

Special Topics in Biostatistics and Bioinformatics

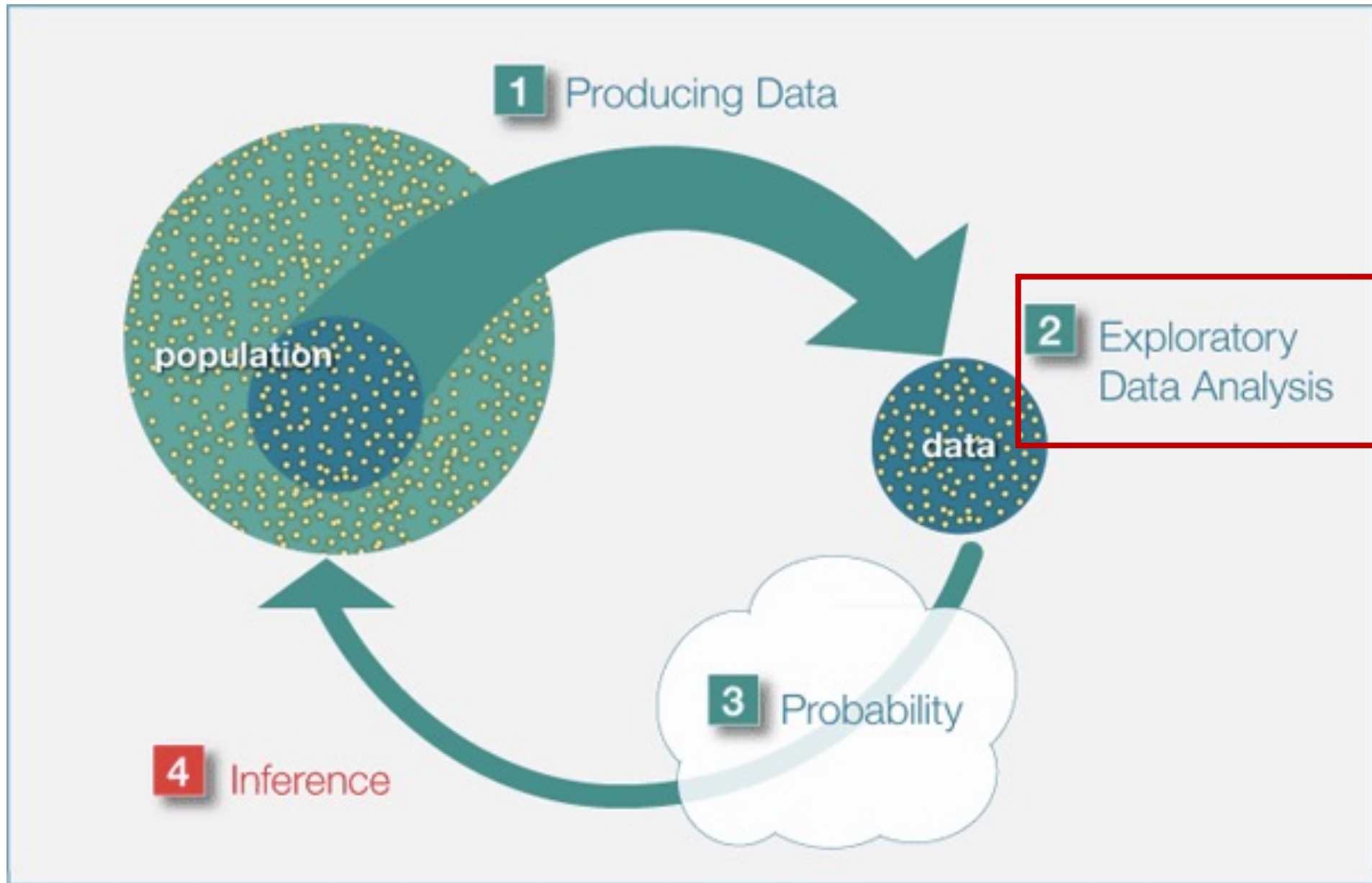
Week II

Ege Ülgen, M.D.

10 March 2022



ACIBADEM
MEHMET ALİ AYDINLAR
ÜNİVERSİTESİ



Reasons to explore data

- To understand data properties
- To find patterns in data
- To suggest modeling strategies
- To debug analyses
- To communicate results

Exploratory Data Analysis

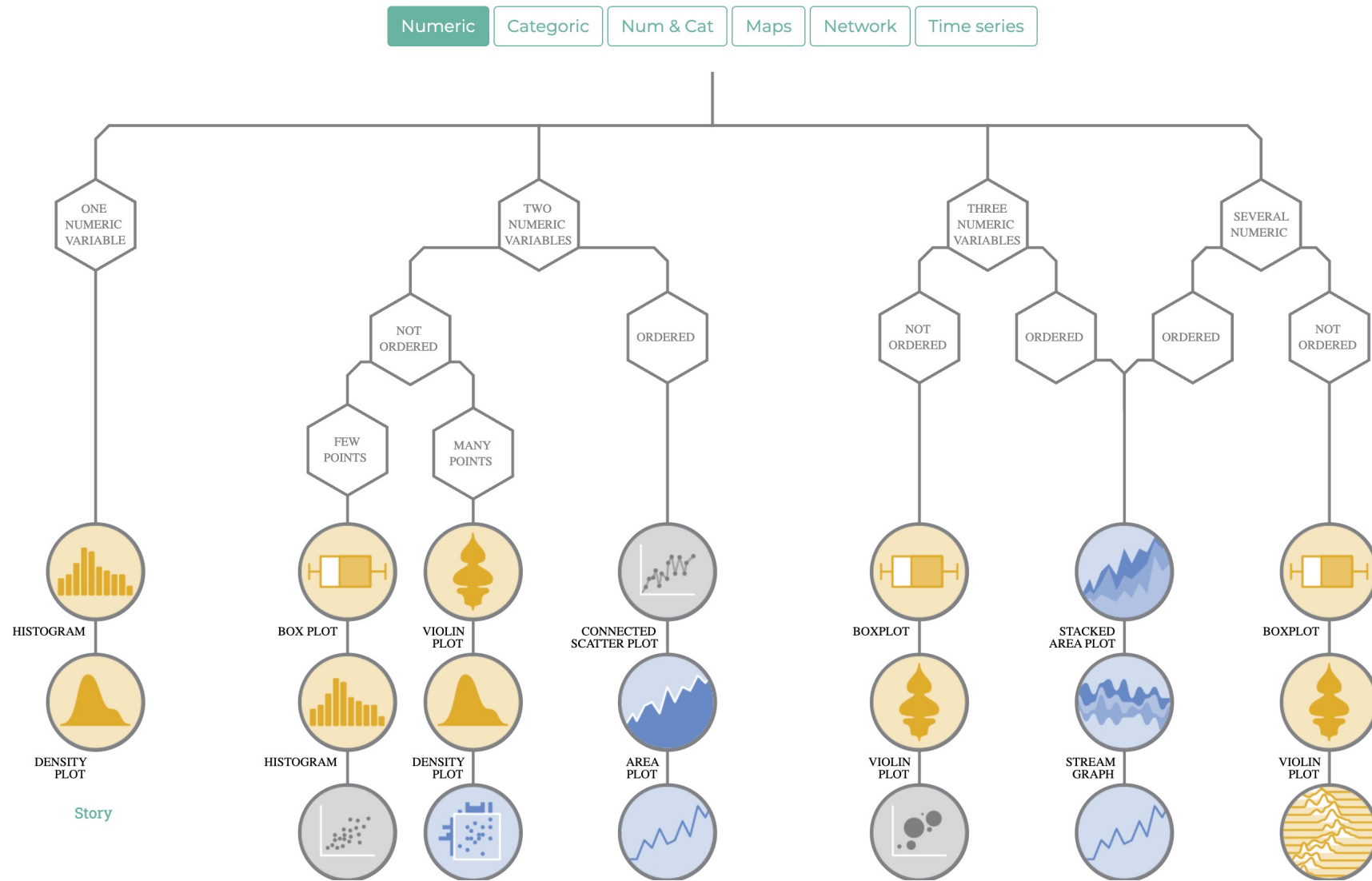
- Visualization
- Summarization
- Showing the data

