

# How I learned to stop visualizing and love statistics

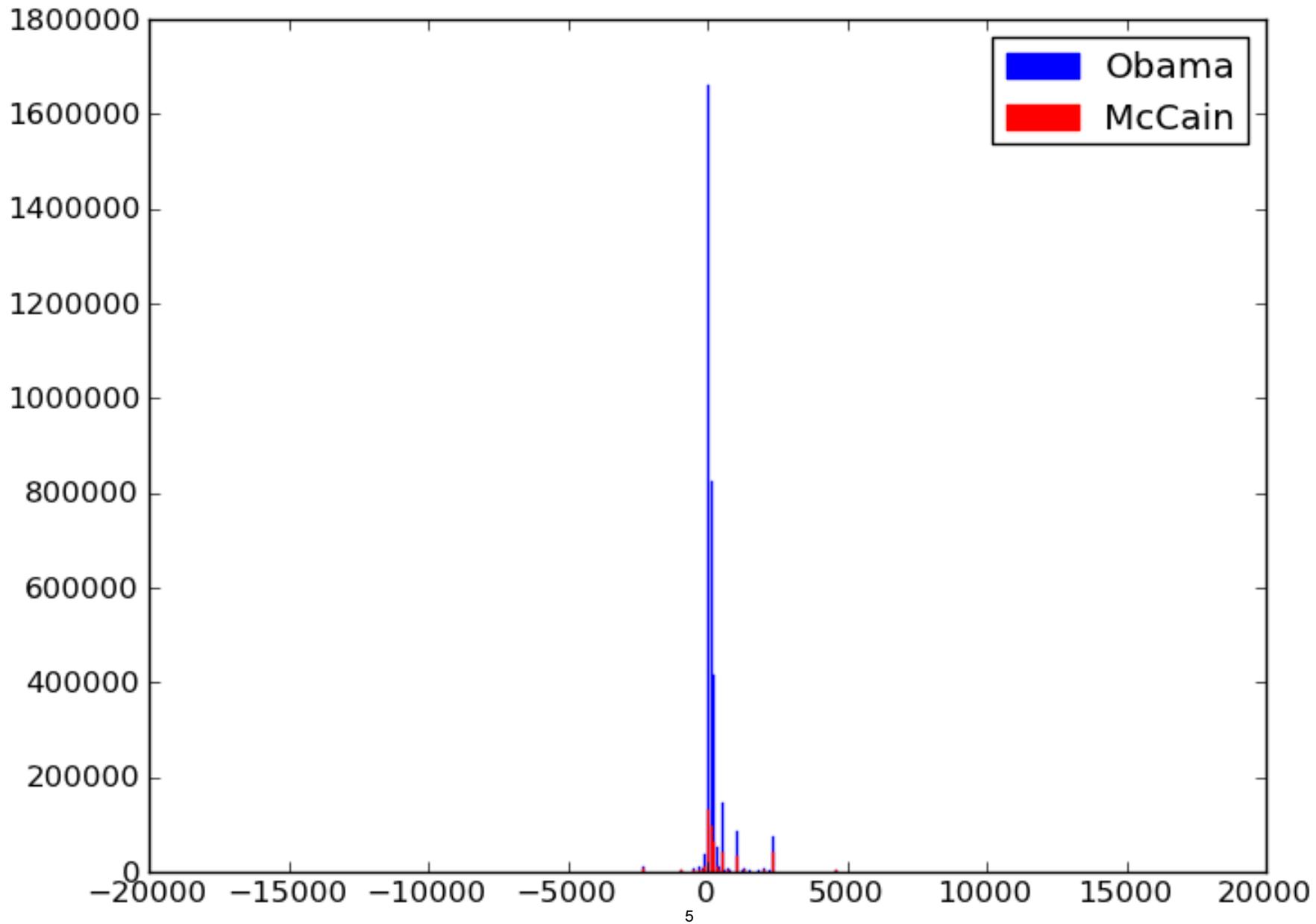
You have a hunch

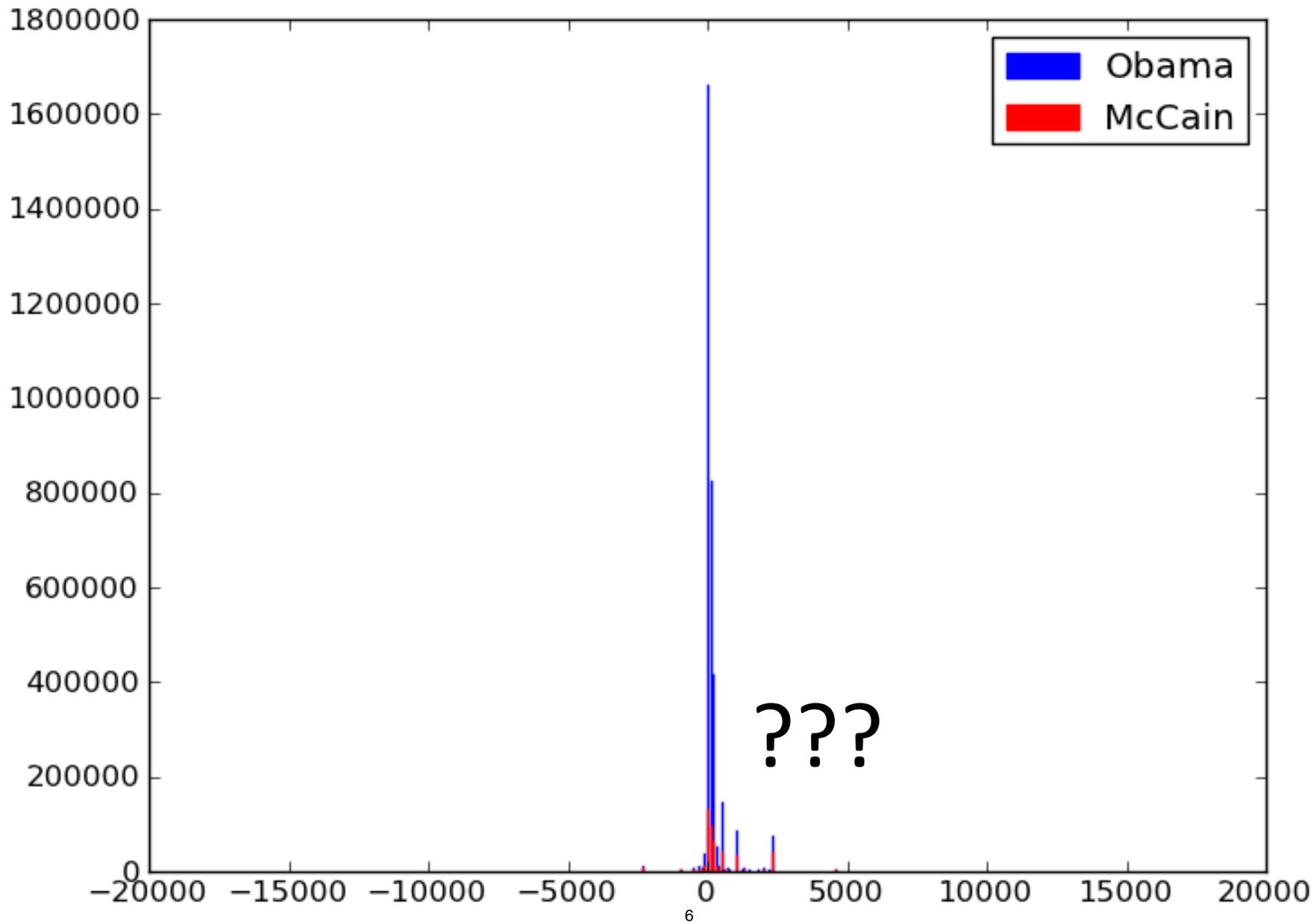
