



Getting to Know Your Data

Jesse Davis

Outline

- Computing simple statistics
- Simple visualizations

Motivation

Exploring data is vitally important as it helps you think about

- What technique to use
- What challenges you may (will) encounter
- What flaws or problems there are with data
- What results to expect

What Should You Look For?

- Good to look at simple statistics of
 - Number of variables
 - Size of data
 - Missing values
 - Class skew
- For each attribute, look at
 - **Discrete:** number of possible values, are they ordered, frequency of each value, etc.
 - **Numeric:** range, mean, min, max, etc.

Summary Statistics

- **Mean:** $\mu_X = \frac{\sum_i x_i}{n}$
- **Mode:** Most common value
- **Median:** Value v s.t. half the values above v , half below v
- **Variance:** $\sigma^2 = \frac{\sum_i (x_i - \mu_X)^2}{n}$
- **Quartile:** Sort X
 - Q1: Value at position $0.25n$
 - Q3: Value at position $0.75n$
 - IQR: $Q3 - Q1$

Quiz

6	1	3	7	10	1	9	3
---	---	---	---	----	---	---	---

What is the:

- Mean:
- Mode:
- Median:
- Q1:
- Q3:

Quiz

6	1	3	7	10	1	9	3
1	1	3	3	6	7	9	10

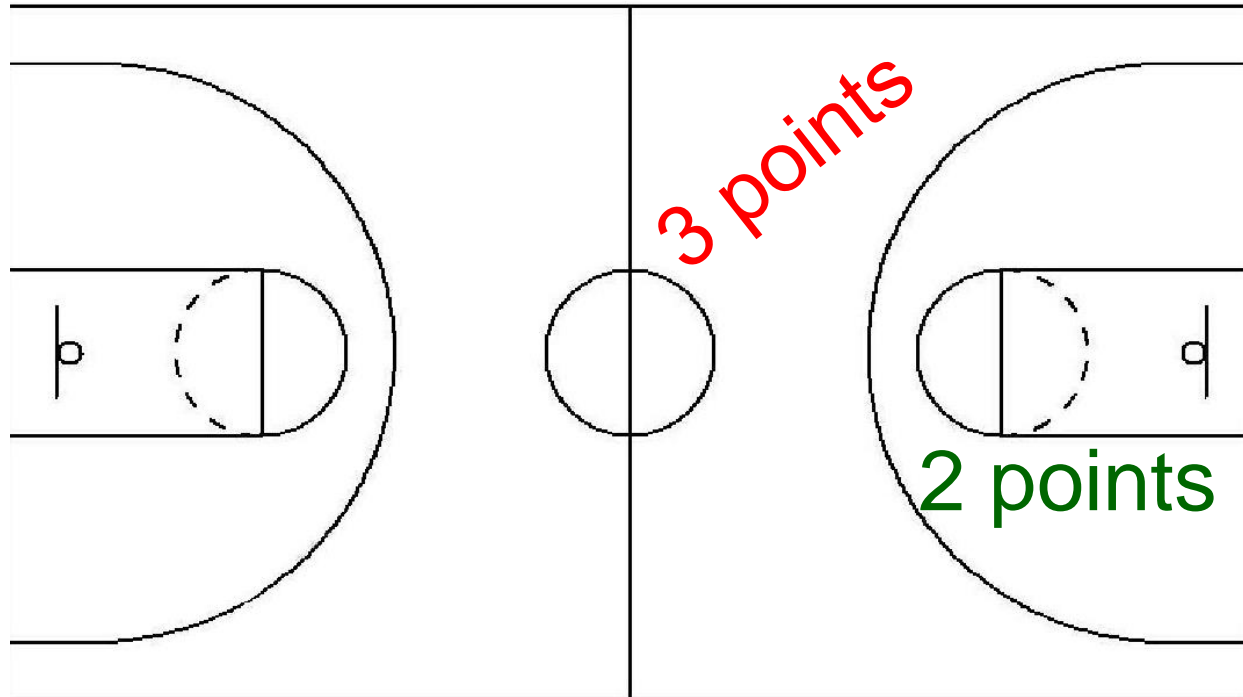
What is the:

- Mean: 5
- Mode: 1, 3
- Median: 4.5
- Q1: 2
- Q3: 8

Weighted Averages

- Mean assumes that each data point is of equal importance
- Often some data points are more important
 - Estimate more reliable
 - More valuable
 - Etc.
- Weighted average: $\mu_x = \frac{\sum_i w_i x_i}{\sum_i w_i}$