**without being misled*

Examples of Visualization

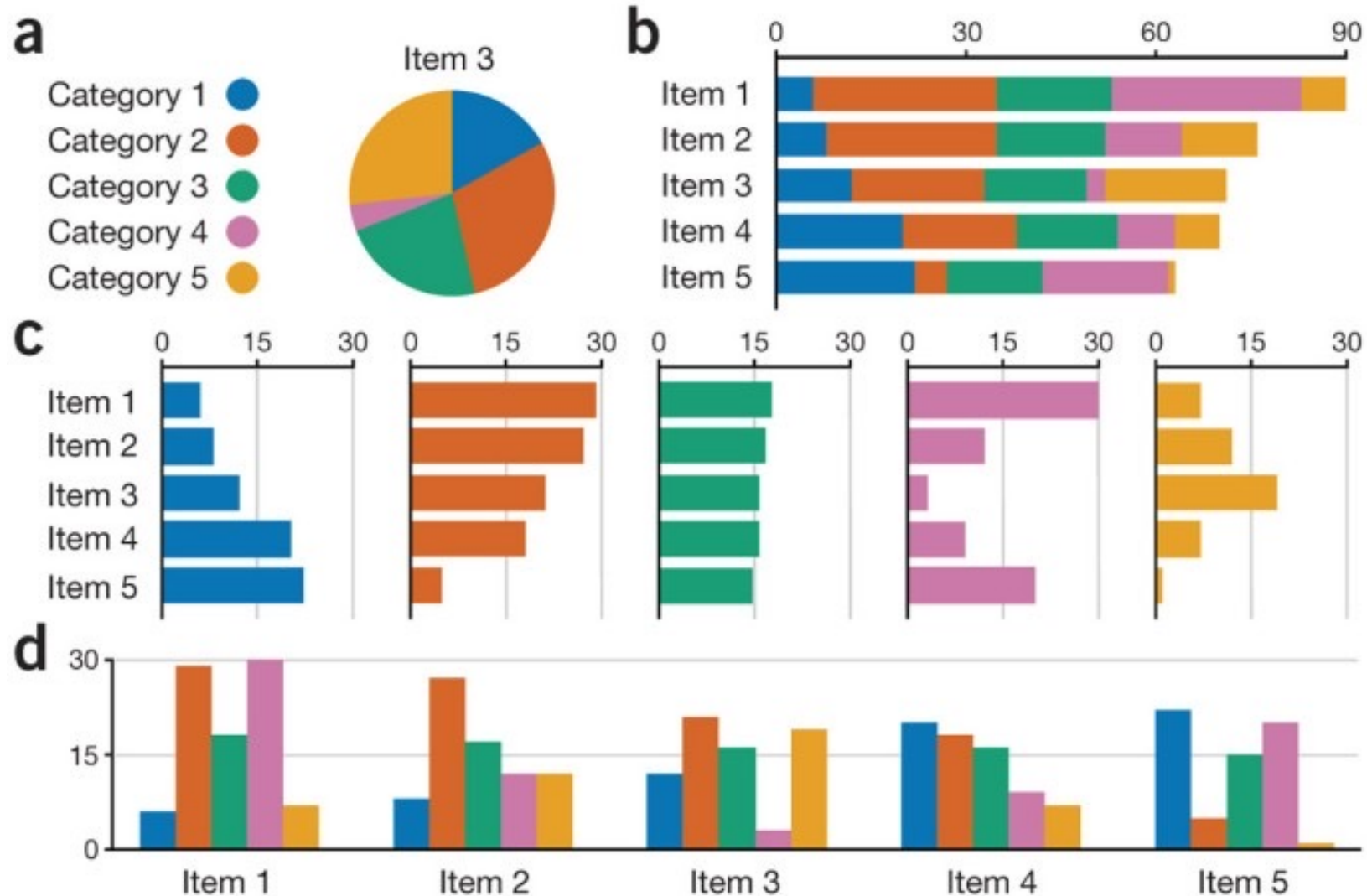


from Data to Viz



Story

Bar Charts



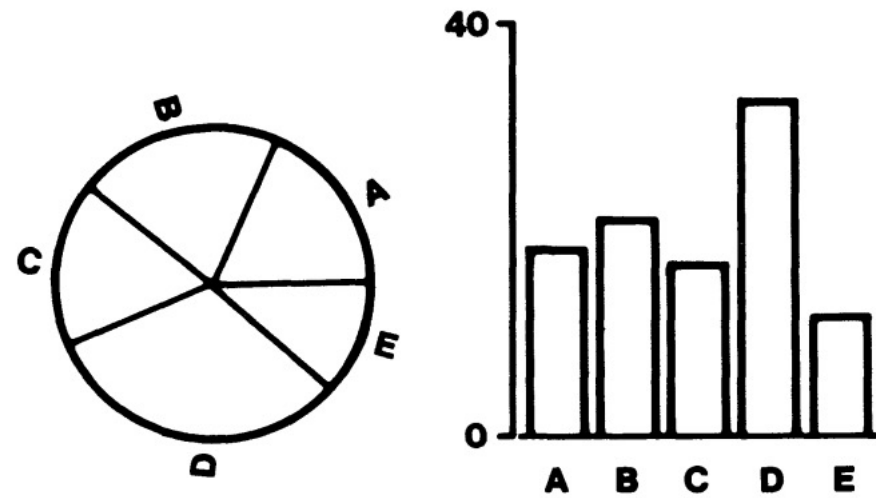


Figure 3. Graphs from position-angle experiment.

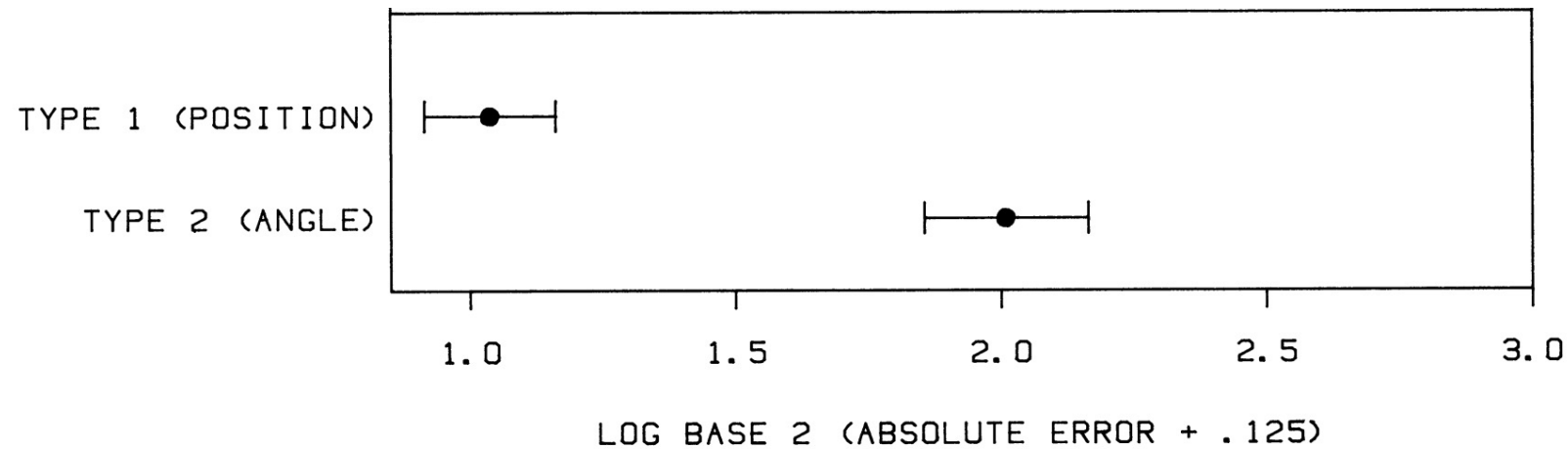
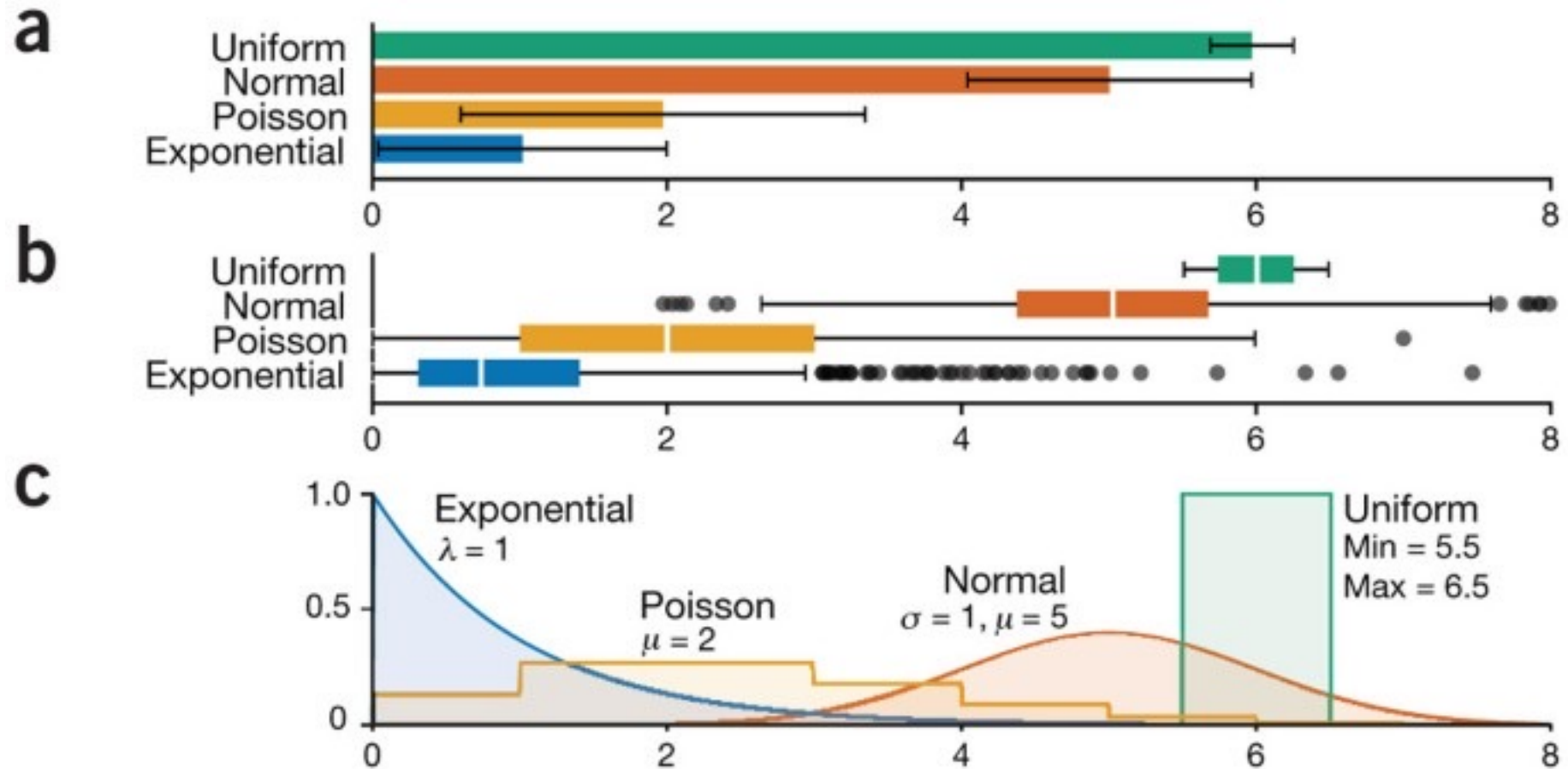


Figure 16. Log absolute error means and 95% confidence intervals for judgment types in position-length experiment (top) and position-angle experiment (bottom).