Visualizations → sanity check

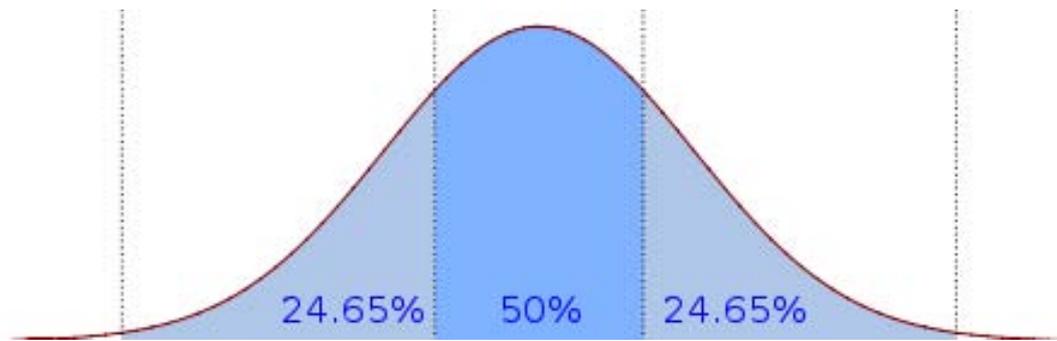
Statistics → quantify the hunch

(Visualizations → storytelling)

