not relevant.

# Feature Transformation

- Mathematical transformations:
  - Square, exponentiation, or take logarithm.



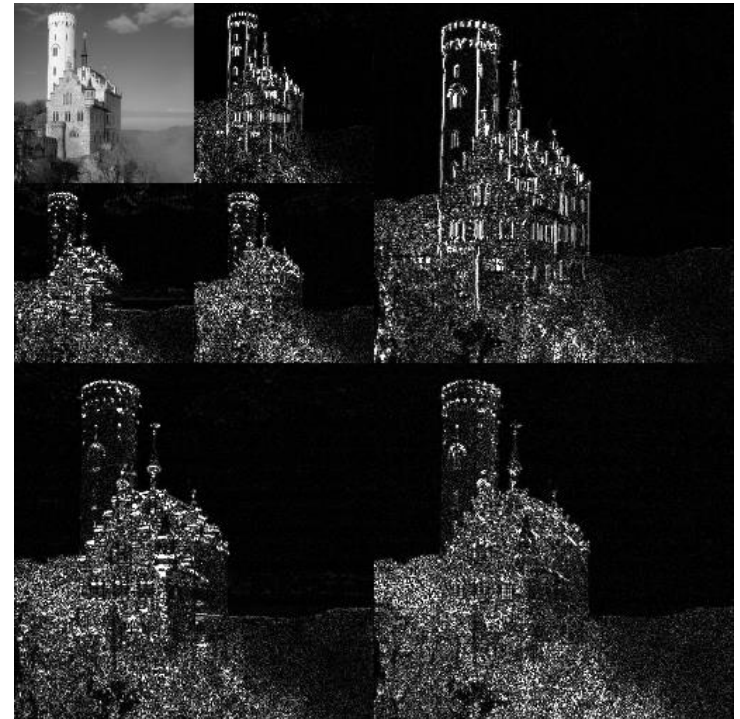
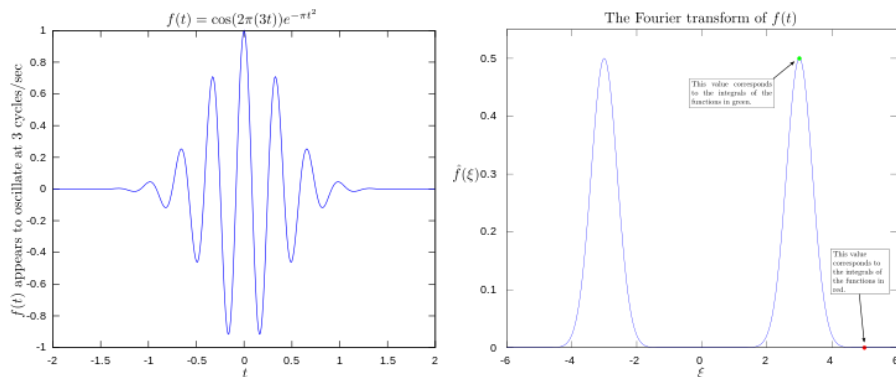
[Settings](#) | [Technicals](#) | [Link to this view](#)



[Settings](#) | [Technicals](#) | [Link to this view](#)

# Feature Transformation

- Mathematical transformations:
  - Square, exponentiation, or take logarithm.
  - Fourier or wavelet transform (signal data).



# Feature Transformation

- Mathematical transformations:
  - Square, exponentiation, or take logarithm.
  - Fourier or wavelet transform (signal data).
  - **Discretization**: turn continuous into discrete.

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
21	0	1	0

# Feature Transformation

- Mathematical transformations:
  - Square, exponentiation, or take logarithm.
  - Fourier or wavelet transform (signal data).
  - **Discretization**: turn continuous into discrete.
  - Scaling: convert variables to comparable scales (E.g., convert kilograms to grams.)

# Outline

- 1) Typical steps in knowledge discovery from data.
- 2) Data Representations
- 3) **Data Exploration**

# Data Exploration

- 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).

# Discrete Summary Statistics

- Summary statistics for a discrete 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
<b>Canada</b>	<b>35,540.4</b>
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.



# Discrete Summary Statistics

- Summary statistics **between** discrete variables:
  - **Simple matching** coefficient:
    - How many times two variables are the same.
    - $SMC = (C_{11} + C_{00}) / (C_{00} + C_{01} + C_{10} + C_{11})$ .
  - **Jaccard** coefficient for binary variables:
    - Intersection divided by union of '1' values.
    - $C_{11} / (C_{01} + C_{10} + C_{11})$ .

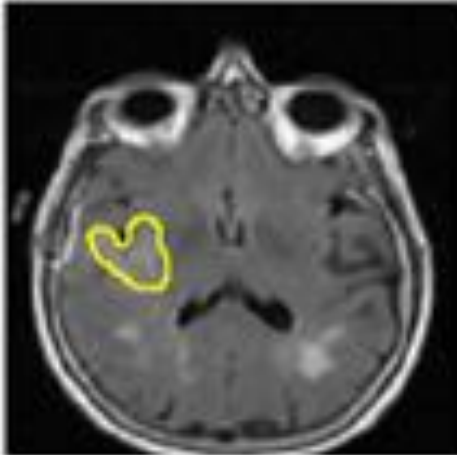
# Simple Matching vs. Jaccard

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

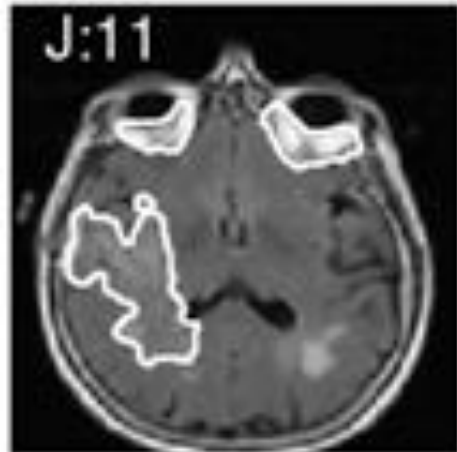
$$\begin{aligned}\text{Sim}(A,B) &= (C_{11} + C_{00}) / (C_{00} + C_{01} + C_{10} + C_{11}) \\ &= (0 + 2) / (2 + 3 + 3 + 0) \\ &= 2/7.\end{aligned}$$

$$\begin{aligned}\text{Jac}(A,B) &= C_{11} / (C_{01} + C_{10} + C_{11}) \\ &= 0 / (3 + 3 + 0) \\ &= 0.\end{aligned}$$

# Simple Matching vs. Jaccard



$$\text{Sim}(A,B) = 0.91$$



$$\text{Jac}(A,B) = 0.11$$

# Continuous Summary Statistics

- Measures of location:
  - **Mean**: average value (sensitive to outliers).
  - **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.
  - **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
- N.B.: mean and std are more sensitive to outliers.