Box plots



Violin Plots

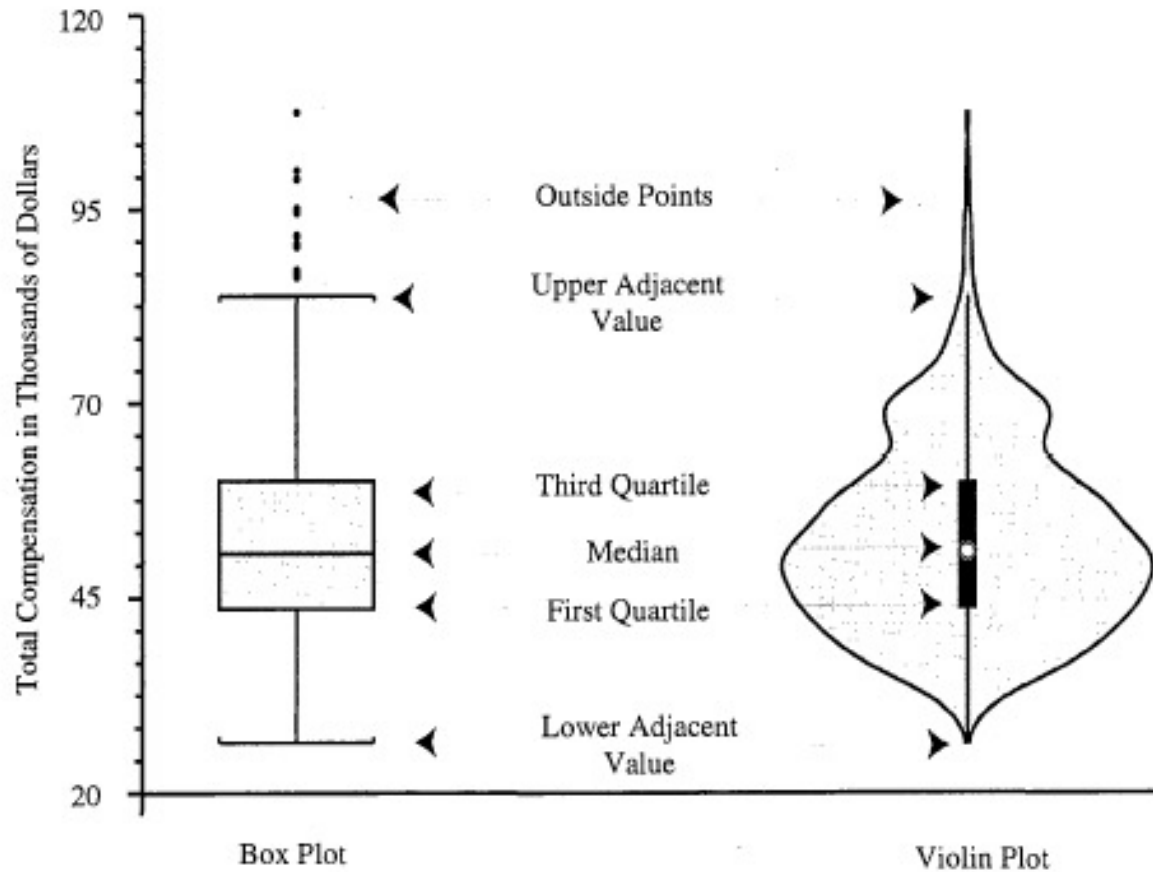


Figure 1. Common Components of Box Plot and Violin Plot. Total compensation for all academic ranks.

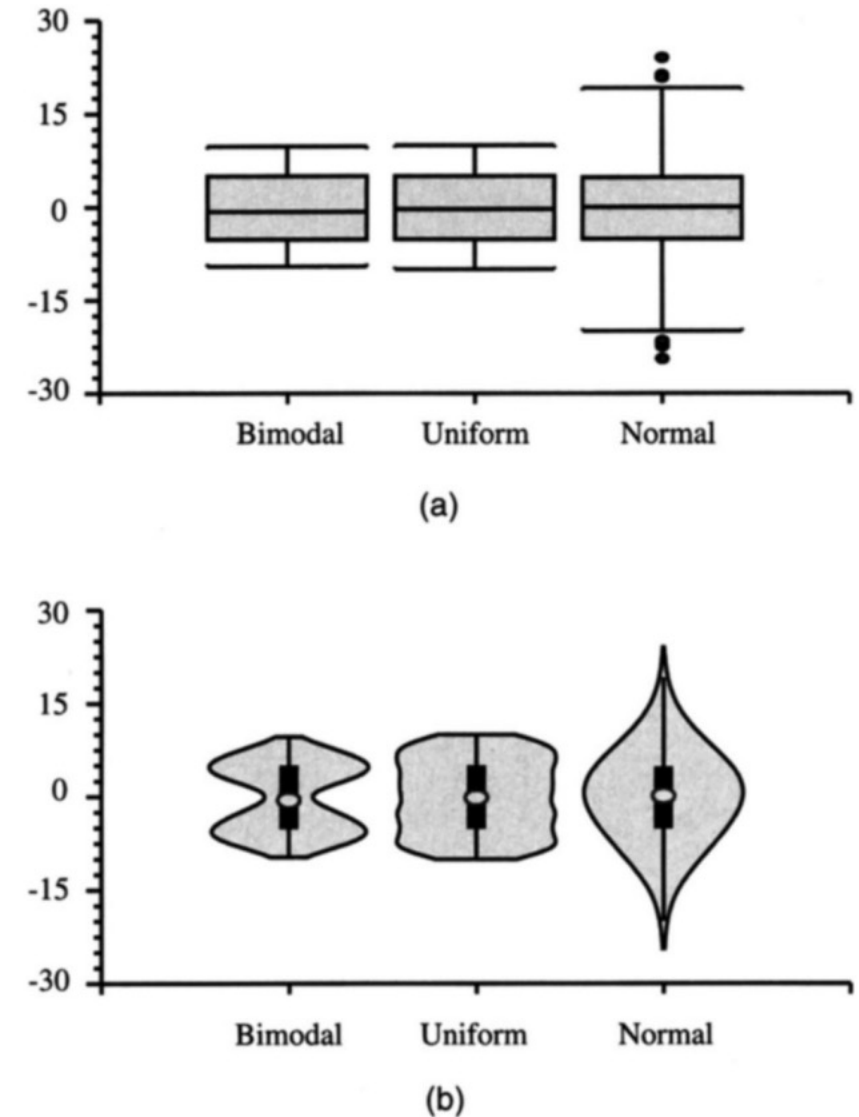
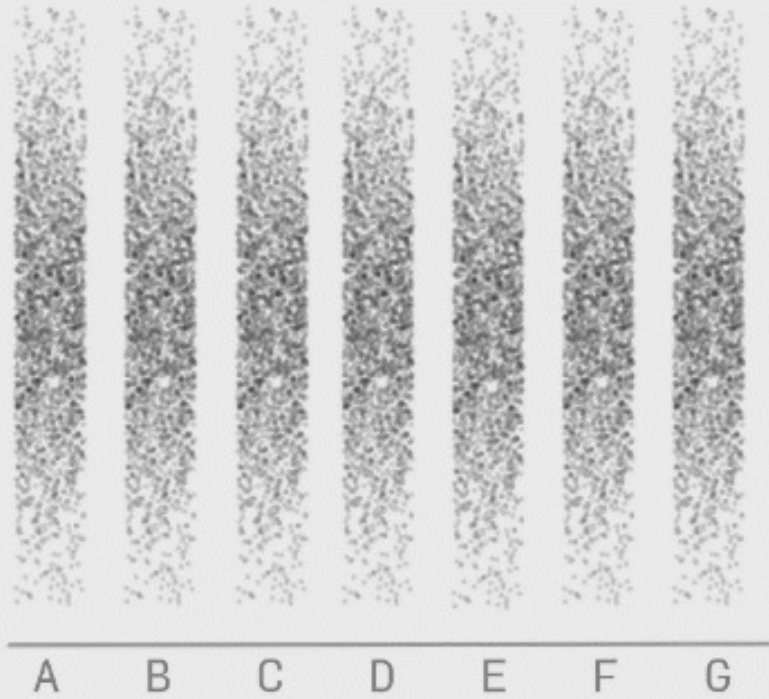
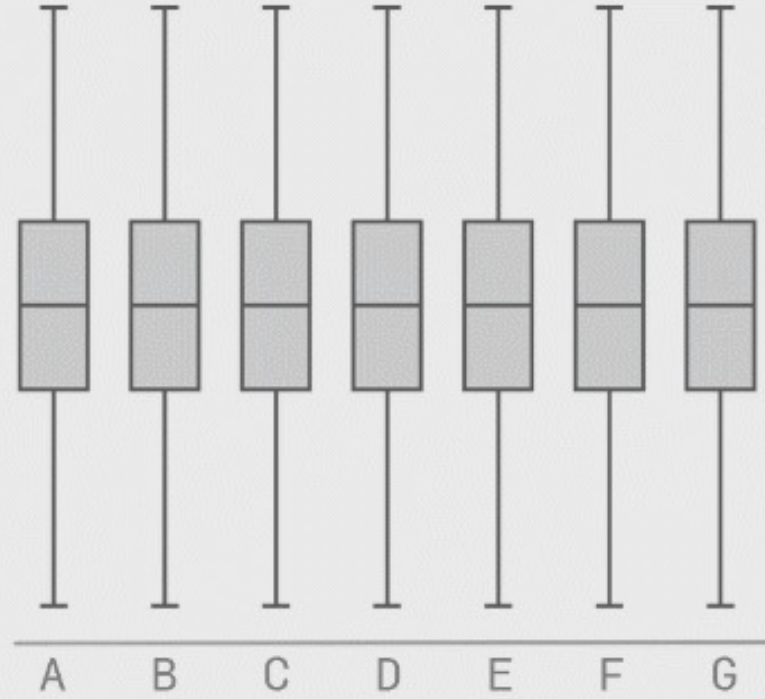


Figure 2. Comparison of Box Plots and Violin Plots to Known Distributions. (a) Box plots; (b) violin plots.