Someone says:  
“Obama got more small campaign  
contributions than McCain”

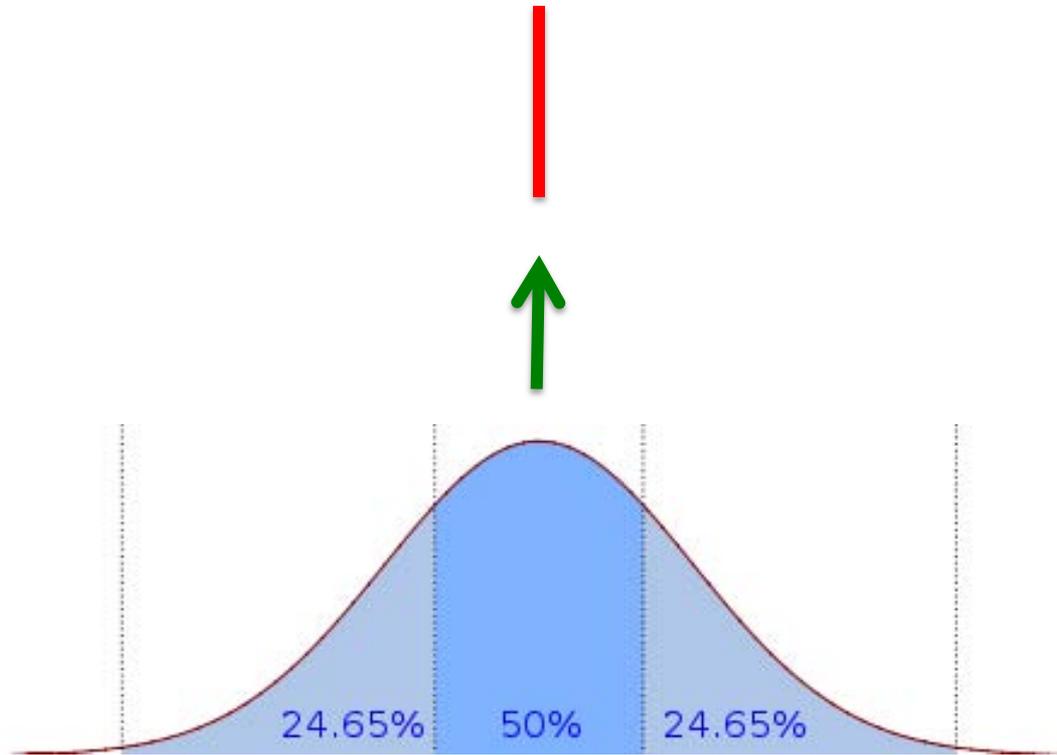




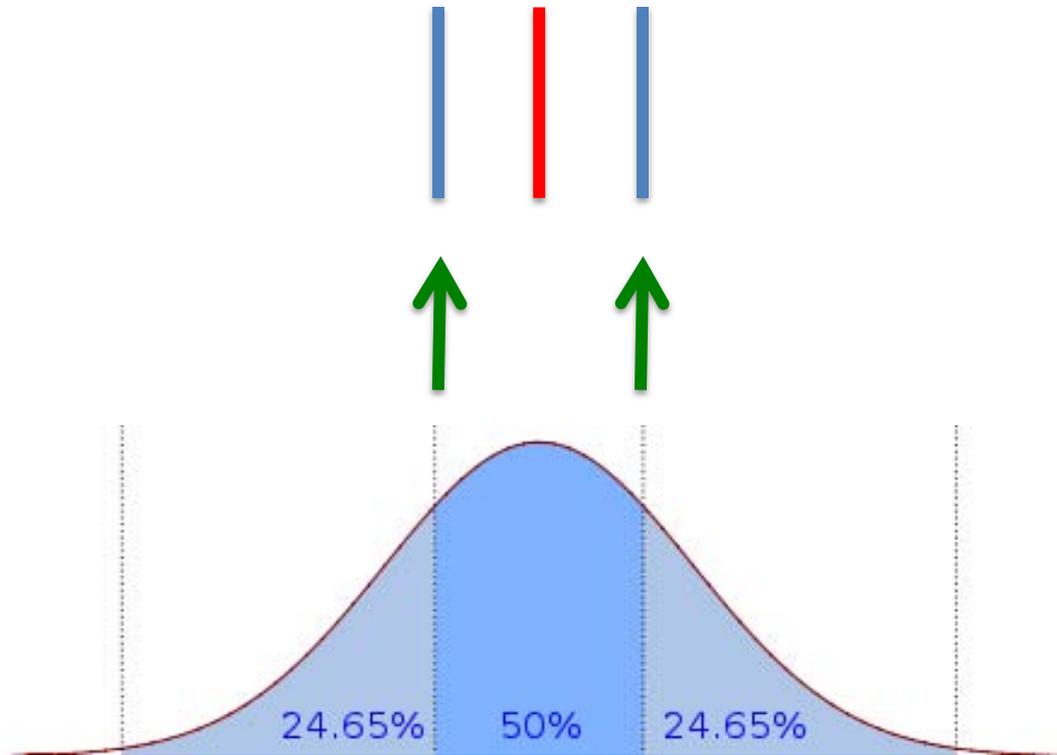


© Jhguch on Wikipedia. All rights reserved. This