

# Exploratory analysis

---

Jeff Leek

@jtleek

[www.jtleek.com](http://www.jtleek.com)

Key ideas

Visualization

Summarization

Showing the data

Not being misled

# Why explore data?

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

# Background

## Perceptual tasks

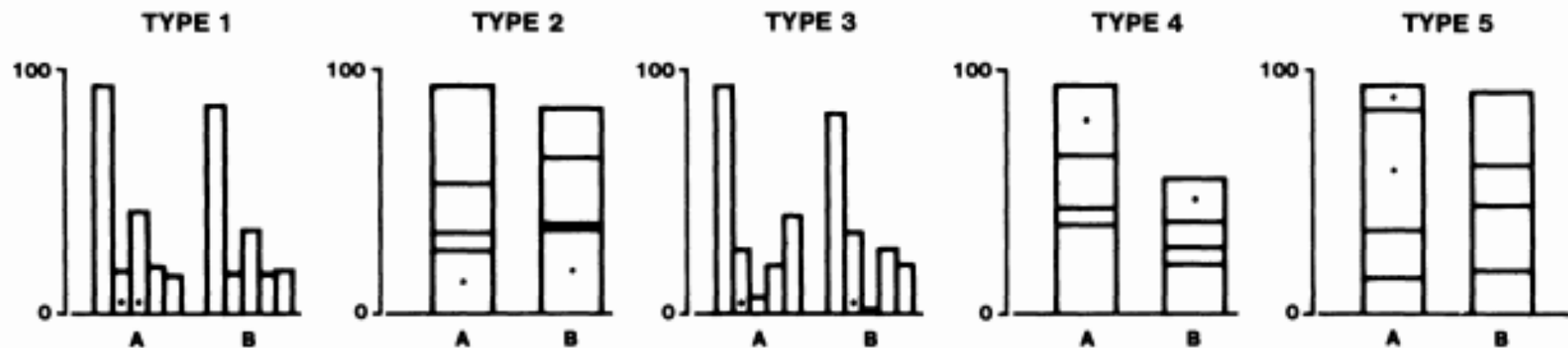
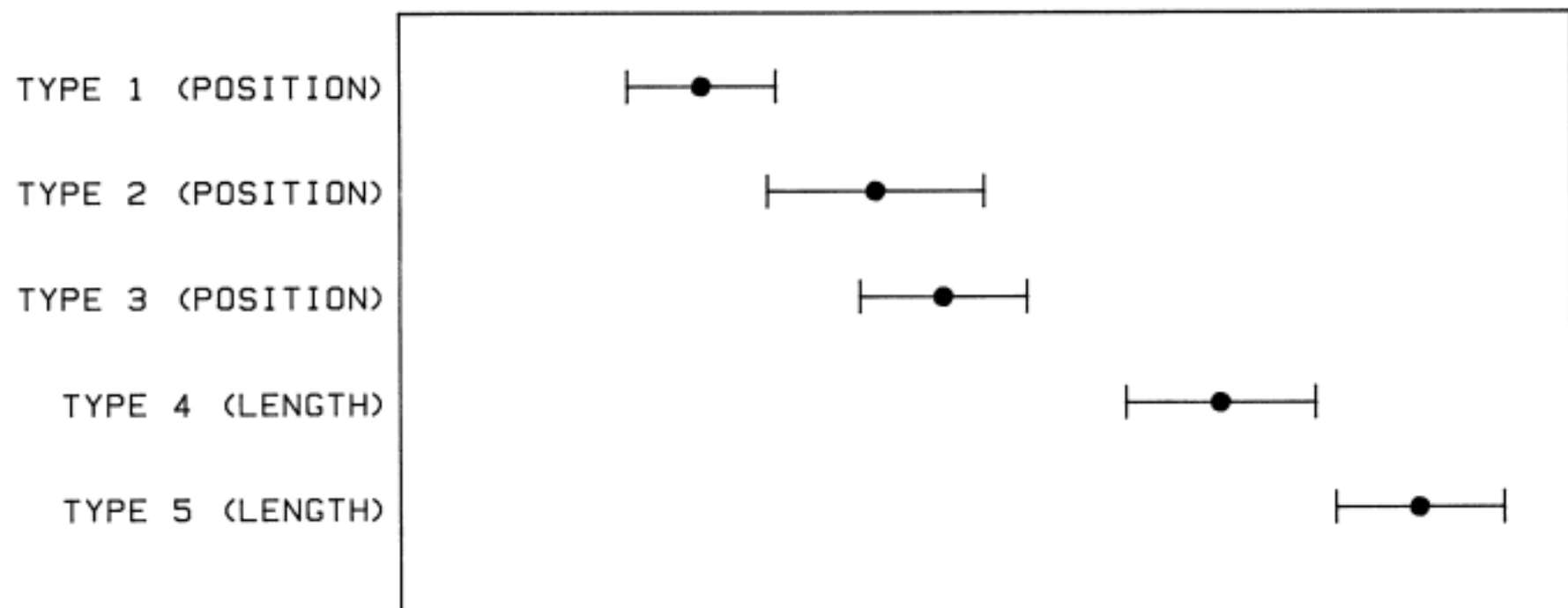
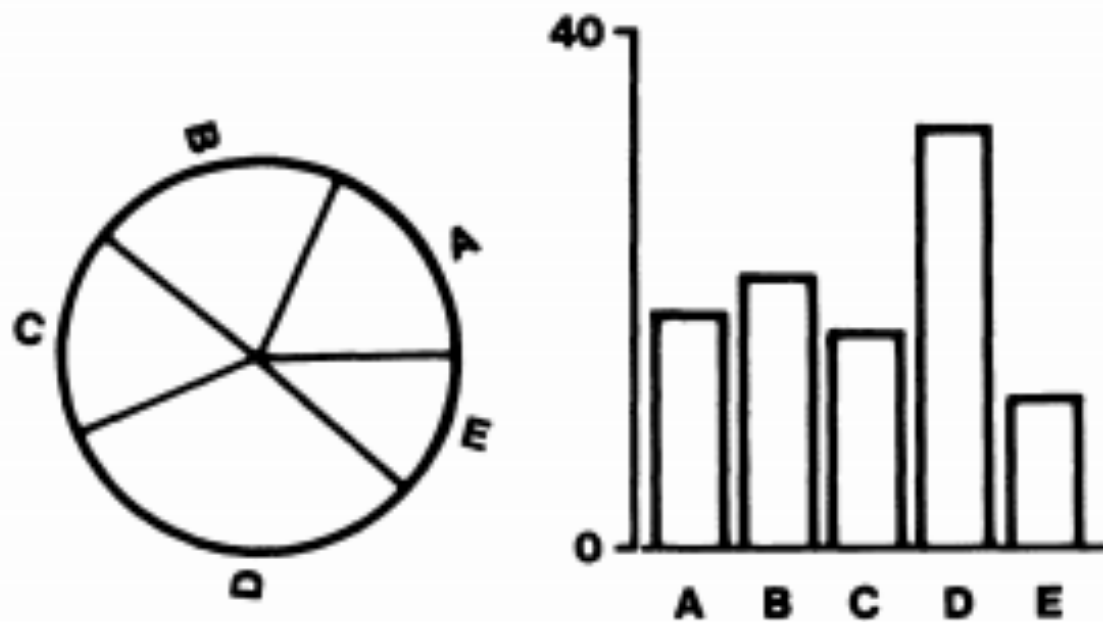
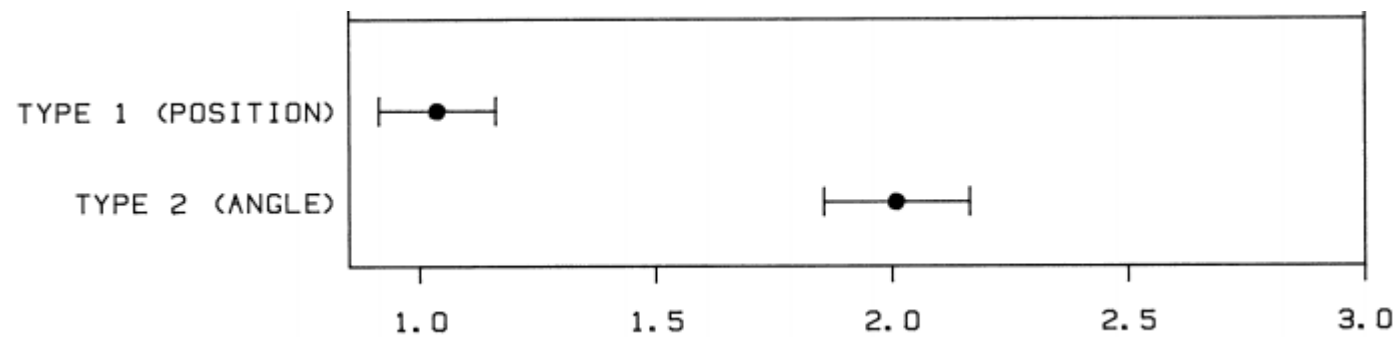