# Continuous Summary Statistics

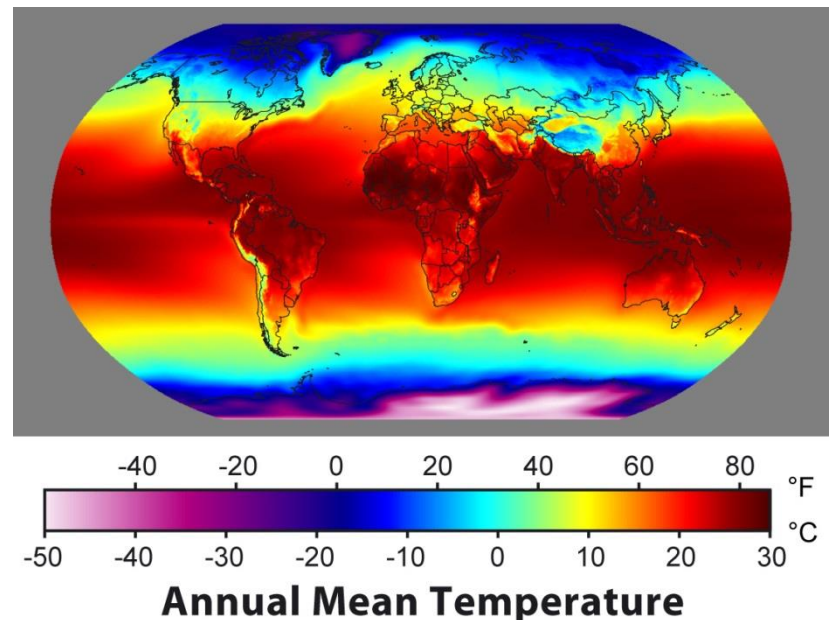
- Measures **between** continuous variables:
  - **Correlation:**
    - Does one increase/decrease proportionally as the other increases?
  - **Rank correlation:**
    - Does one increase/decrease as the other increases?
  - **Euclidean distance:**
    - How far apart are the values?
  - **Cosine similarity:**
    - What is the angle between them?

# Visualization

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

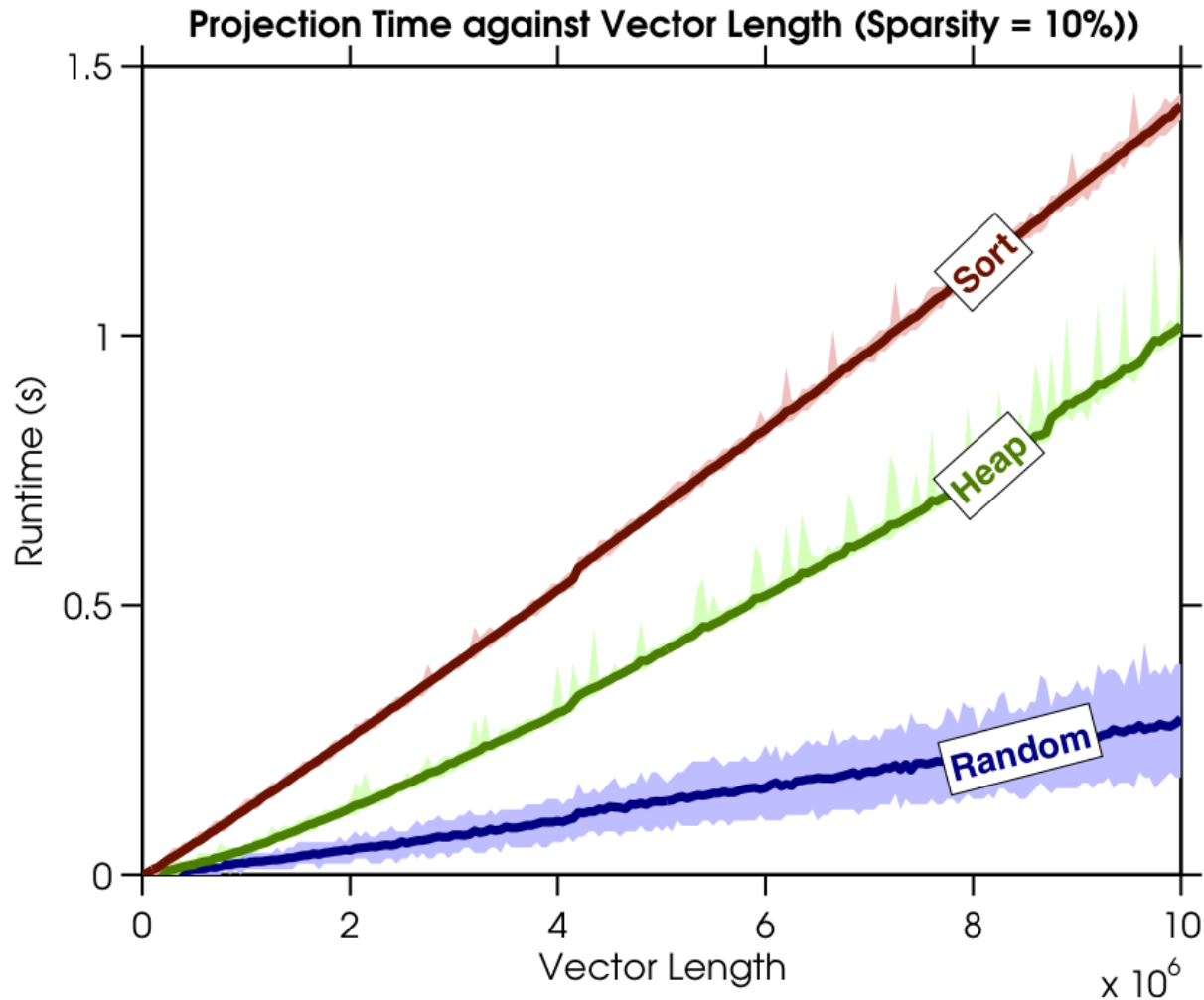
Lat	Long	Temp
0	0	30.1
0	1	29.8
0	2	29.9
0	3	30.1
0	4	29.9
...	...	...

vs.