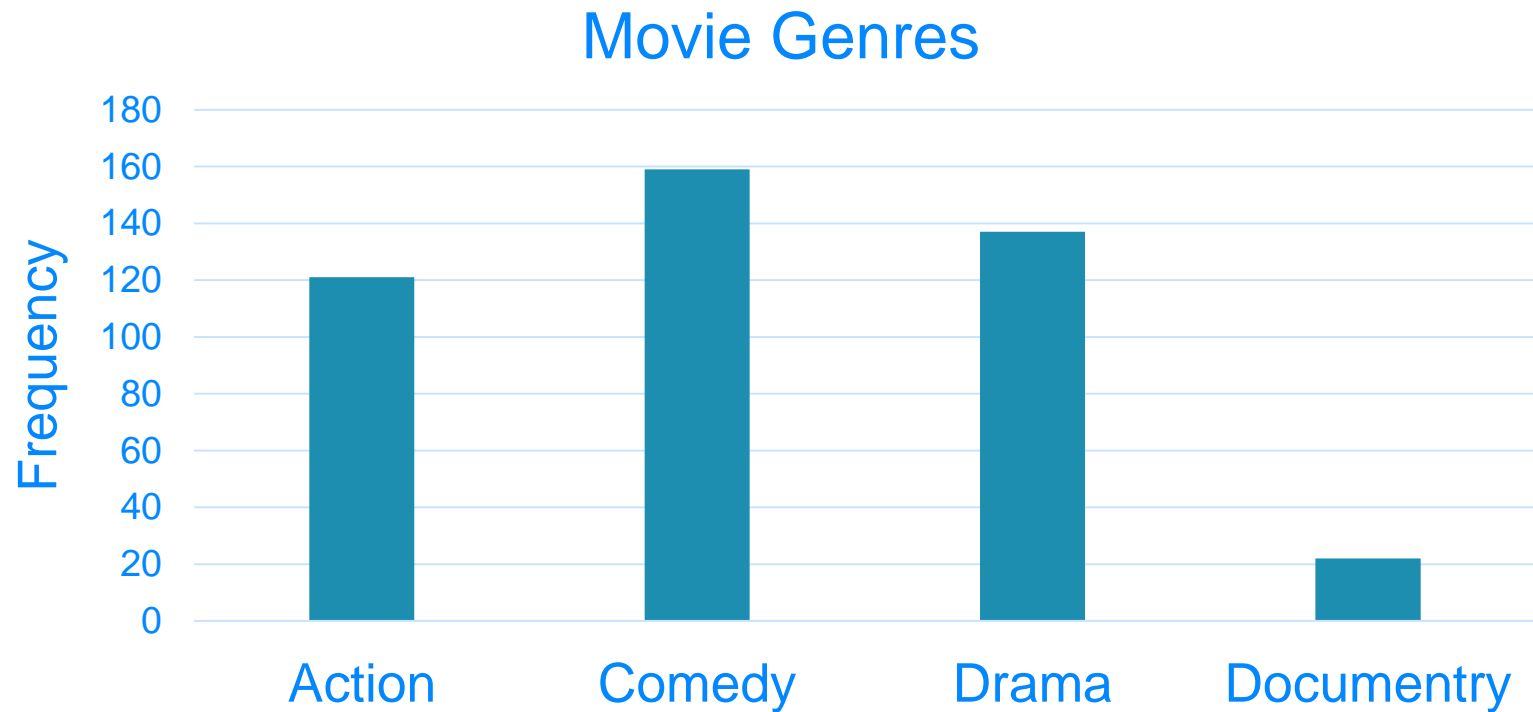
Example: Basketball