content is excluded from our Creative Commons license. For more information, see <http://ocw.mit.edu/fairuse>.

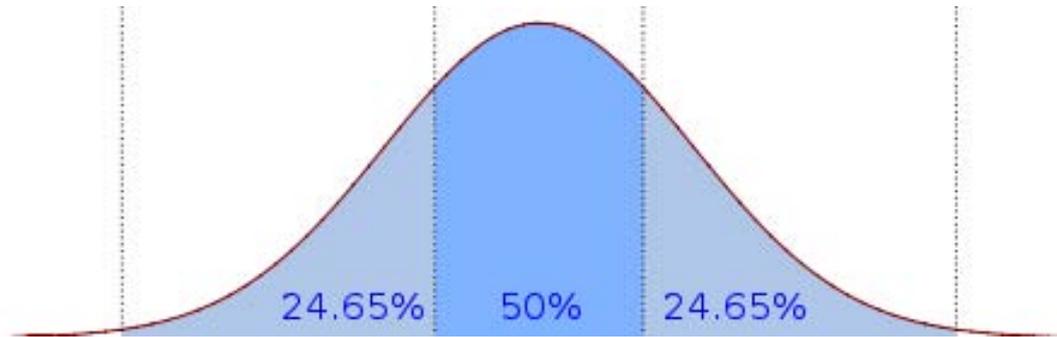
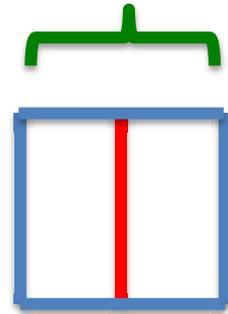
# Median