# Basic Plot

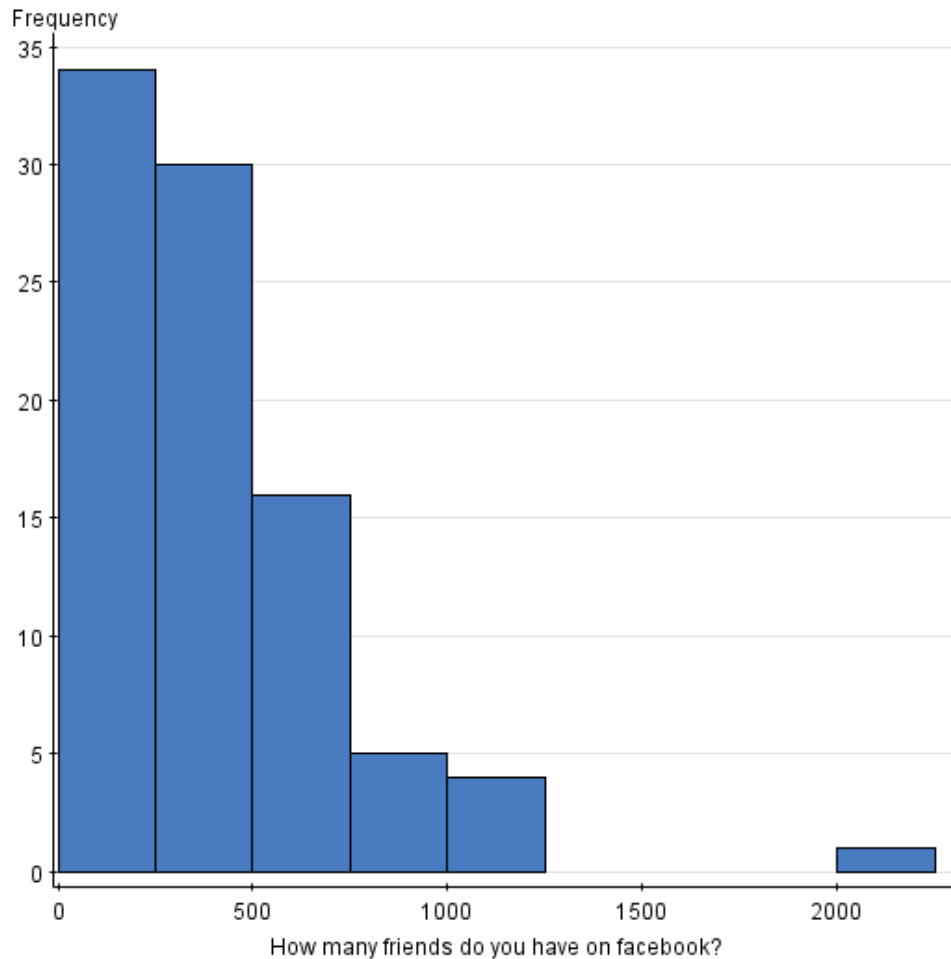
- Visualize one variable as a function of another.



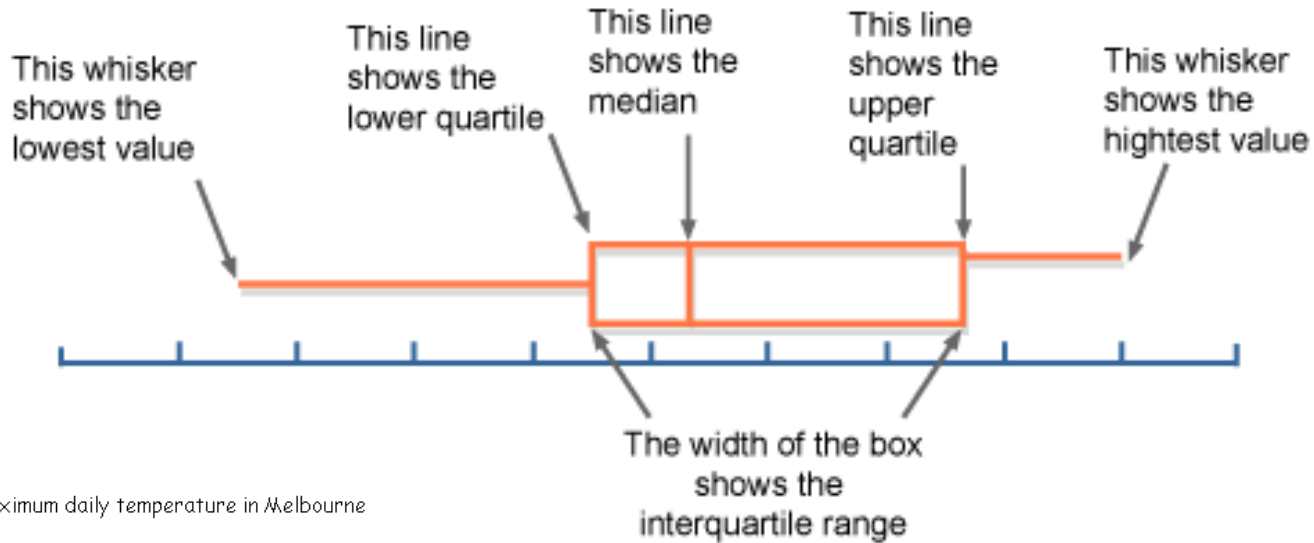


# Histogram

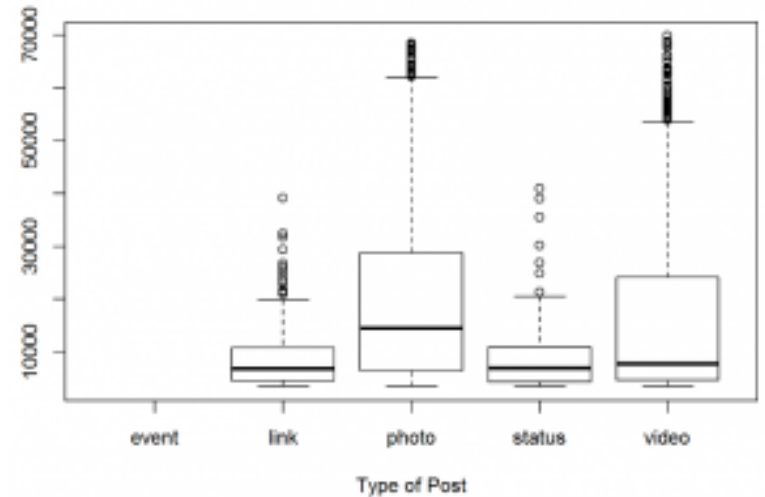
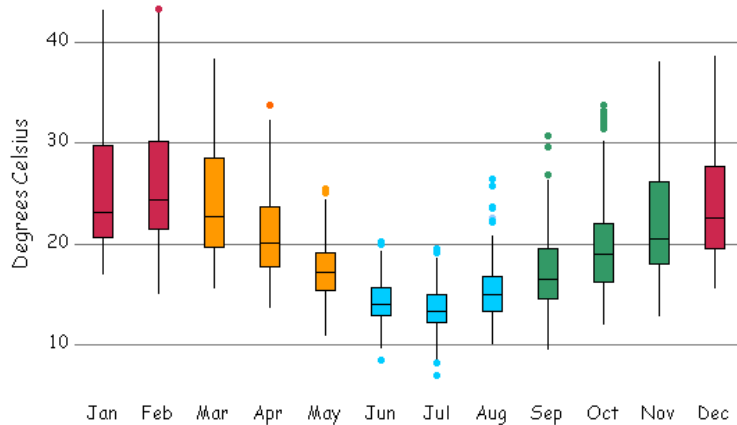
- Histograms display distribution of a variable.