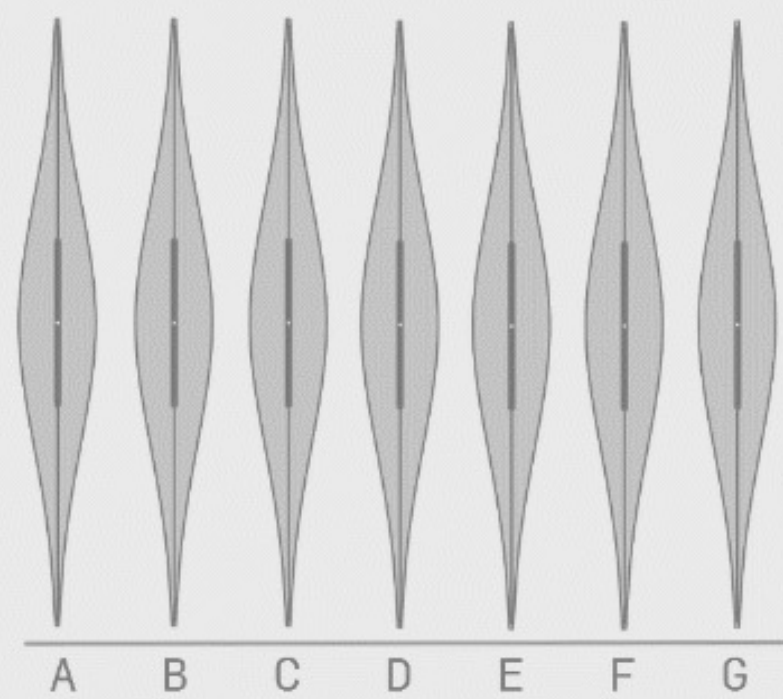
Raw Data



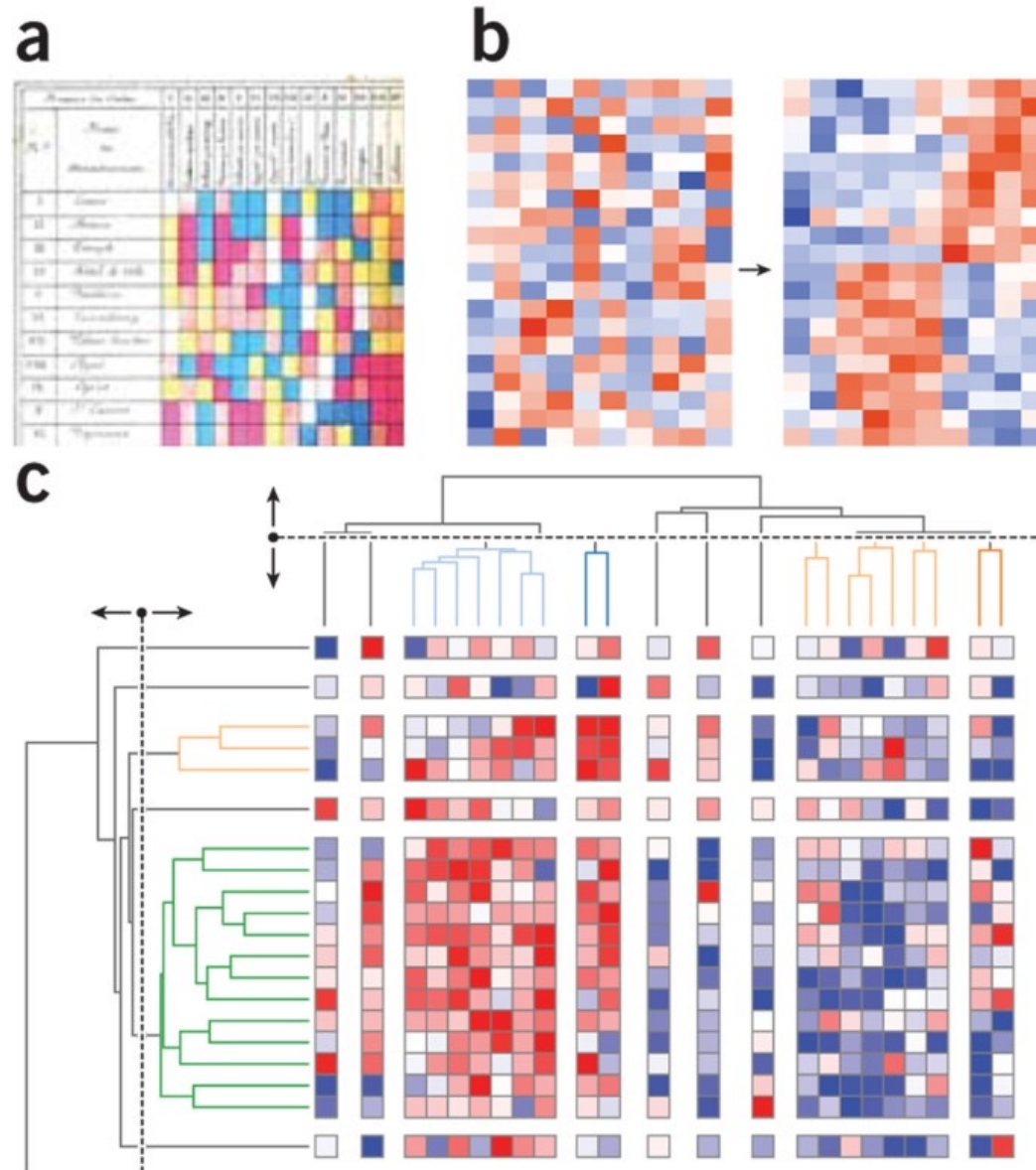
Box-plot of the Data



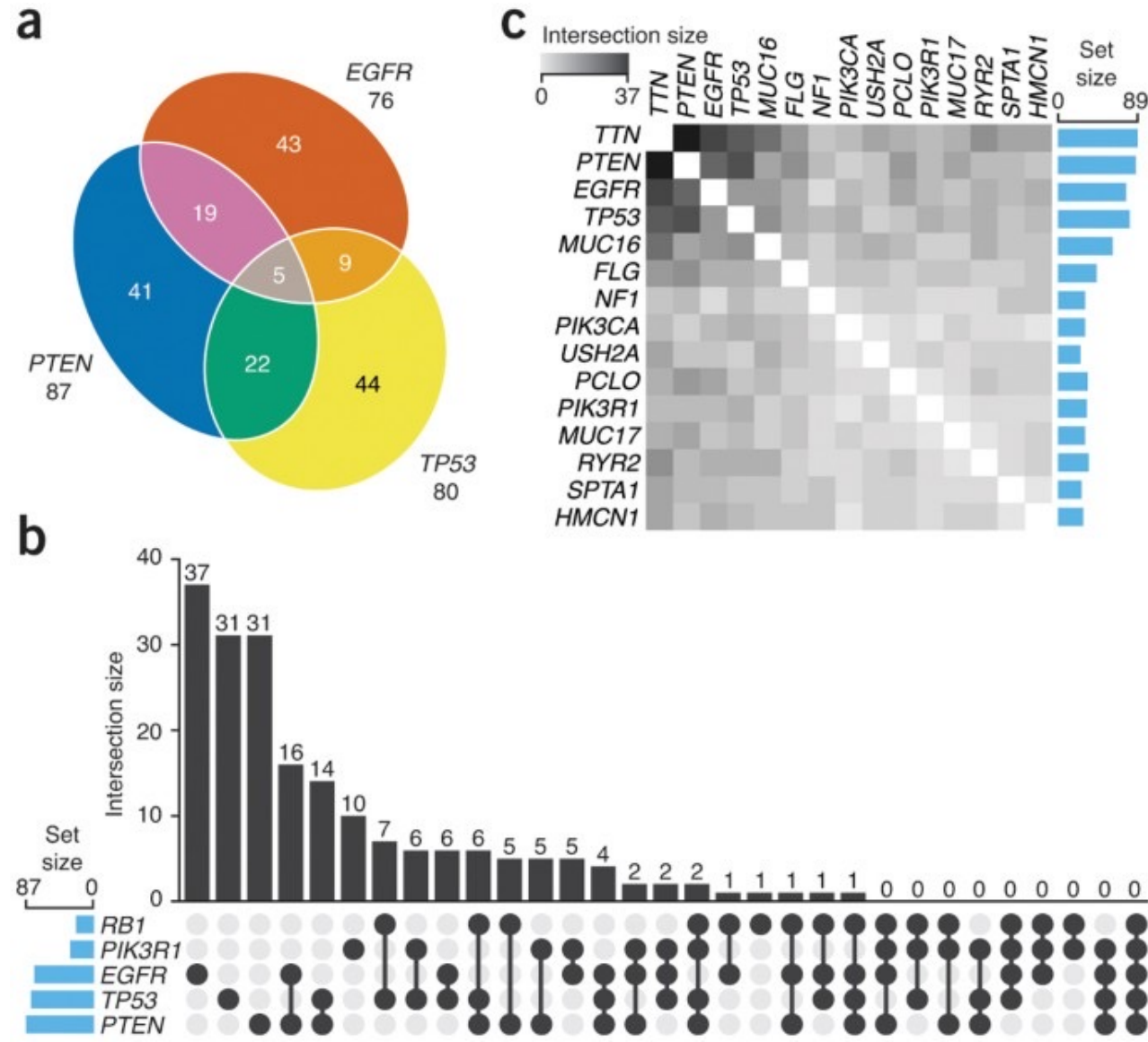
Violin-plot of the Data



Heat maps



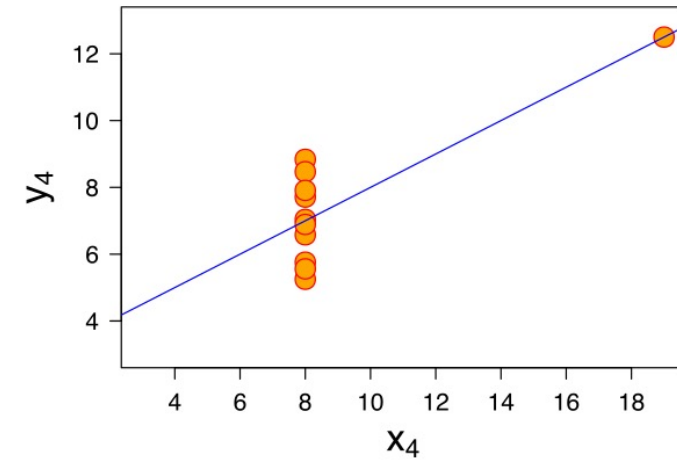
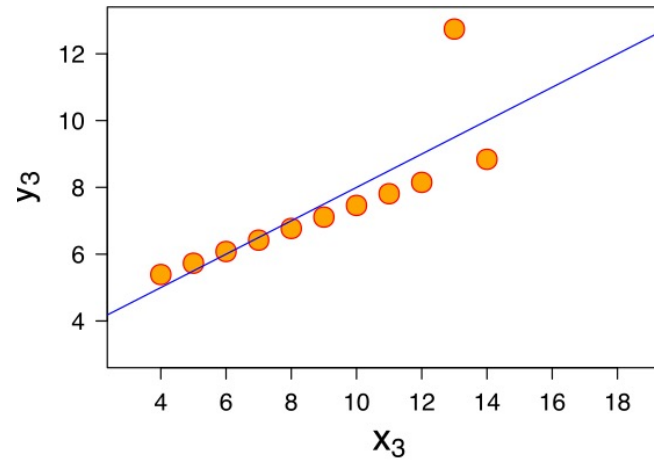
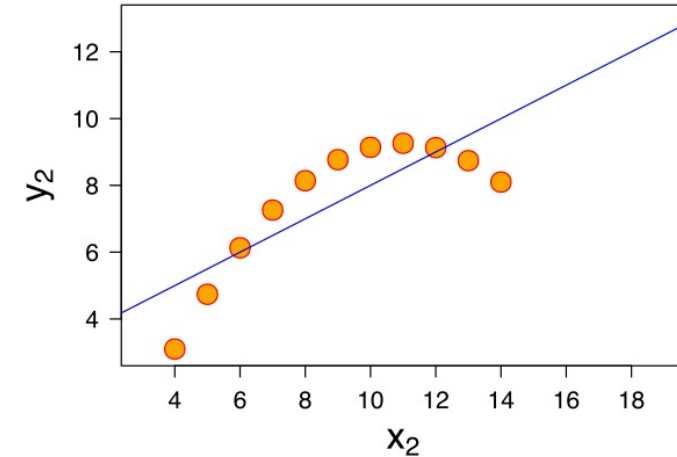
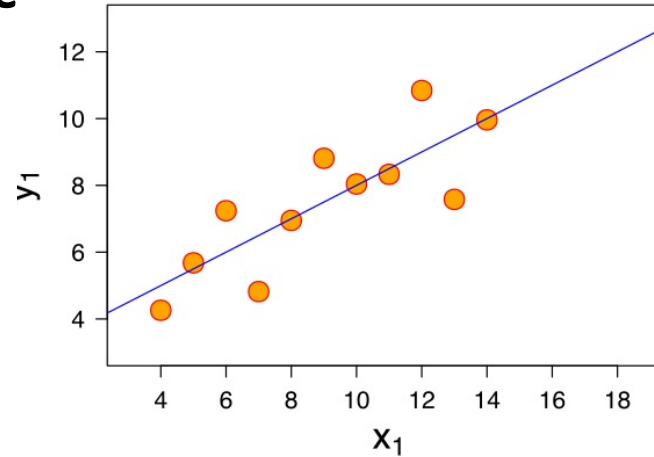
Sets and intersections