Figure 4. Graphs from position-length experiment.





*Figure 3. Graphs from position-angle experiment.*

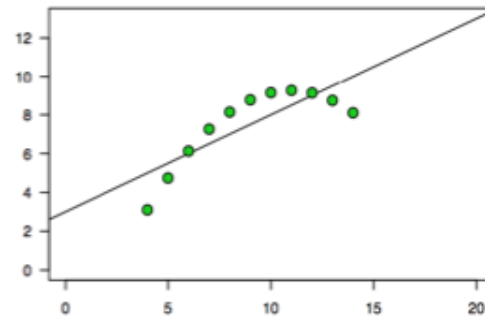
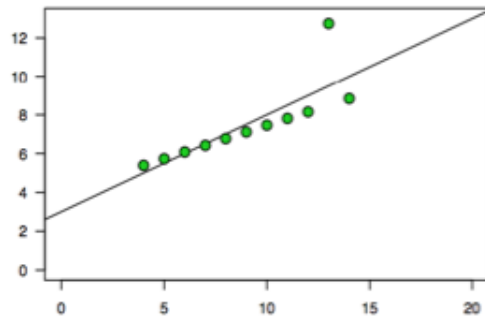
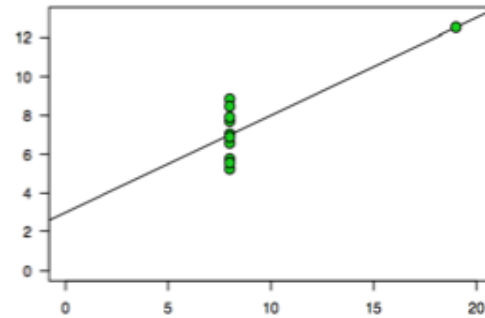
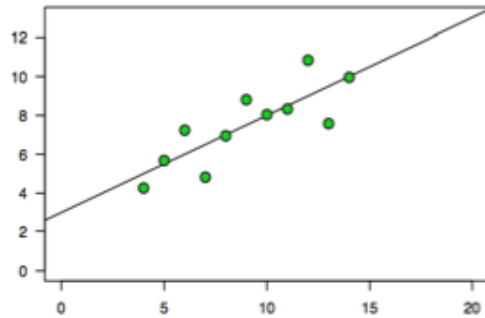