# Box Plot



Maximum daily temperature in Melbourne



<http://www.bbc.co.uk/schools/gcsebitesize/maths/statistics/representingdata3hirev6.shtml>

<http://www.scc.ms.unimelb.edu.au/whatisstatistics/weather.html>

<http://r.ramganalytics.com/r/facebook-likes-and-analytics/>

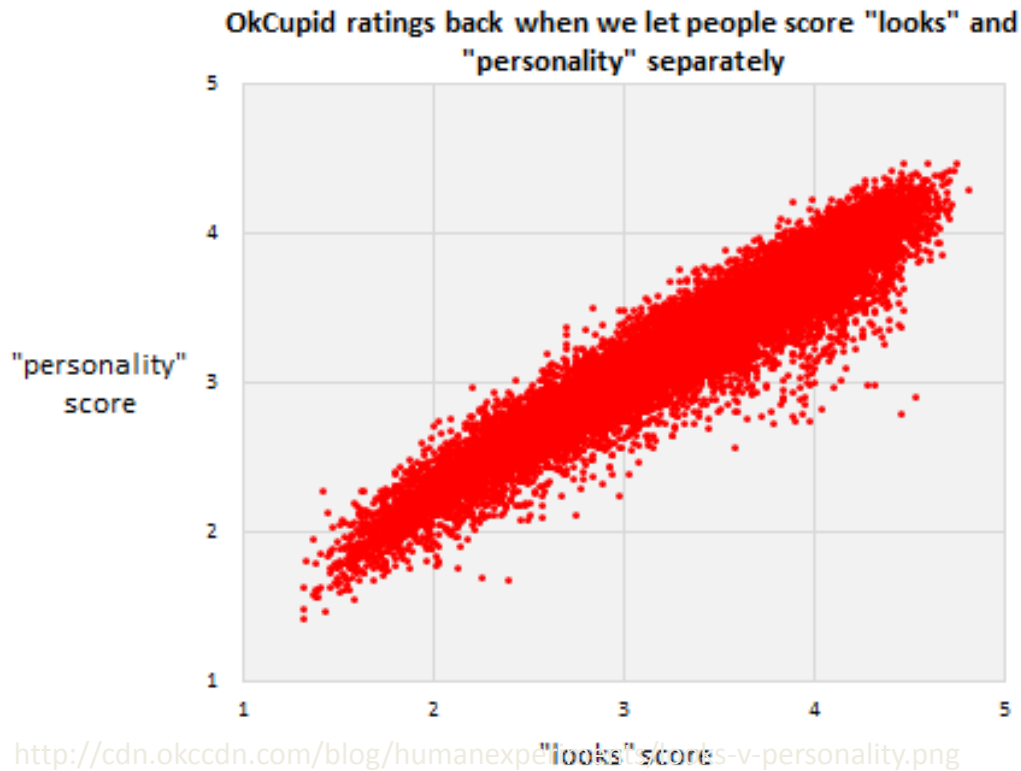
# Box Plot

- Photo from CTV Olympic coverage in 2010:



# Scatterplot

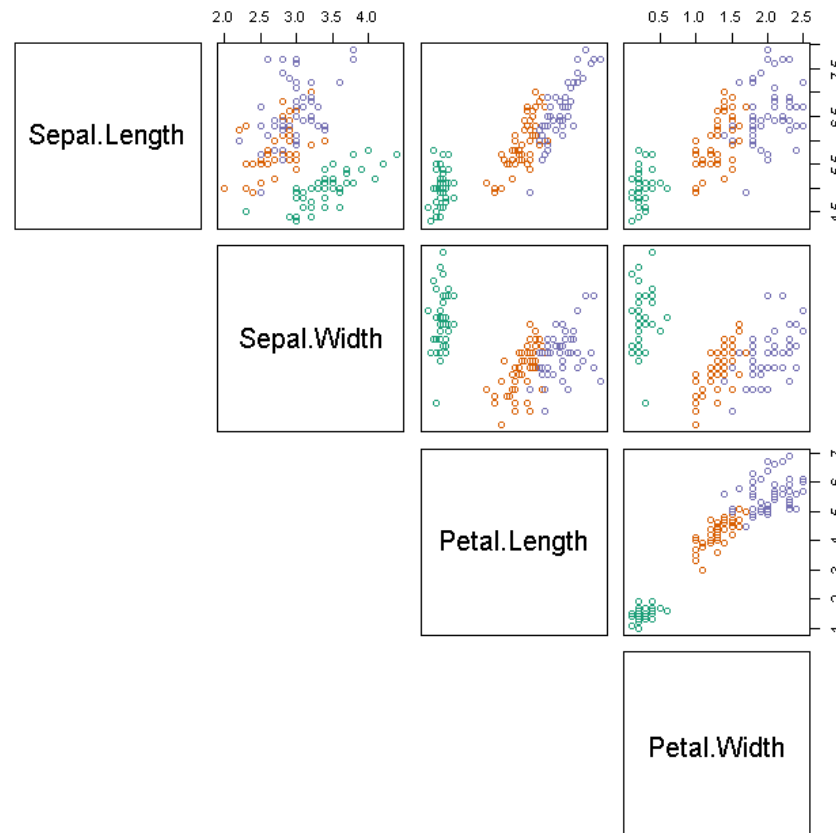
- Look at distribution of two features:
  - Feature 1 on x-axis.
  - Feature 2 on y-axis.