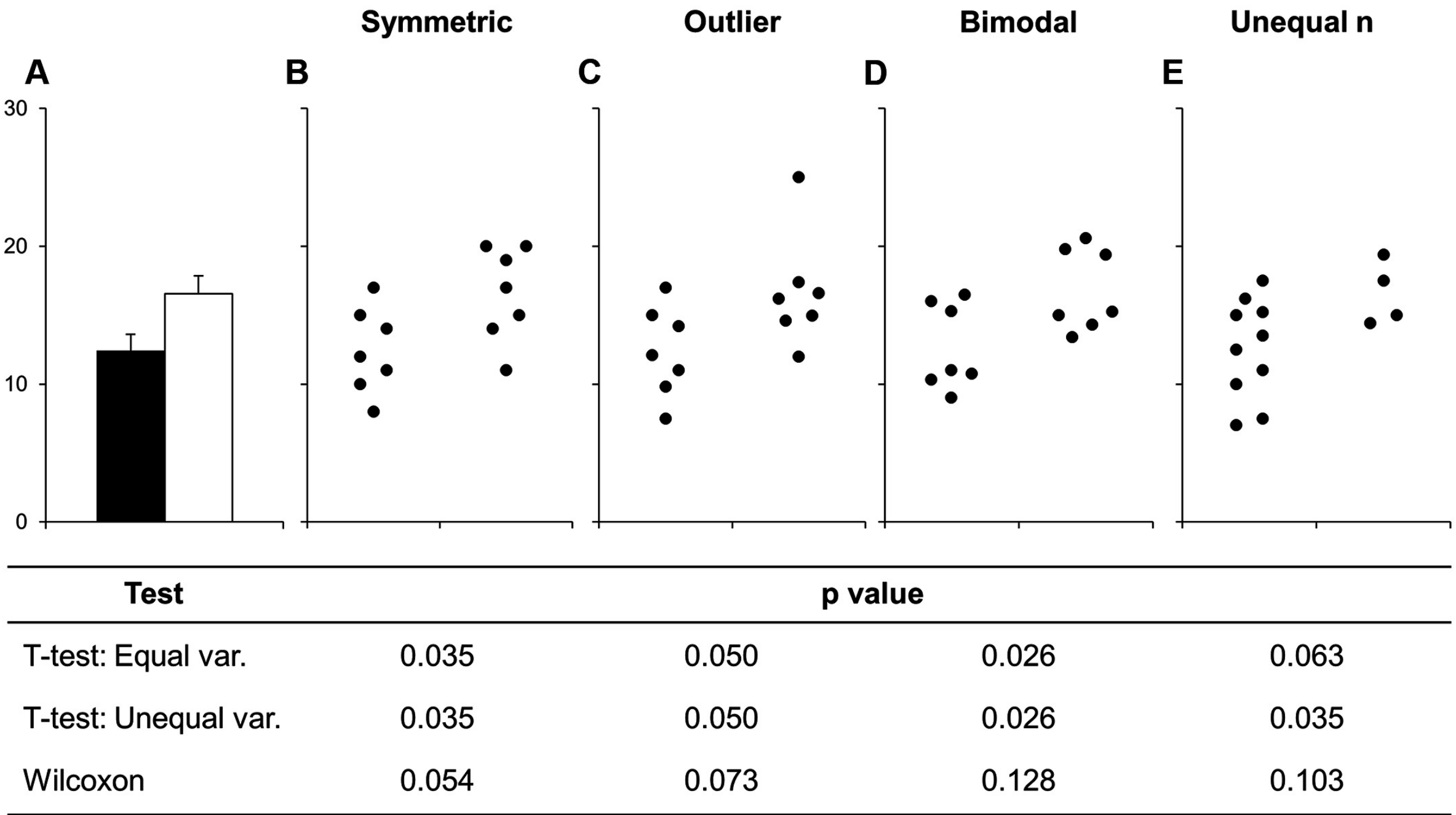


Avoid being misled

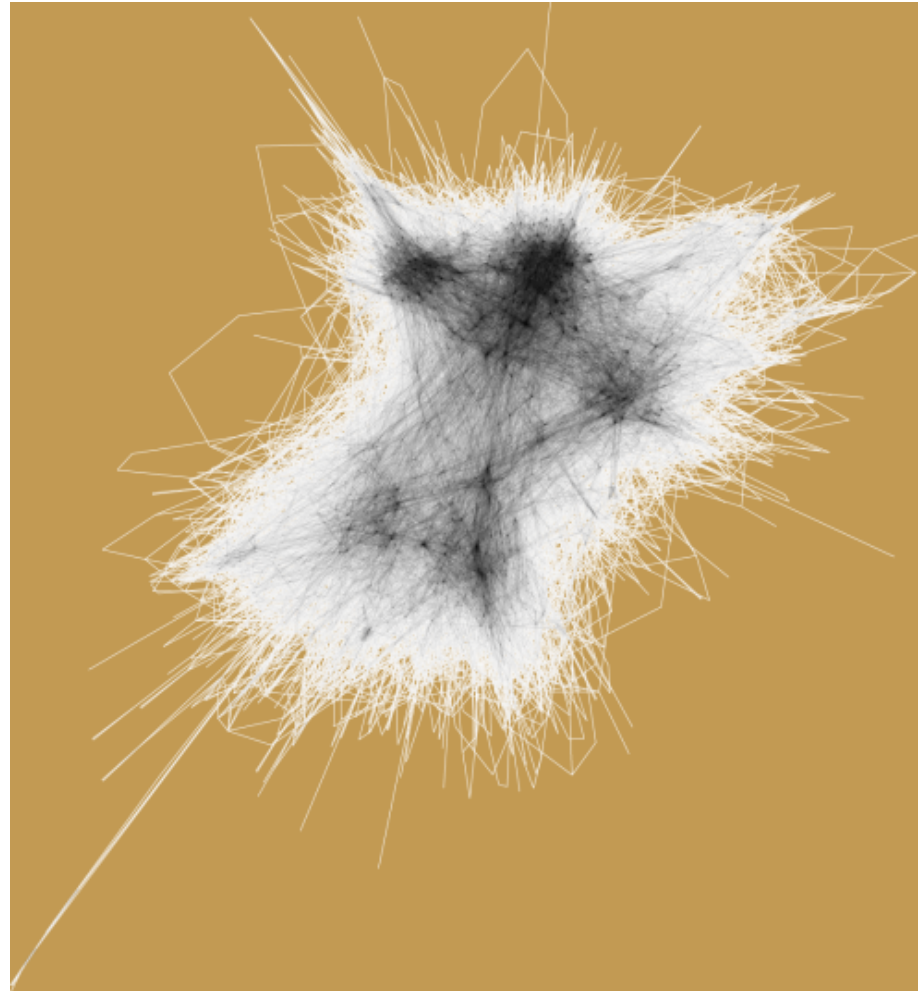
Anscombe's quartet

Property	Value
Mean of x	9
Sample variance of $x : s_x^2$	11
Mean of y	7.50
Sample variance of $y : s_y^2$	4.125
Correlation between x and y	0.816
Linear regression line	$y = 3.00 + 0.500x$
Coefficient of determination of the linear regression : R^2	0.67