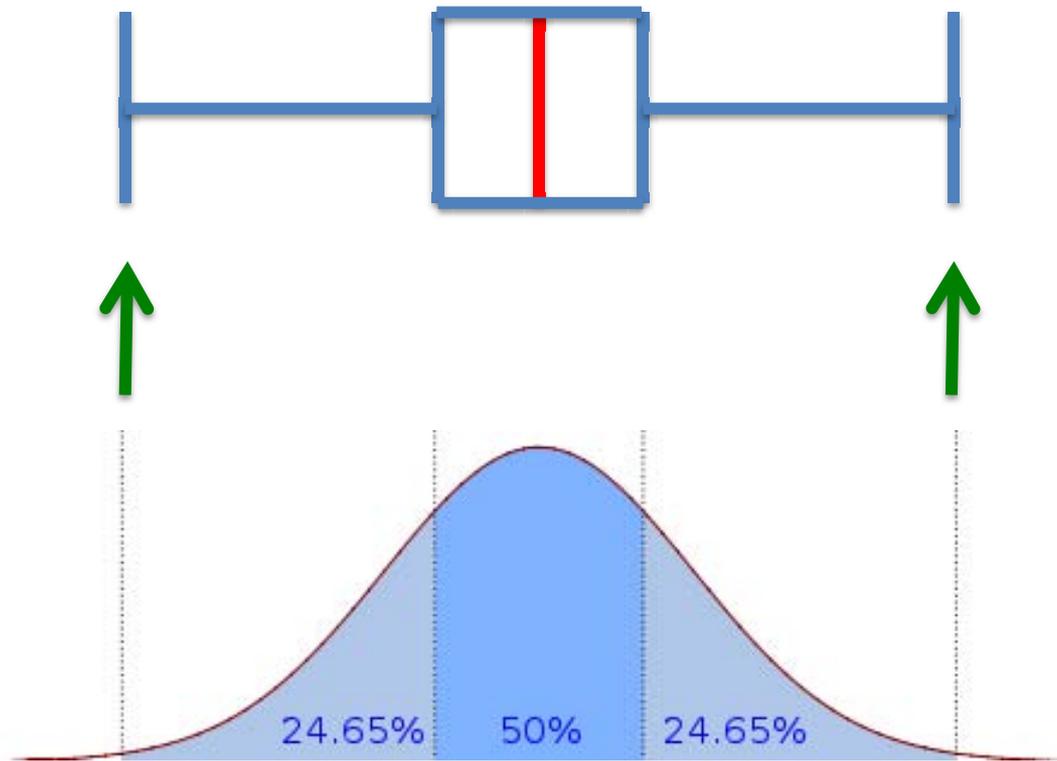
25% 75%



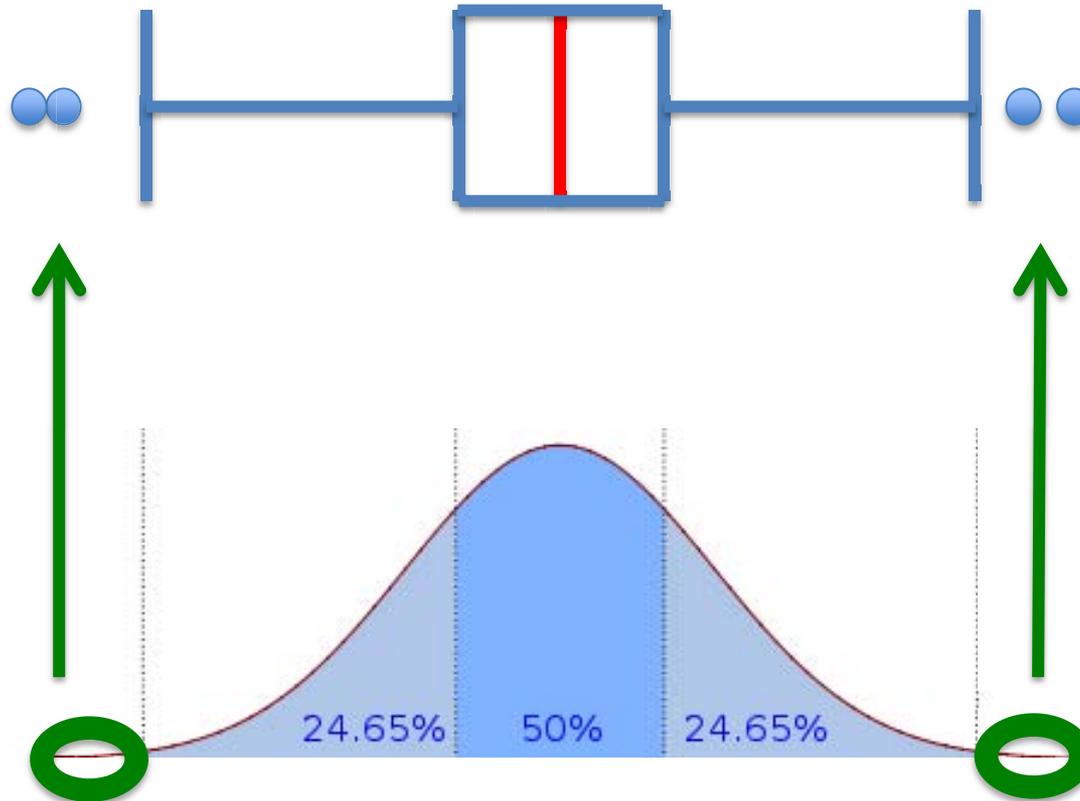
# Inner Quartile Range



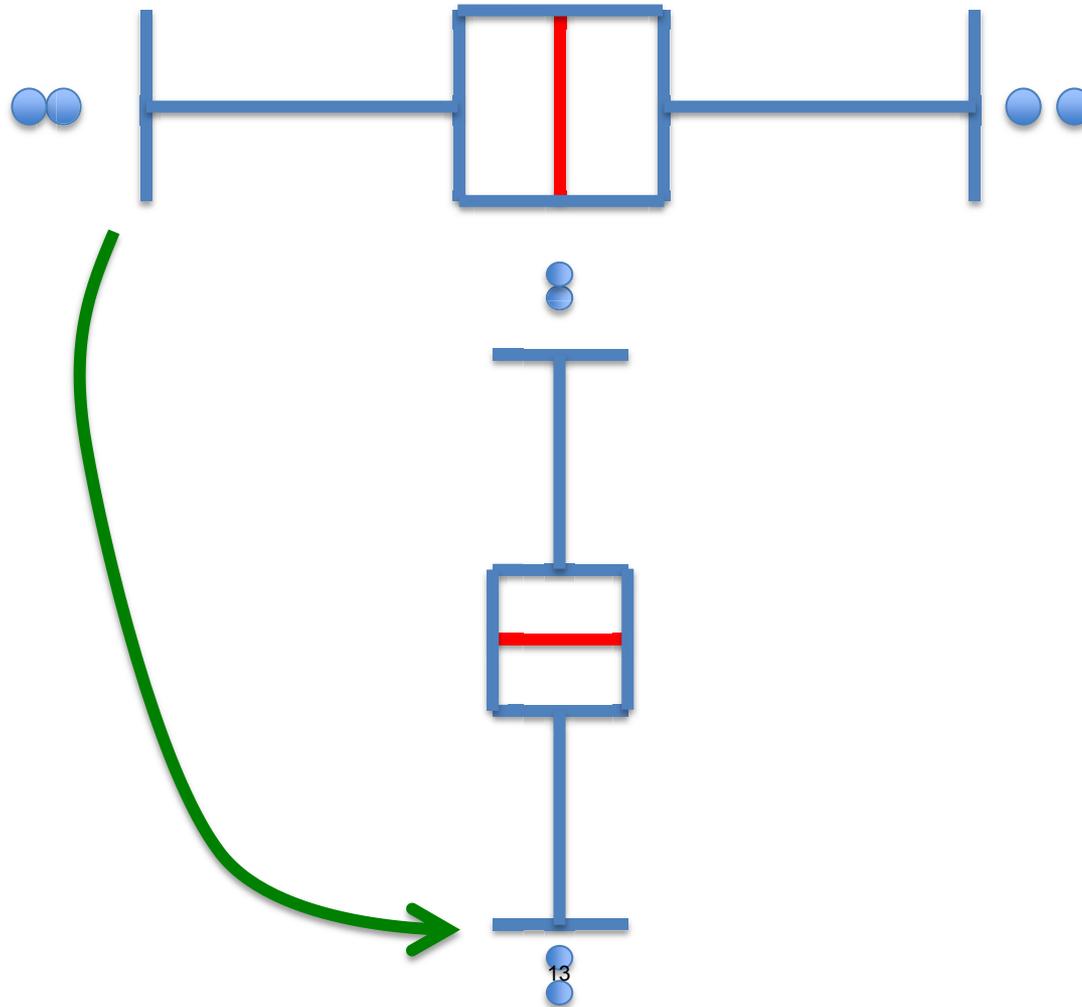