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

# Scatterplot Arrays

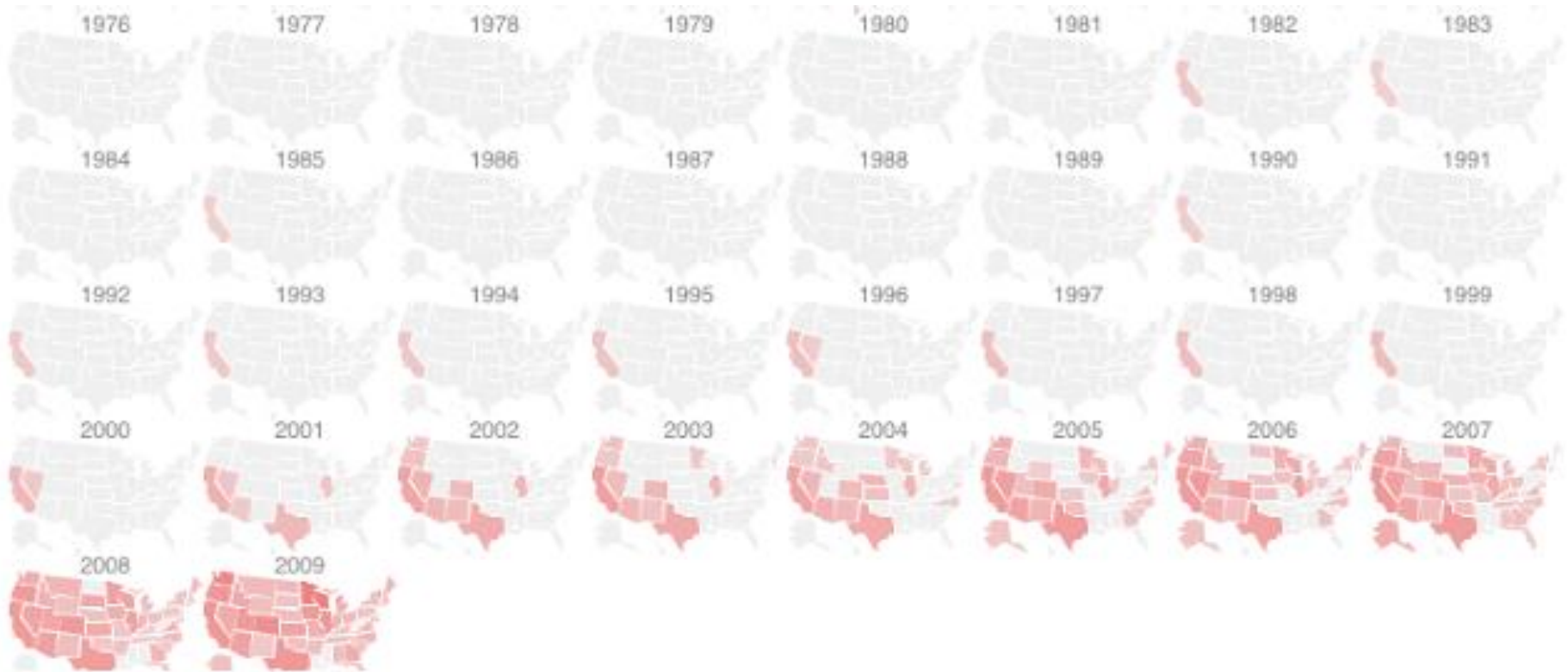
- Can plot multiple variables in an array.
- Colors can indicate a third categorical variable.



# Map Coloring

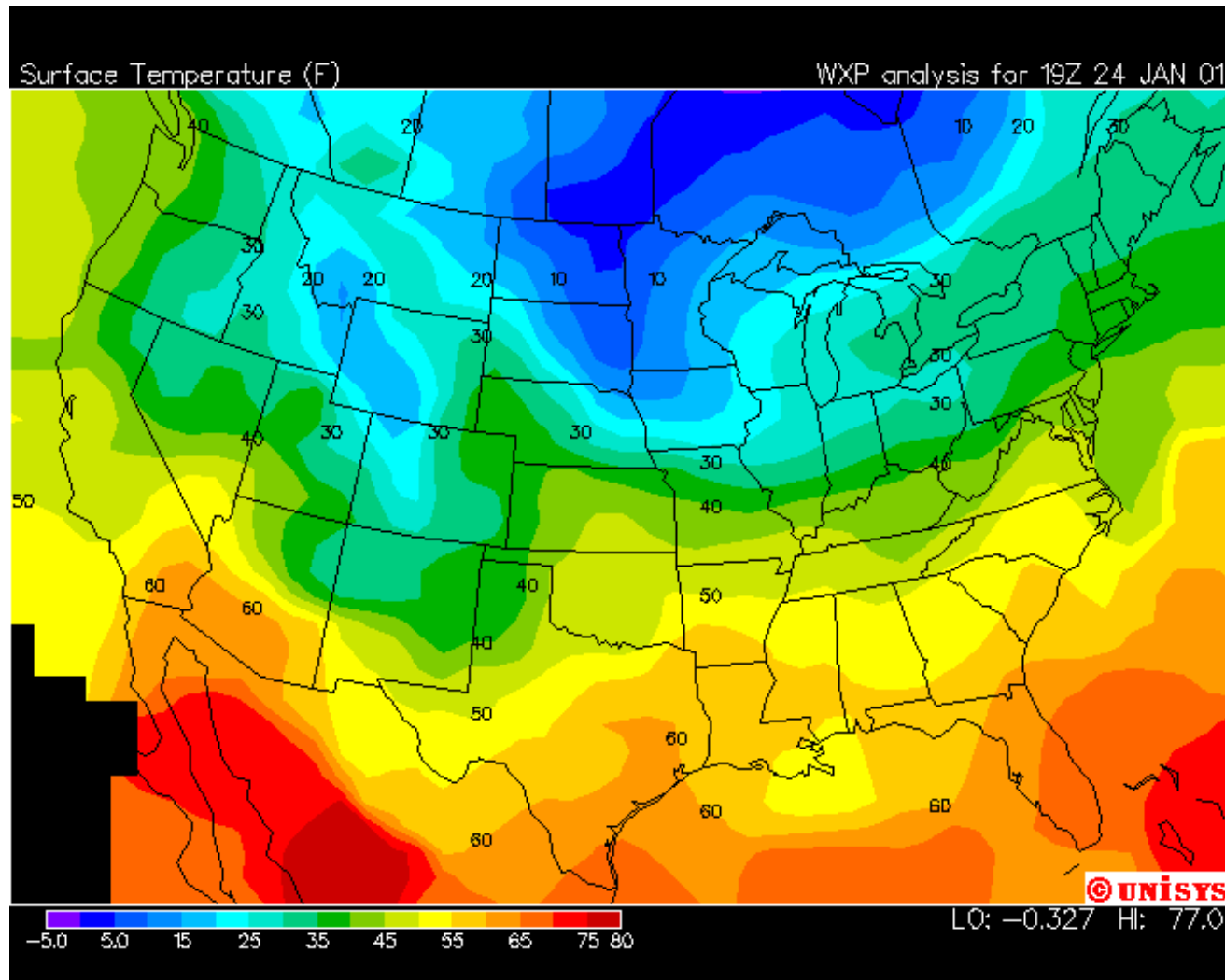
- Color/intensity can represent feature of region.