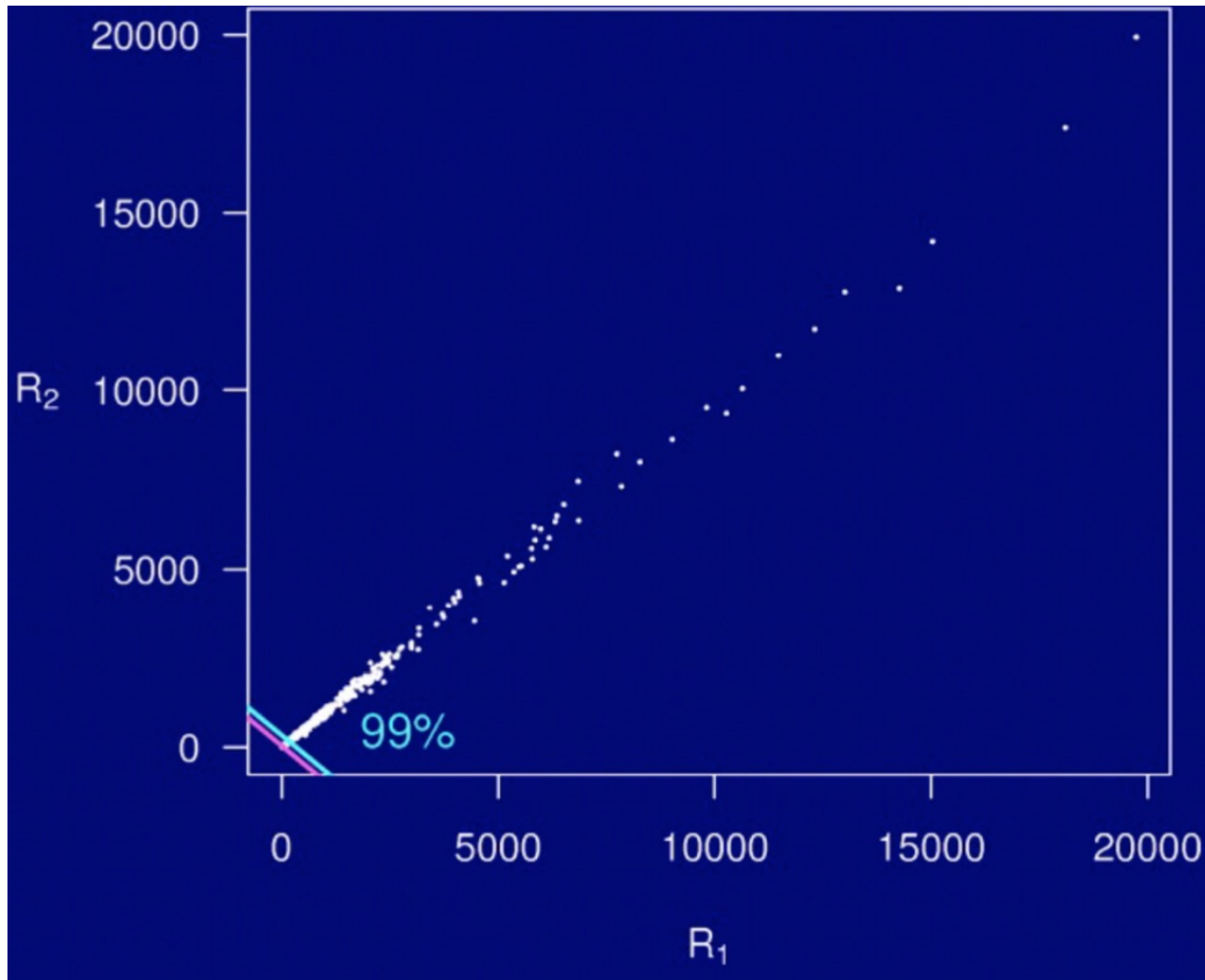


Avoid “Ridiculograms”



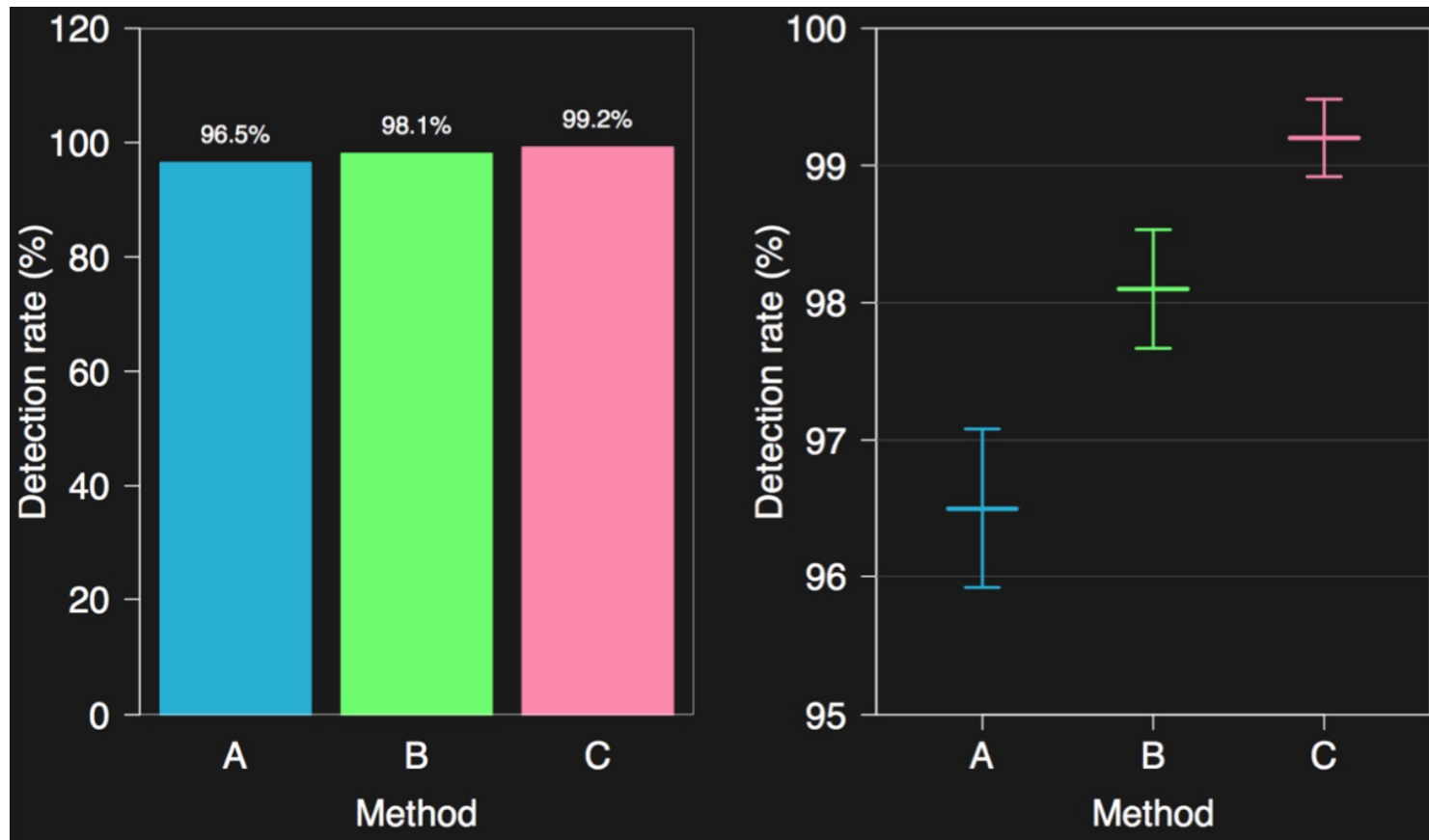
Scale - Basic Principles

- Be careful with scale



Scale - Basic Principles

- Use common scales and start at 0



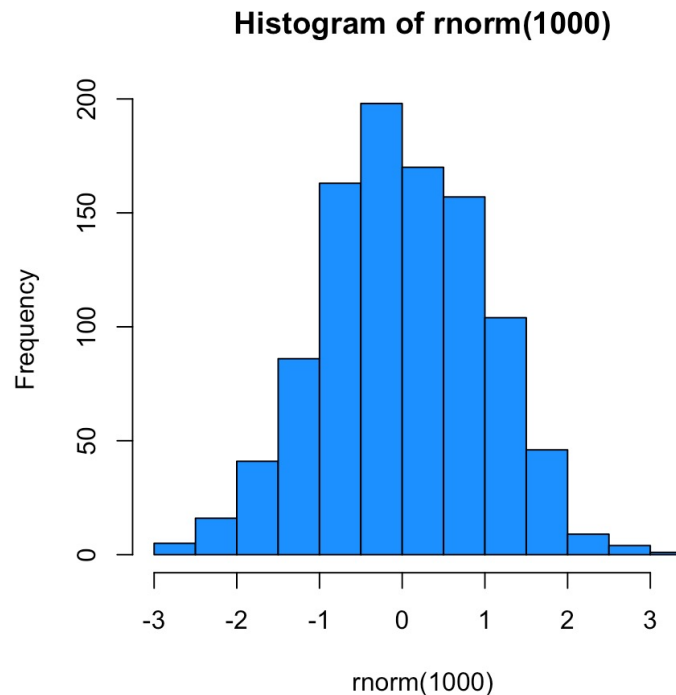
A nice Resource

- <http://blogs.nature.com/methagora/2013/07/data-visualization-points-of-view.html>

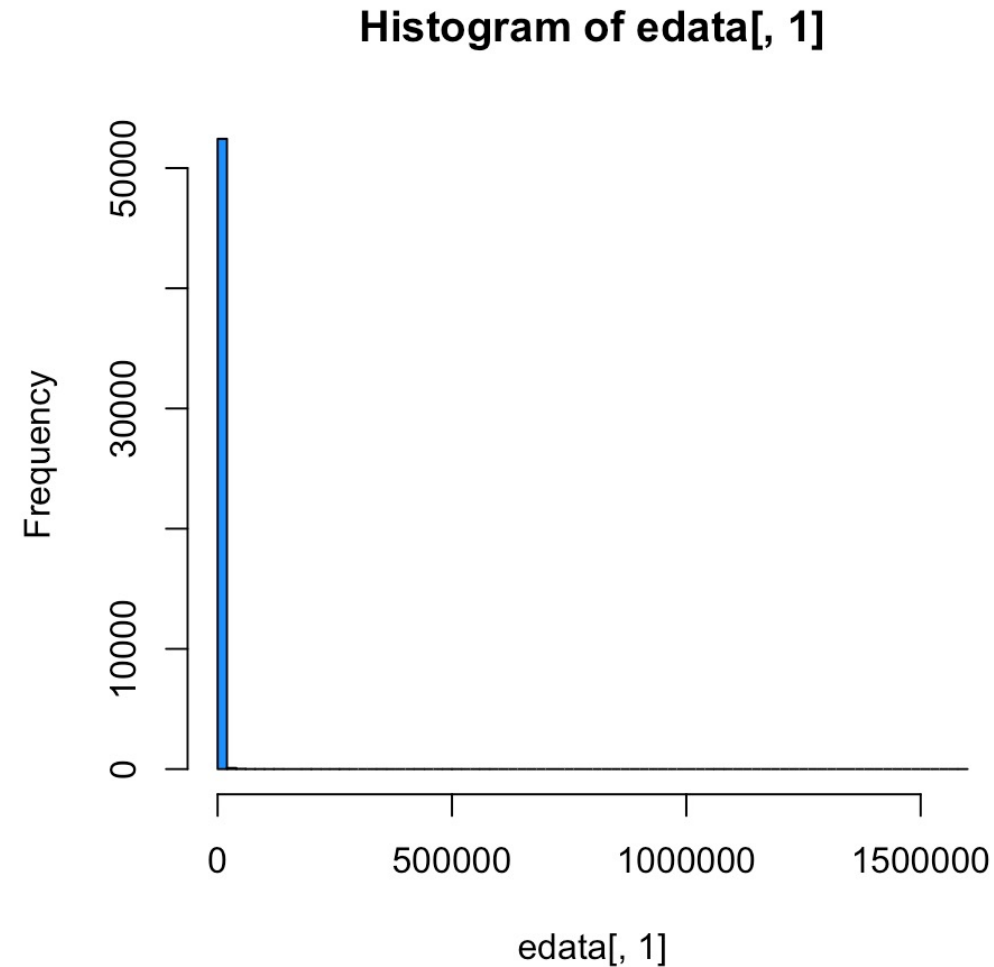
Data Transformations

Reasons for need for symmetric data

- plots are easier to see this way
- most statistical methods are designed to work better for non-skewed data