**Background**

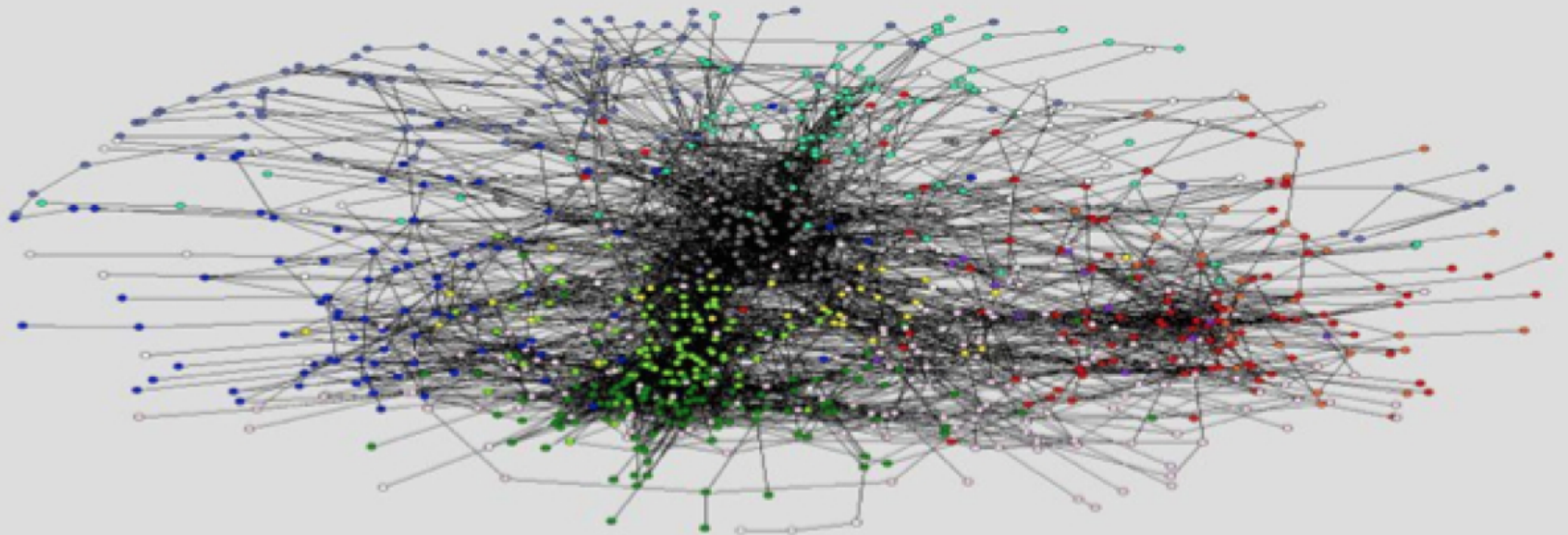
**Graphs reveal structure  
summaries don't**

$\hat{\beta}_0 = 3.0$ ,  $\hat{\beta}_1 = 0.5$ , p-value (slope) = 0.002,  $R^2 = 0.67$ .