Evelyn's popularity over time:



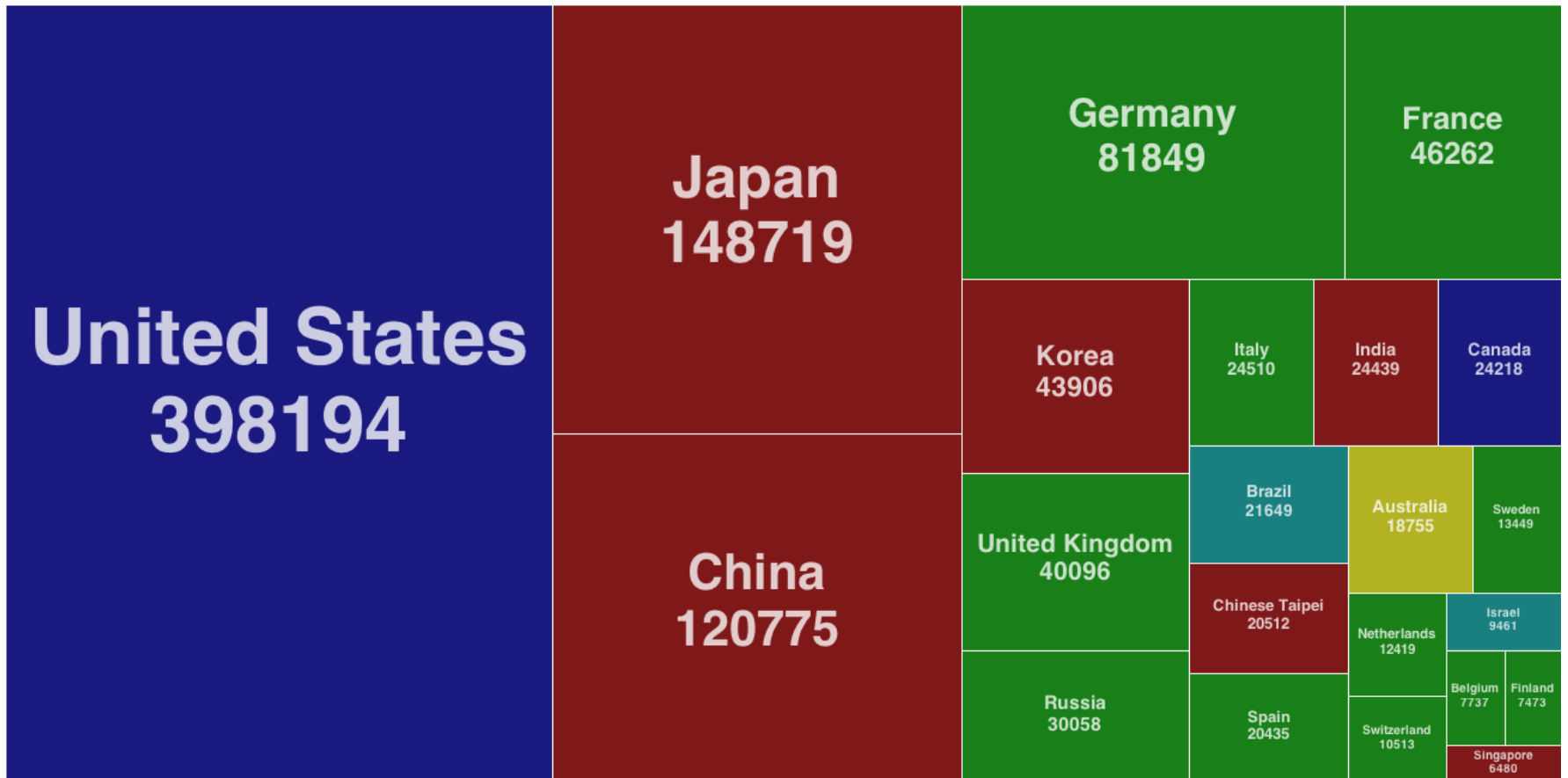
**[babynamewizard.com](http://babynamewizard.com) (via [waitbutwhy.com](http://waitbutwhy.com))**

# Contour Plot



# Treemaps

- Area represents attribute value:





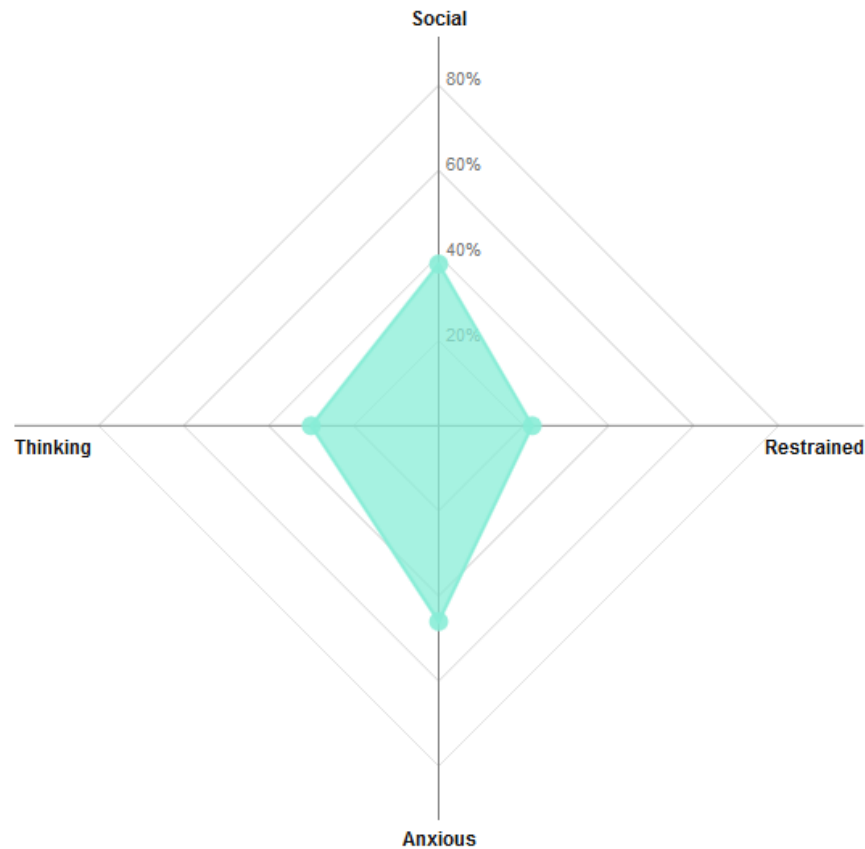
# Cartogram

- Fancier version of treemaps:

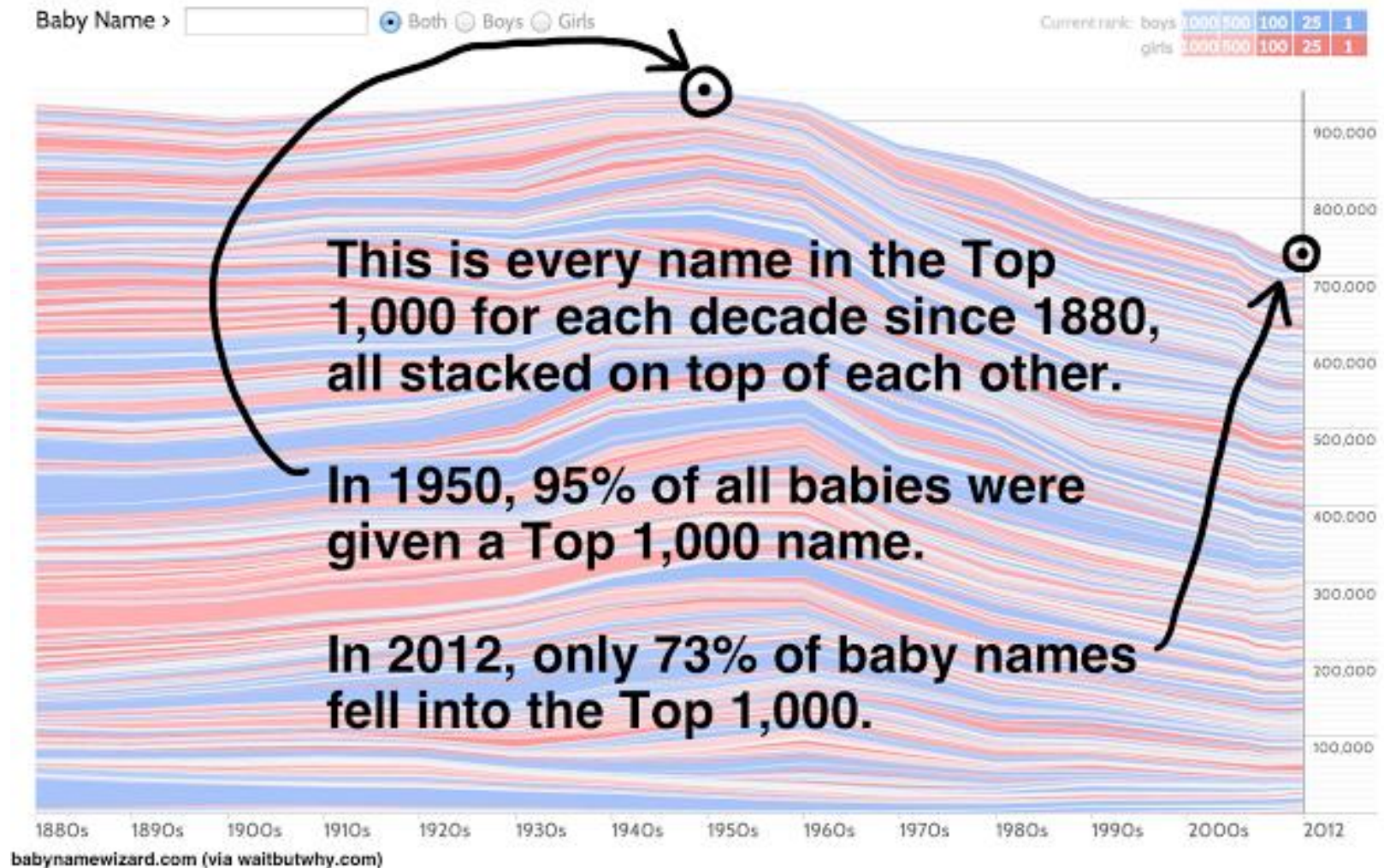


# Star Plots.

- Display values of several numerical variables.