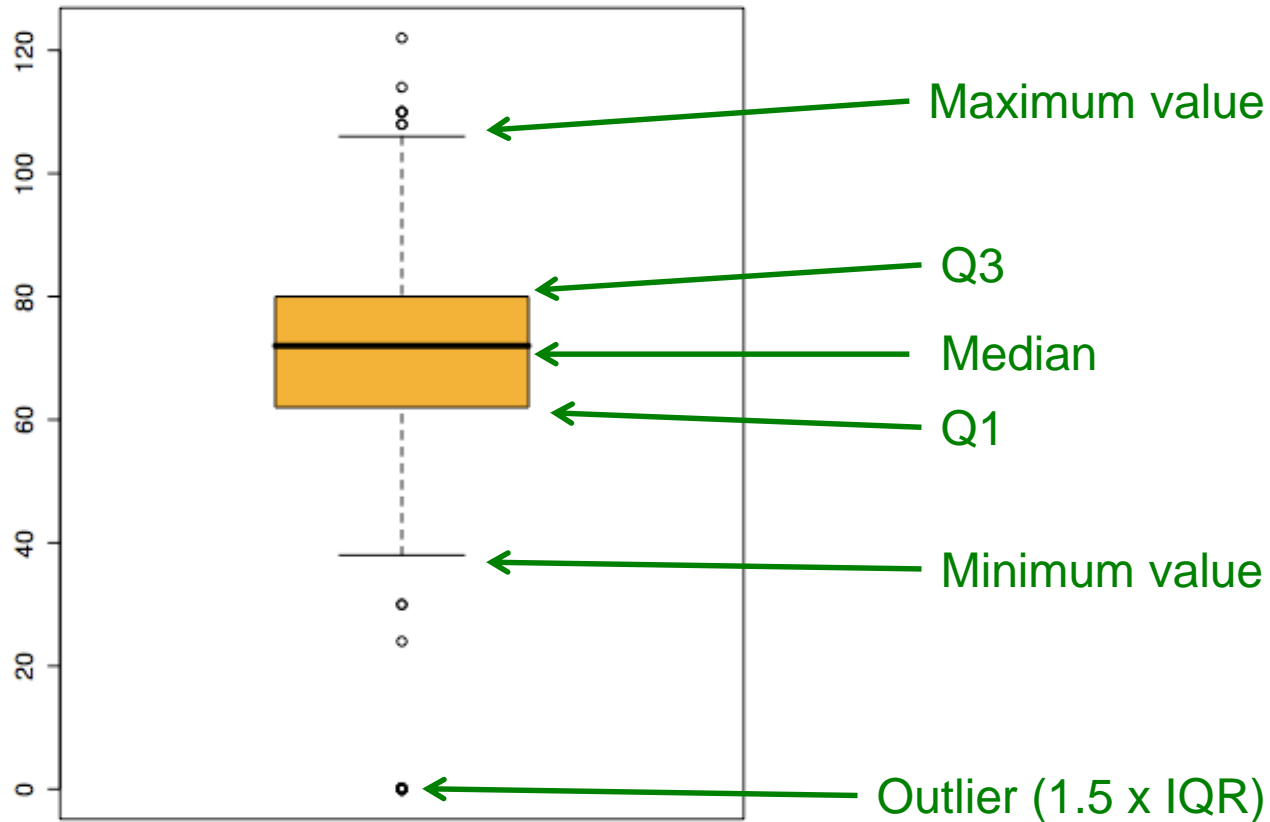
$$\text{FG \%} = \frac{2p + 3p}{2pa + 3pa}$$