# Whiskers / Extremes



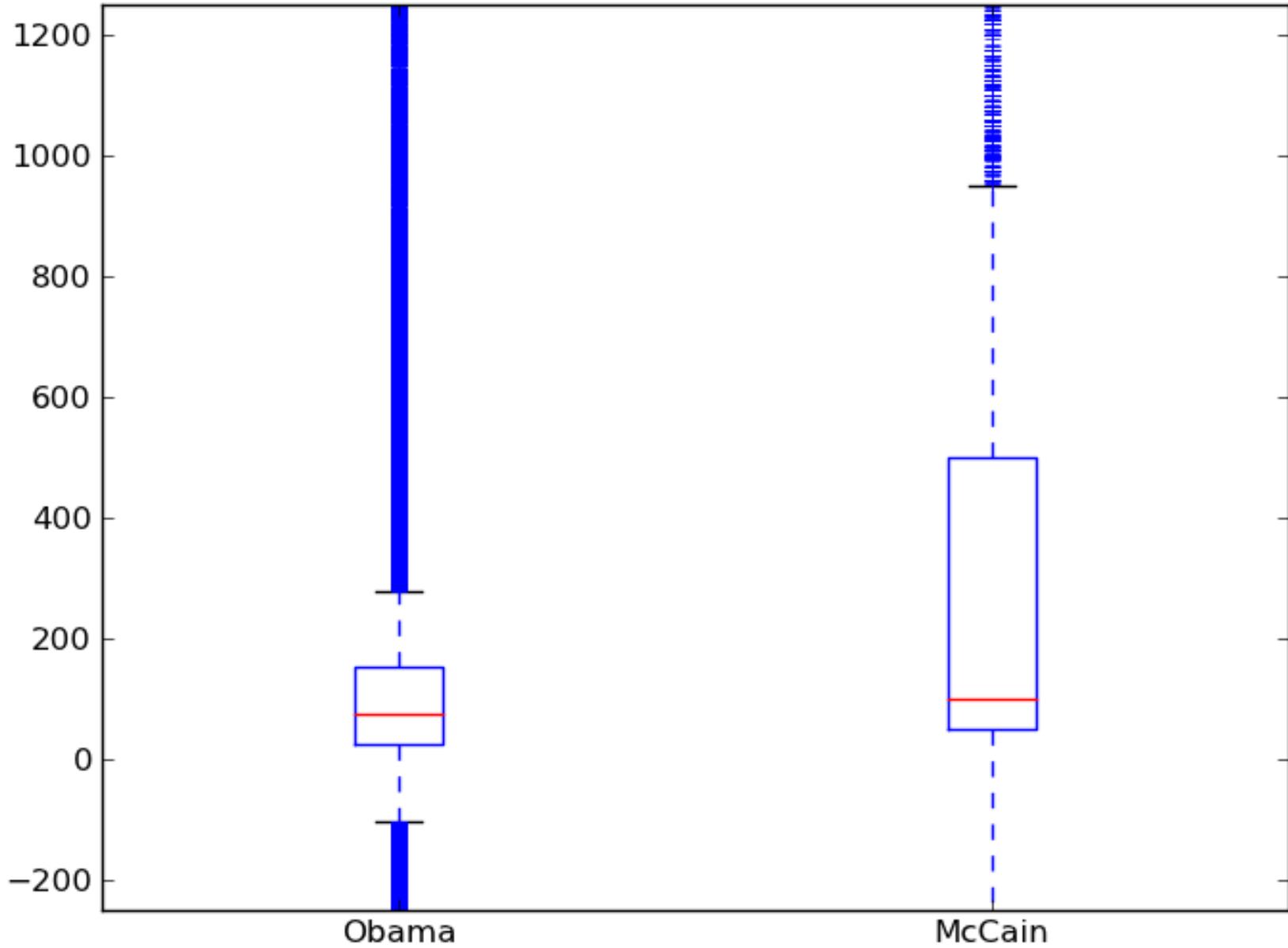
# Outliers



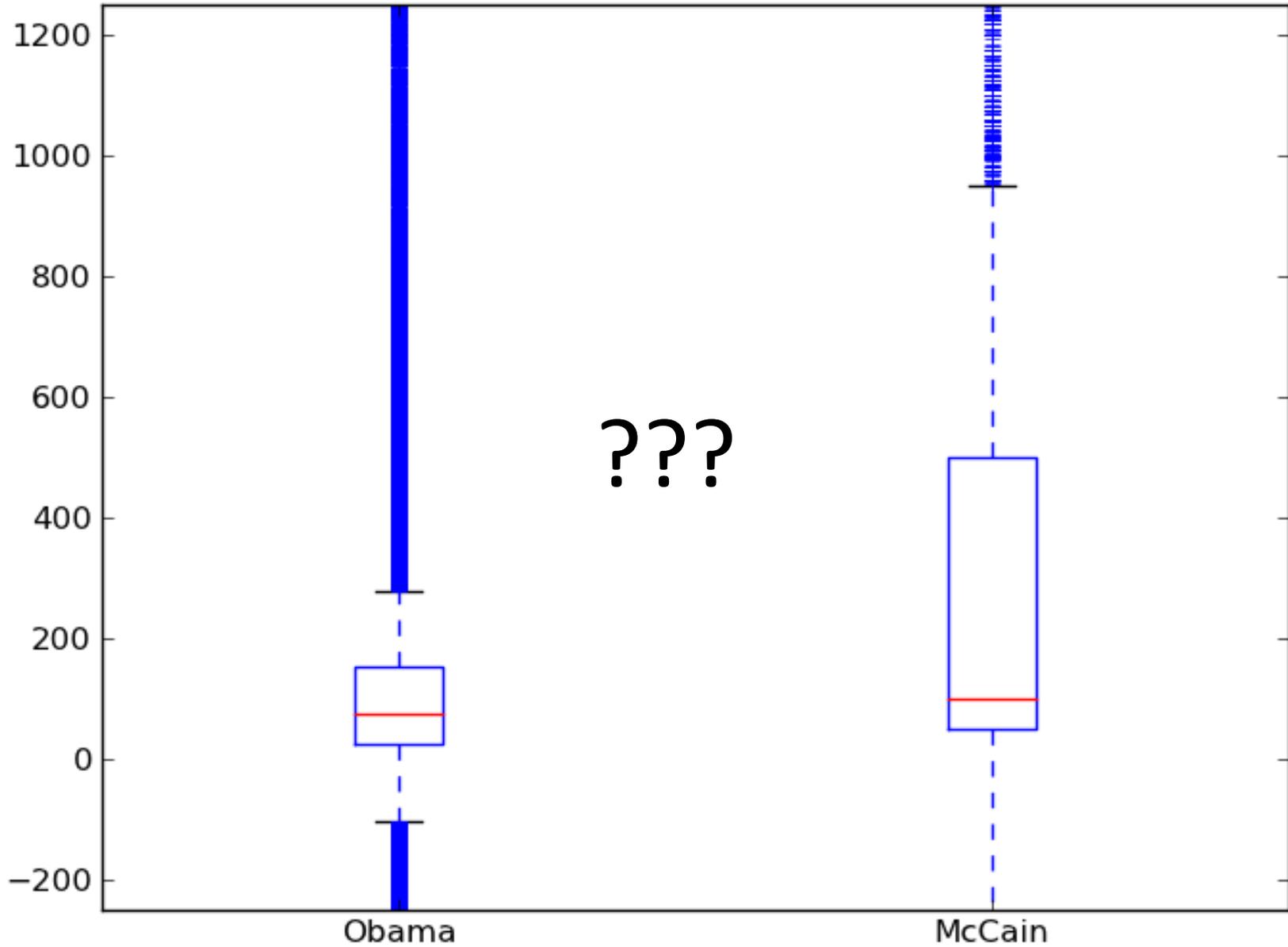