# Stream Graph





# Stream Graph

Baby Name >

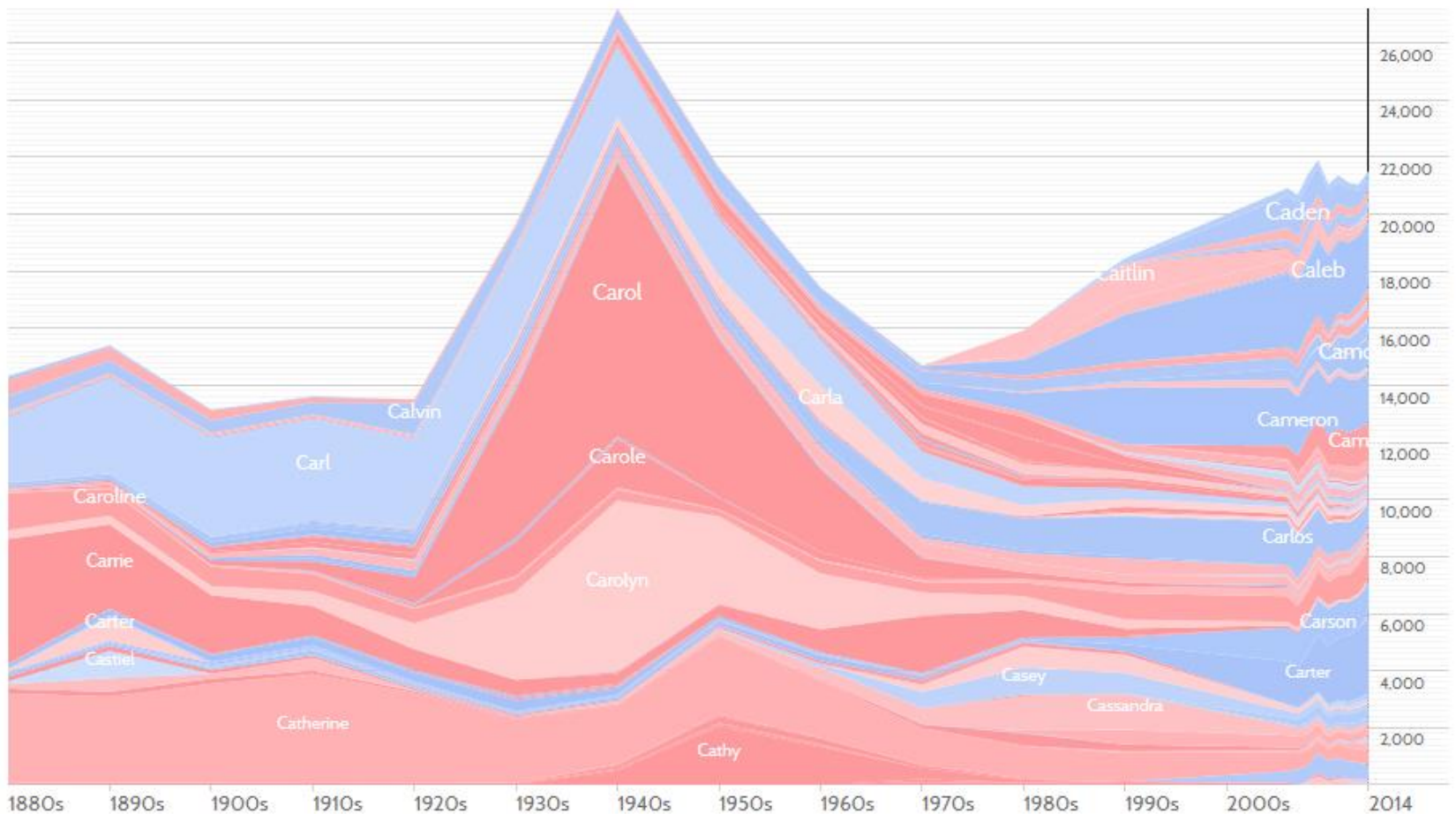
Both  Boys  Girls

boys	1000	500	100	25	1
girls	1000	500	100	25	1

Current rank:

per million births

Names starting with 'CA' per million babies



# Stream Graph

Baby Name >   Both  Boys  Girls

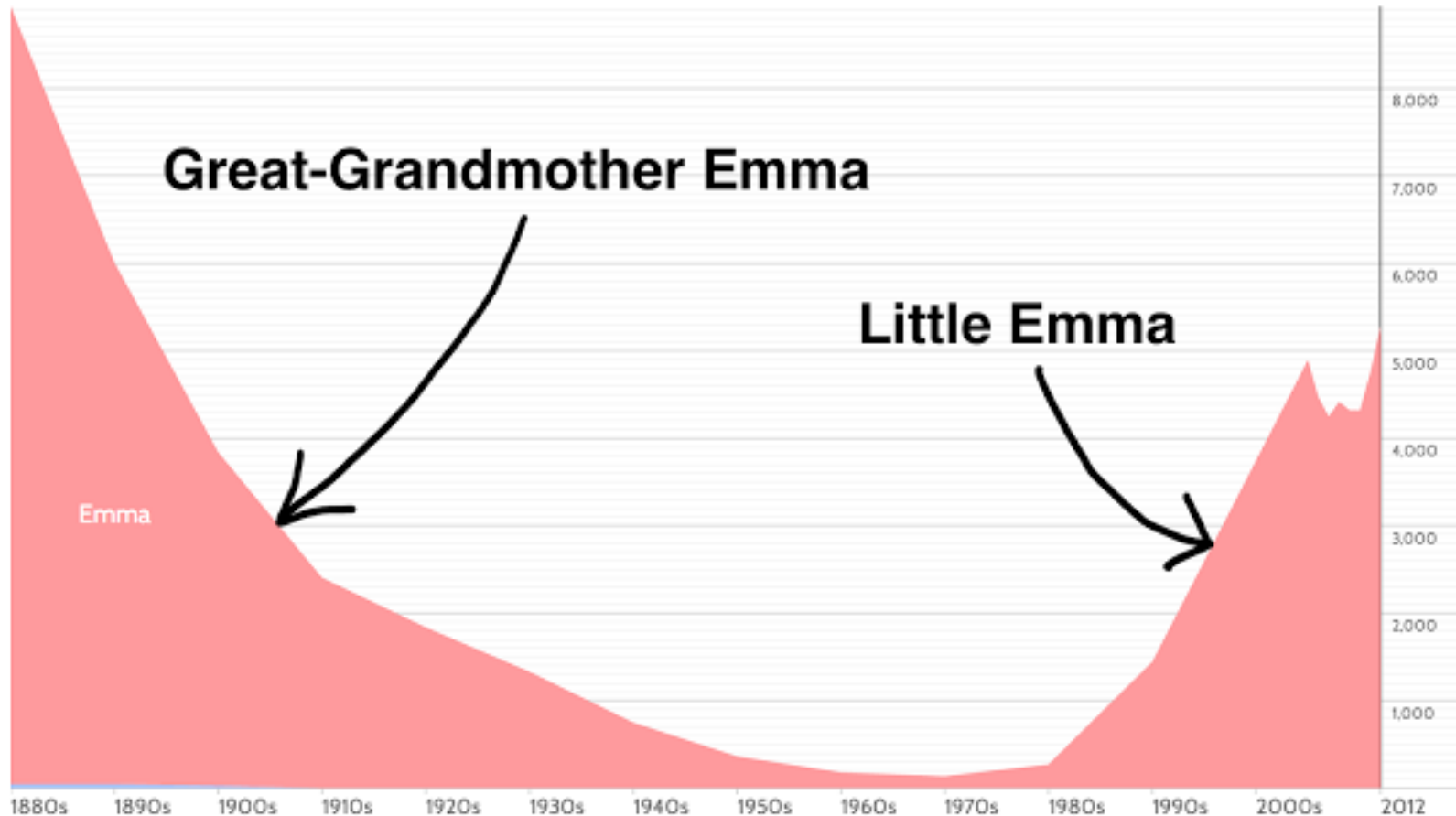
Current rank: boys 

1000	500	100	25	1
------	-----	-----	----	---

  
girls 

1000	500	100	25	1
------	-----	-----	----	---

Babies named 'EMMA' per million babies.



babynamewizard.com (via waitbutwhy.com)

# Summary

- 1) Typical data mining steps.
- 2) Representing data:
  - Object-feature representation.
  - Discrete vs. numerical features.
- 3) Preprocessing data:
  - IID Assumption
  - Data cleaning.
  - Feature transformations.
- 4) Exploring data:
  - Summary statistics.
  - Visualization.

Next week: let's start some machine learning...