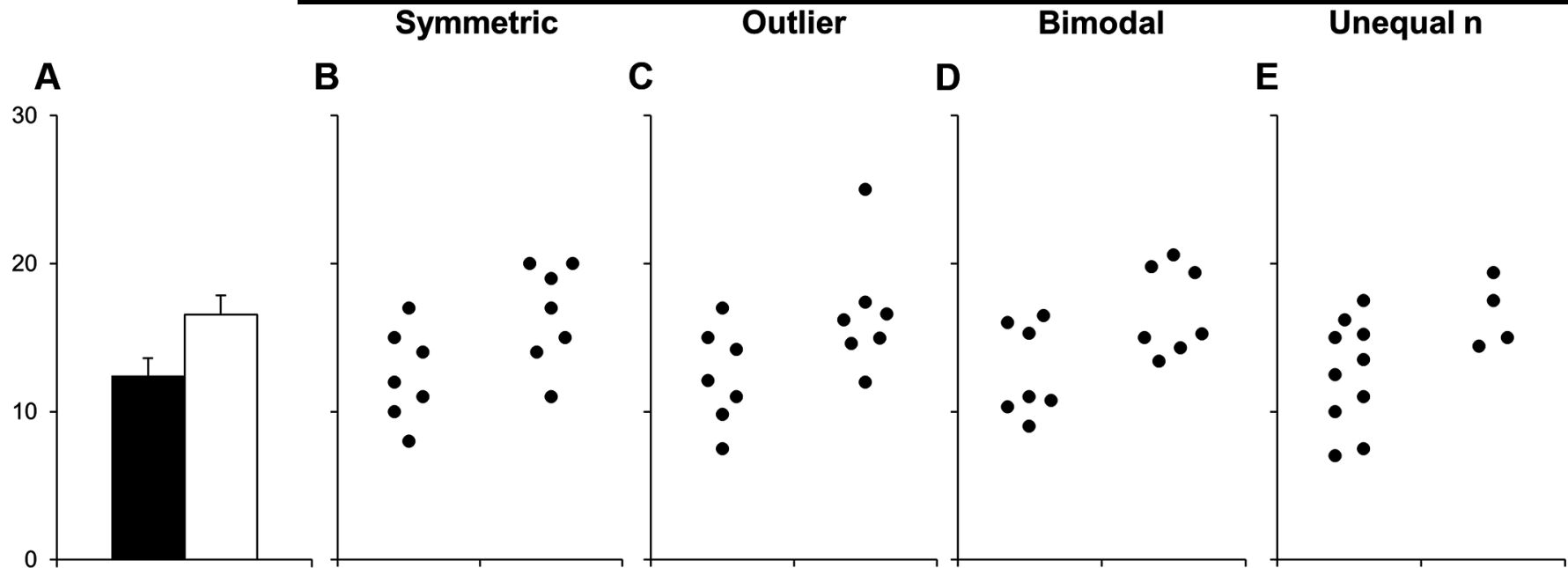
**Basic principles**

**Avoid ridiculograms**



*Ridiculogram: meaningless albeit visually impressive image of a network*

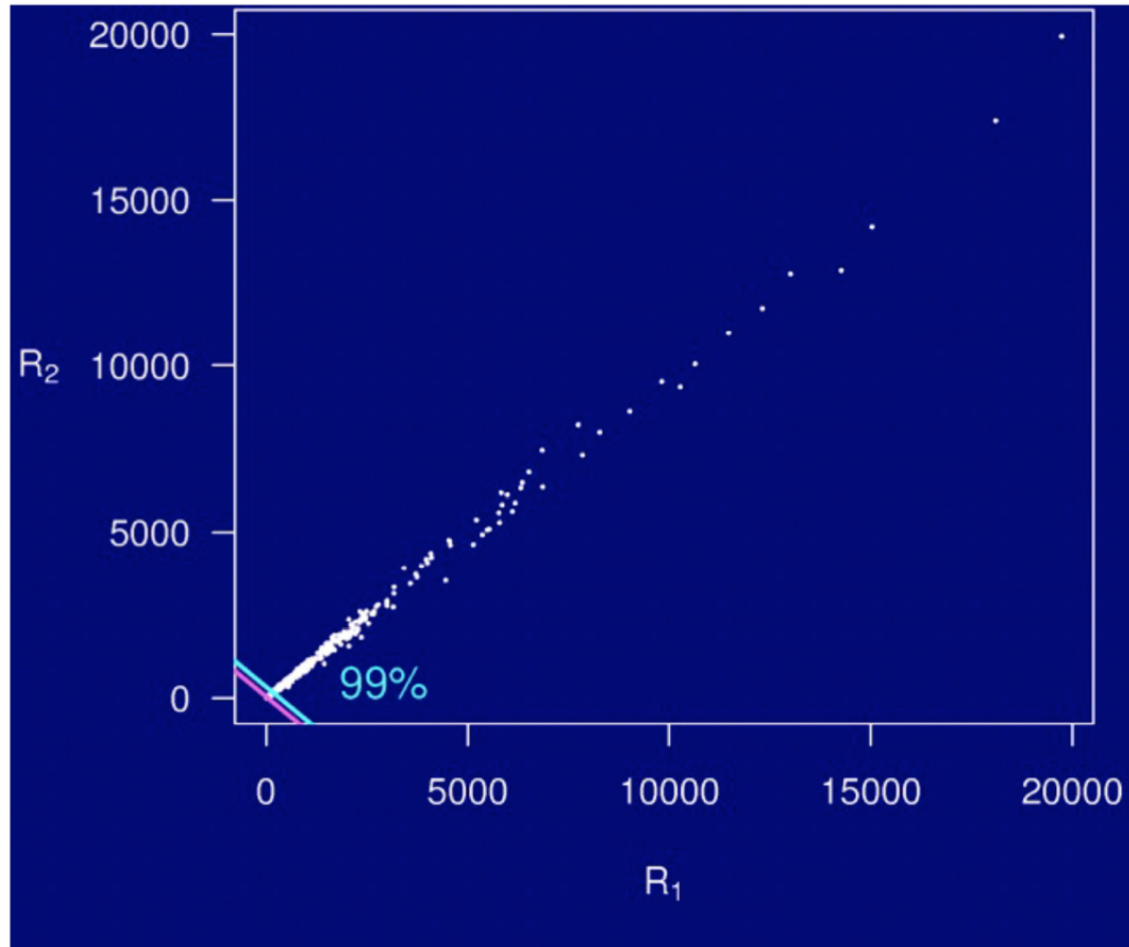
**Basic principles**  
**Show the data**