$$\text{eFG \%} = \frac{2p + 1.5 * 3p}{2pa + 3pa}$$

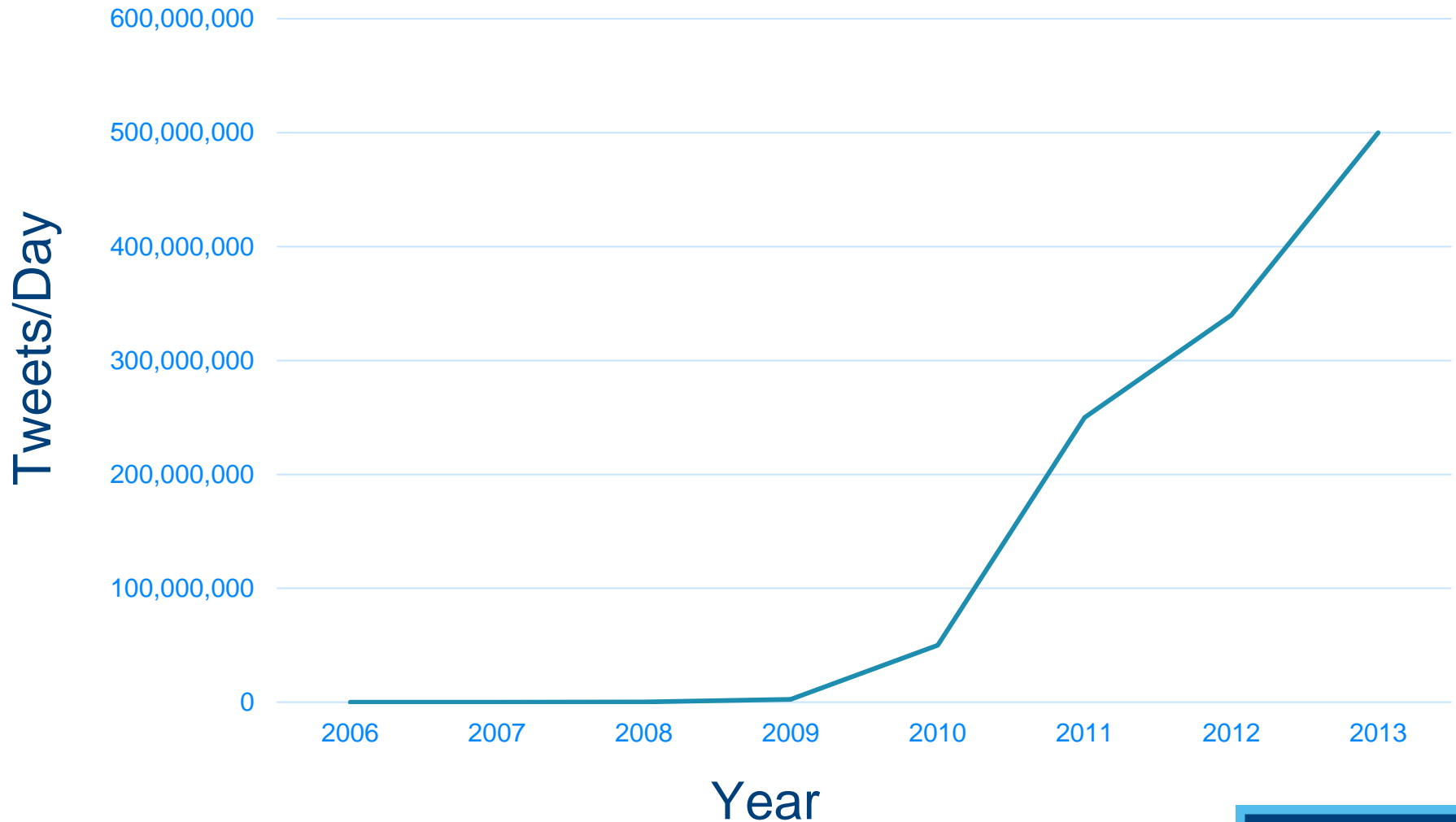
Histograms



Boxplots



Time Series

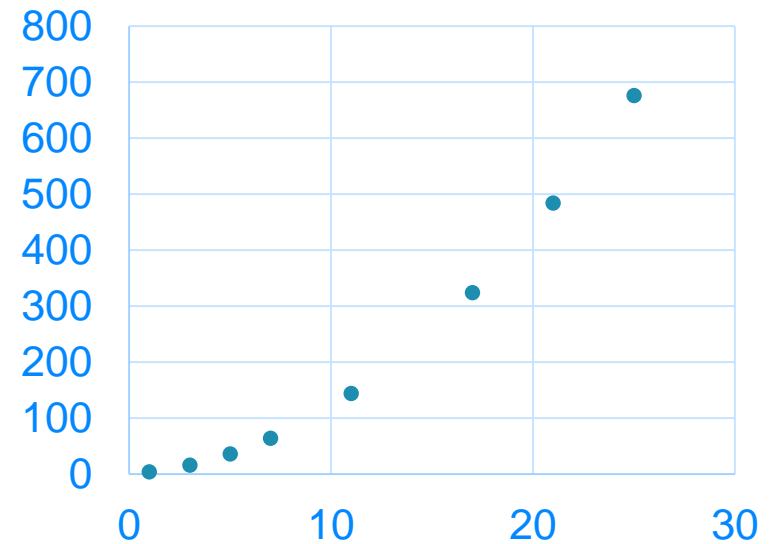
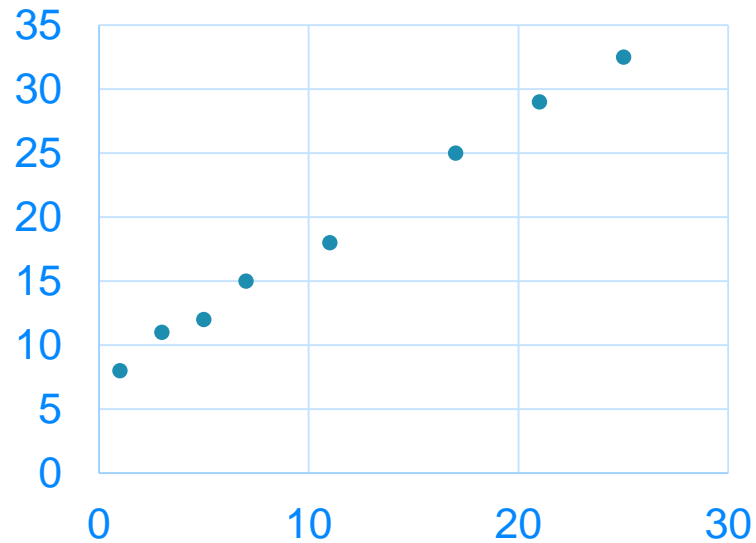