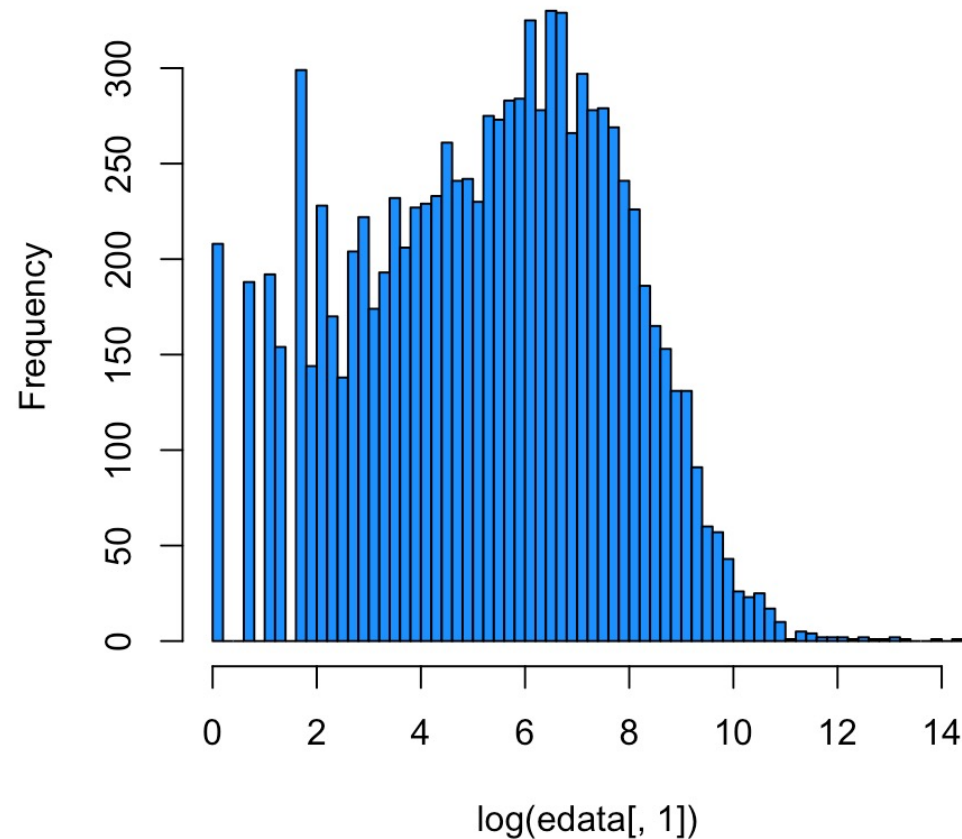


Realistically, most omics data is skewed

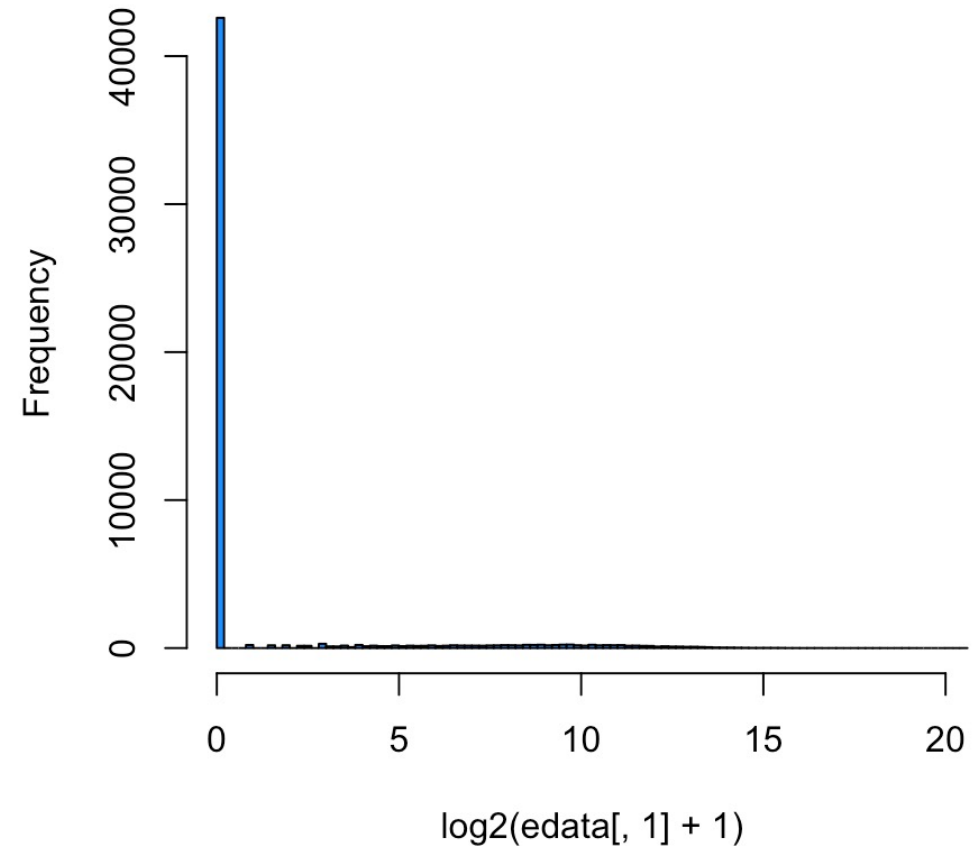


Realistically, most omics data is skewed

Histogram of $\log(\text{edata[, 1]})$

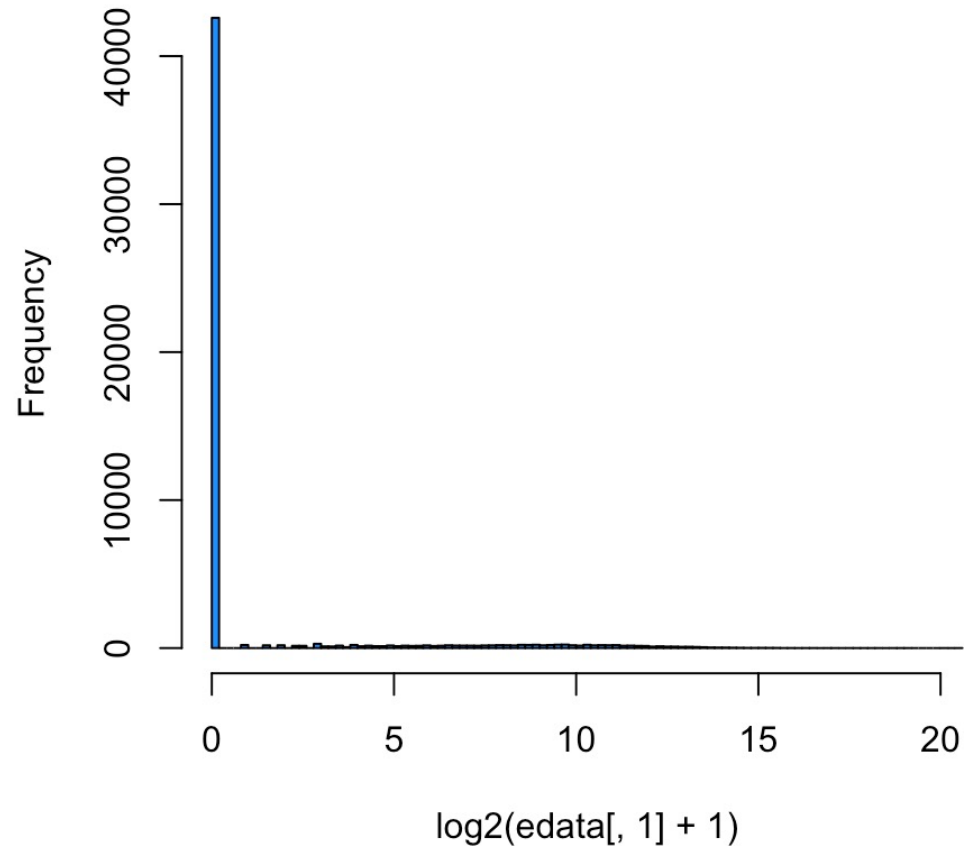


Histogram of $\log_2(\text{edata[, 1]} + 1)$

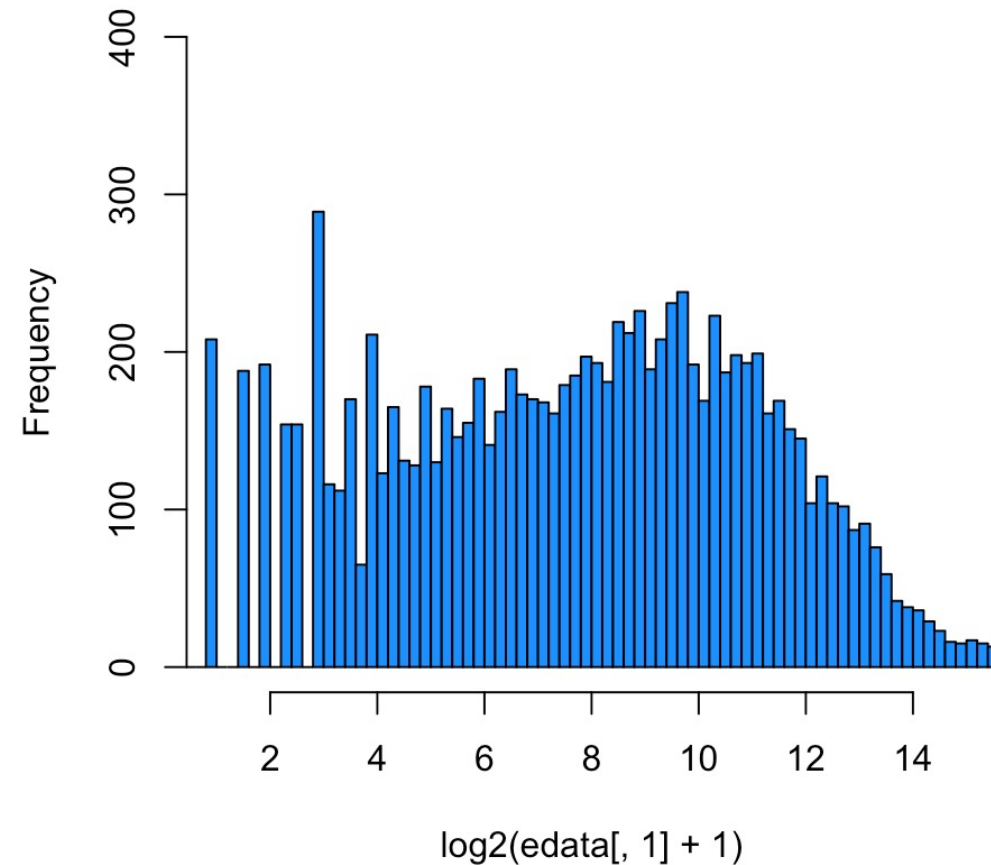


Realistically, most omics data is skewed

Histogram of $\log_2(\text{edata[, 1]} + 1)$



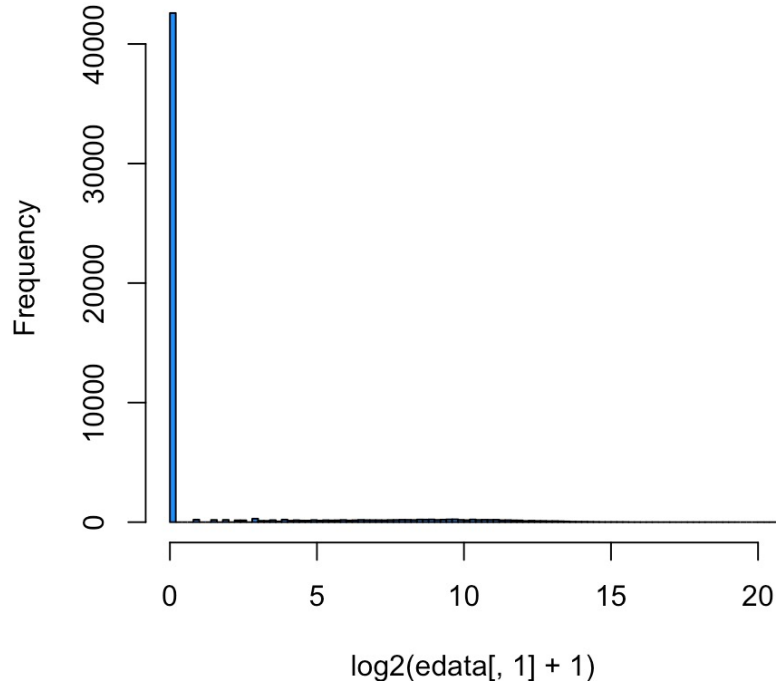
Histogram of $\log_2(\text{edata[, 1]} + 1)$



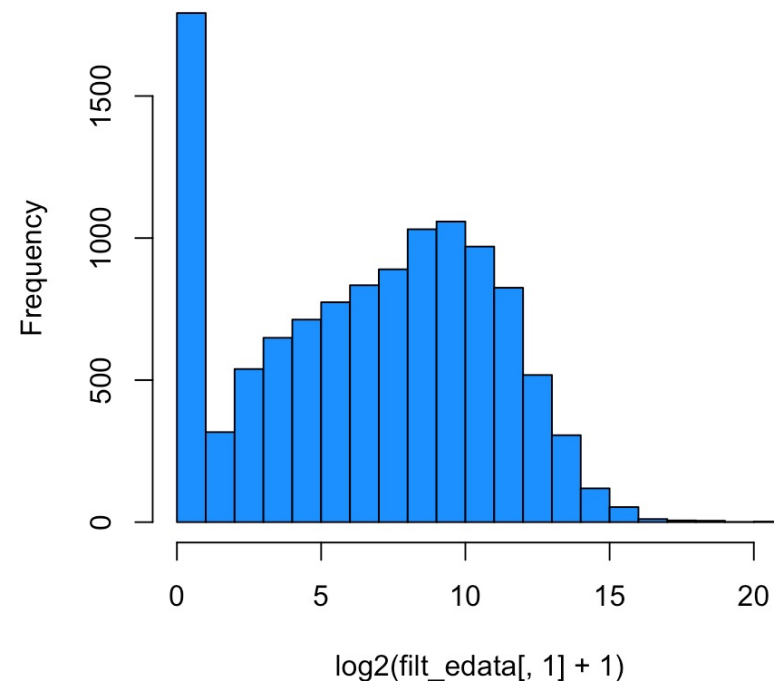
A common pre-processing technique is to remove features that don't have much data

```
low_genes = rowMeans(edata) < 5  
filt edata = filter(edata,!low_genes)
```

Histogram of $\log_2(\text{edata[, 1]} + 1)$



Histogram of $\log_2(\text{filt_edata[, 1]} + 1)$



Other transformations

- **Variance stabilizing transforms**

which seek to remove a mean variance relationship among the data

- **Box-Cox transforms**

which seek to make the data approximately Normally distributed

- **rlog transform**

unique to genomics count data, this is a regularized version of the log transform that seeks to minimize differences at low count levels

Preprocessing

- Convert raw data to “processed”, trying to remove technical artifacts
 - GC content bias
 - PCR duplicates
 - Probe sequence and fragment length
 - ...
- highly platform/problem dependent

Normalization

- Remove technological biases
- Make samples comparable
- highly platform/problem dependent

Scaling data - Standardization

$$X'_i = \frac{X_i - \textit{center}(X)}{\textit{scale}(X)}$$

Center: mean/median
Scale: sd/IQR/MAD

Scaling data – Min-Max Scaling

$$X'_i = \frac{X_i - \min(X)}{\max(X) - \min(X)}$$

Scaling in any interval $[a,b]$

$$x' = a + \frac{(x - \min(x))(b - a)}{\max(x) - \min(x)}$$

Quantile Normalization

Raw data	Order values within each sample (or column)	Average across rows and substitute value with average	Re-order averaged values in original order																																																																																
<table><tr><td>2</td><td>4</td><td>4</td><td>5</td></tr><tr><td>5</td><td>14</td><td>4</td><td>7</td></tr><tr><td>4</td><td>8</td><td>6</td><td>9</td></tr><tr><td>3</td><td>8</td><td>5</td><td>8</td></tr><tr><td>3</td><td>9</td><td>3</td><td>5</td></tr></table>	2	4	4	5	5	14	4	7	4	8	6	9	3	8	5	8	3	9	3	5	<table><tr><td>2</td><td>4</td><td>3</td><td>5</td></tr><tr><td>3</td><td>8</td><td>4</td><td>5</td></tr><tr><td>3</td><td>8</td><td>4</td><td>7</td></tr><tr><td>4</td><td>9</td><td>5</td><td>8</td></tr><tr><td>5</td><td>14</td><td>6</td><td>9</td></tr></table>	2	4	3	5	3	