Test	p value			
T-test: Equal var.	0.035	0.050	0.026	0.063
T-test: Unequal var.	0.035	0.050	0.026	0.035
Wilcoxon	0.054	0.073	0.128	0.103

**Basic principles**

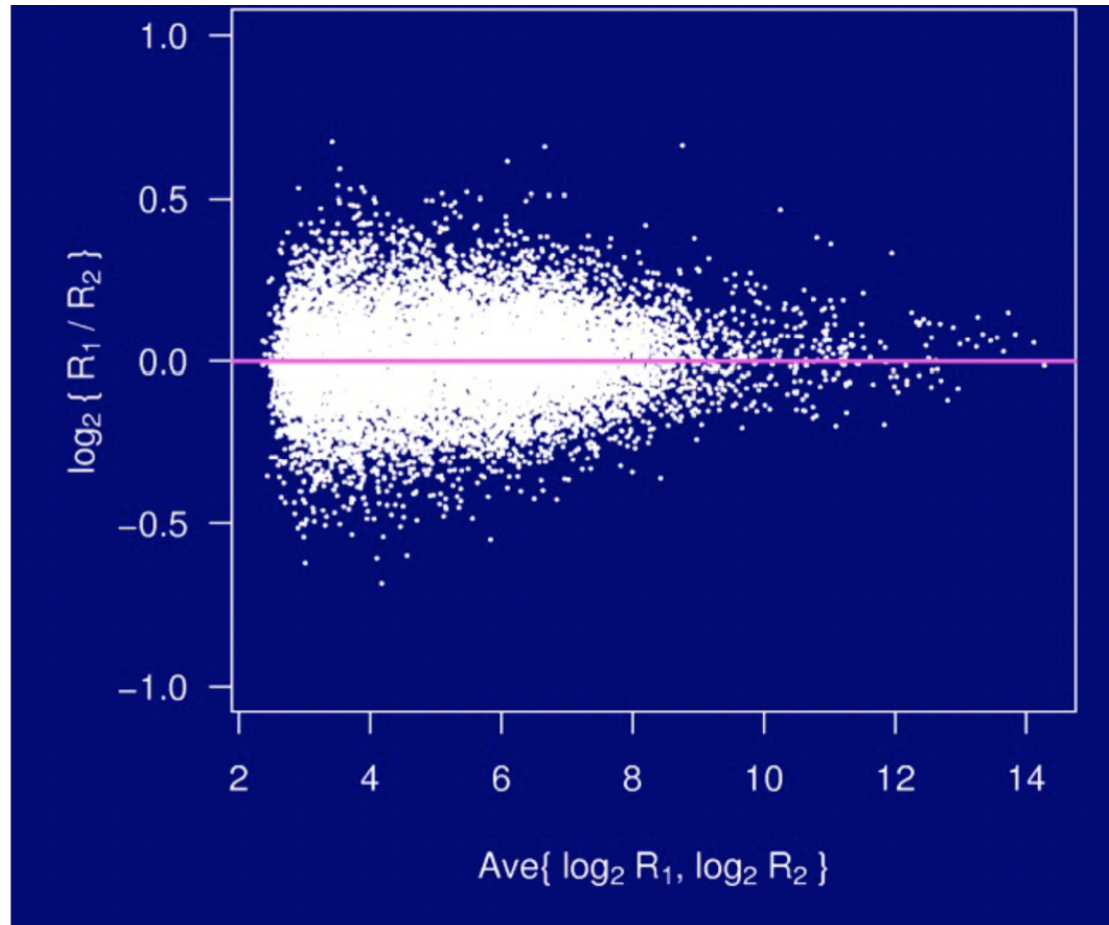
**Be careful with scale**





**Basic principles**

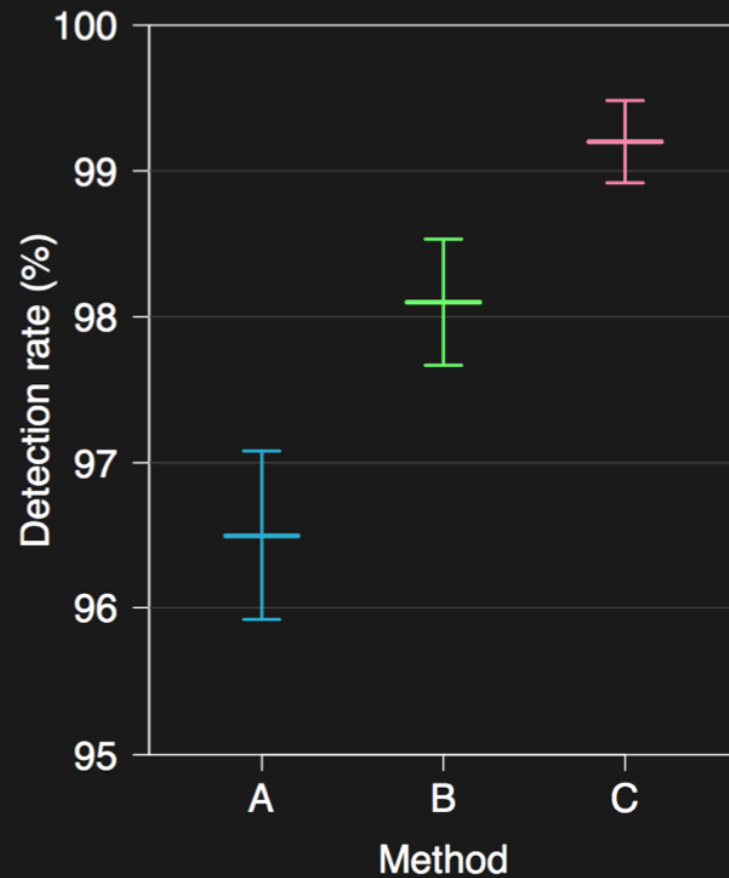
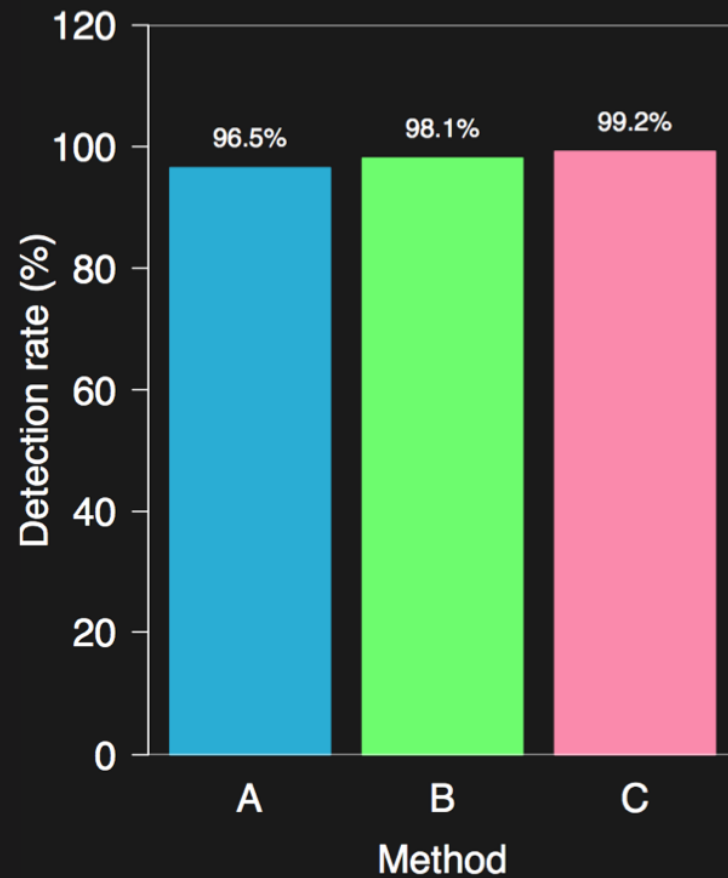
**Compare things directly**



**Basic principles**

**Use common scales**

**Start at zero**



## Further resources

- Karl Broman's guide to displaying data
  - [https://www.biostat.wisc.edu/~kbroman/presentations/IowaState2013/graphs\\_combined.pdf](https://www.biostat.wisc.edu/~kbroman/presentations/IowaState2013/graphs_combined.pdf)
- Data visualization at Nature
  - <http://blogs.nature.com/methagora/2013/07/data-visualization-points-of-view.html>