Time Series

Inglourious Basterds IMDB Rating Over Time



Scatterplots



Spatial

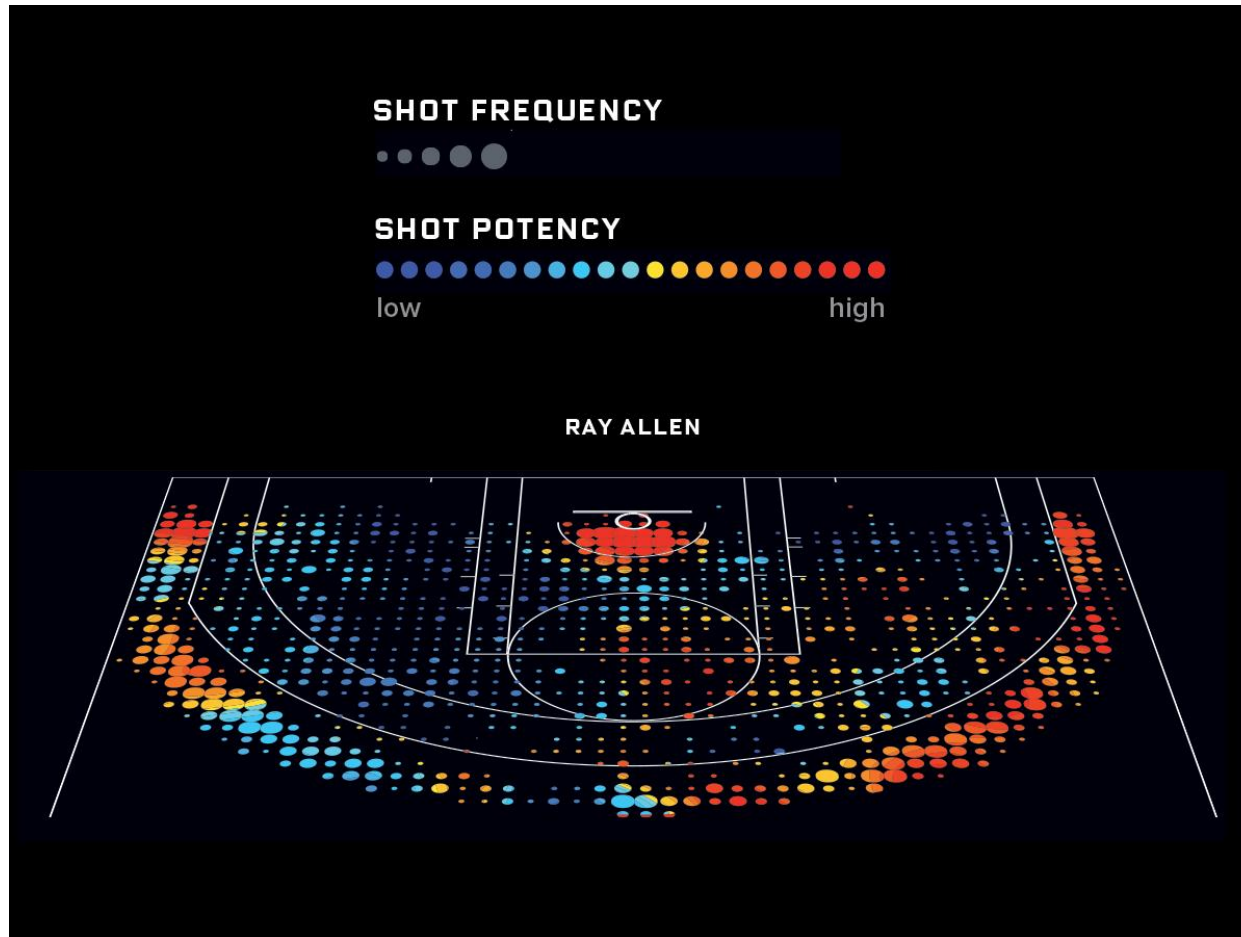


Image from: <http://www.wired.com/2014/10/faster-higher-stronger/>