8	4	5	3	8	4	7	4	9	5	8	5	14	6	9	<table><tr><td>3.5</td><td>3.5</td><td>3.5</td><td>3.5</td></tr><tr><td>5.0</td><td>5.0</td><td>5.0</td><td>5.0</td></tr><tr><td>5.5</td><td>5.5</td><td>5.5</td><td>5.5</td></tr><tr><td>6.5</td><td>6.5</td><td>6.5</td><td>6.5</td></tr><tr><td>8.5</td><td>8.5</td><td>8.5</td><td>8.5</td></tr></table>	3.5	3.5	3.5	3.5	5.0	5.0	5.0	5.0	5.5	5.5	5.5	5.5	6.5	6.5	6.5	6.5	8.5	8.5	8.5	8.5	<table><tr><td>3.5</td><td>3.5</td><td>5.0</td><td>5.0</td></tr><tr><td>8.5</td><td>8.5</td><td>5.5</td><td>5.5</td></tr><tr><td>6.5</td><td>5.0</td><td>8.5</td><td>8.5</td></tr><tr><td>5.0</td><td>5.5</td><td>6.5</td><td>6.5</td></tr><tr><td>5.5</td><td>6.5</td><td>3.5</td><td>3.5</td></tr></table>	3.5	3.5	5.0	5.0	8.5	8.5	5.5	5.5	6.5	5.0	8.5	8.5	5.0	5.5	6.5	6.5	5.5	6.5	3.5	3.5
2	4	4	5																																																																																
5	14	4	7																																																																																
4	8	6	9																																																																																
3	8	5	8																																																																																
3	9	3	5																																																																																
2	4	3	5																																																																																
3	8	4	5																																																																																
3	8	4	7																																																																																
4	9	5	8																																																																																
5	14	6	9																																																																																
3.5	3.5	3.5	3.5																																																																																
5.0	5.0	5.0	5.0																																																																																
5.5	5.5	5.5	5.5																																																																																
6.5	6.5	6.5	6.5																																																																																
8.5	8.5	8.5	8.5																																																																																
3.5	3.5	5.0	5.0																																																																																
8.5	8.5	5.5	5.5																																																																																
6.5	5.0	8.5	8.5																																																																																
5.0	5.5	6.5	6.5																																																																																
5.5	6.5	3.5	3.5																																																																																

Quantile Normalization