# Box-and-Whiskers Plot



# Obama vs. McCain Contributions



# Obama vs. McCain Contributions



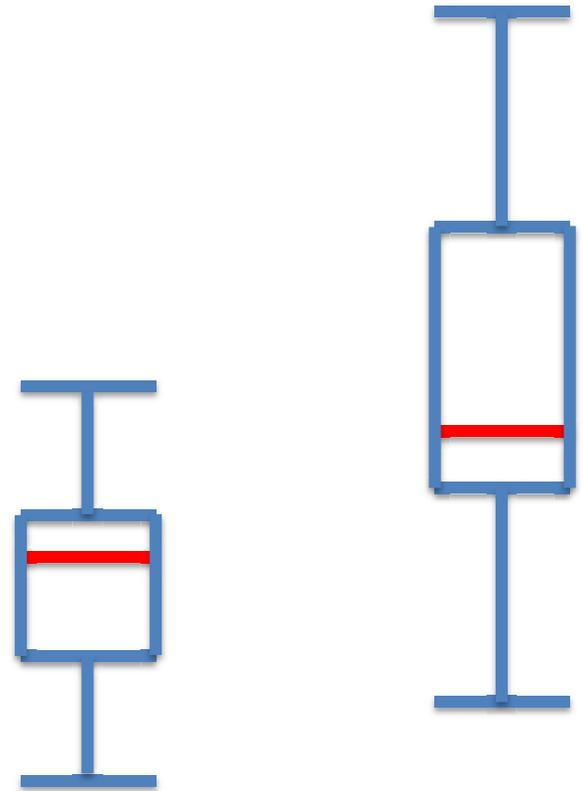
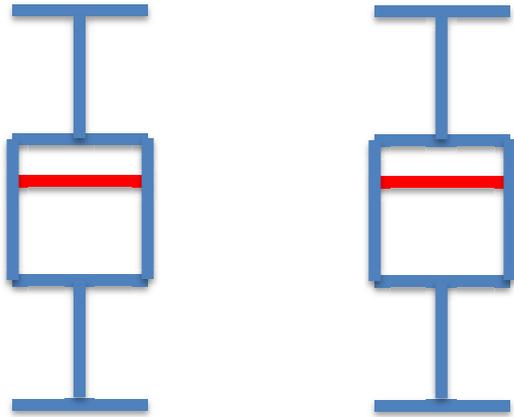
Are they actually different?



T-Test

# Assume

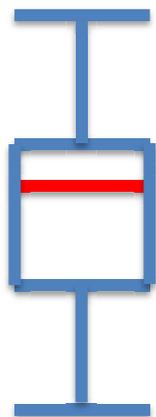
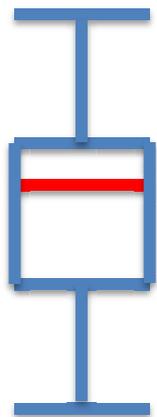
# Reality



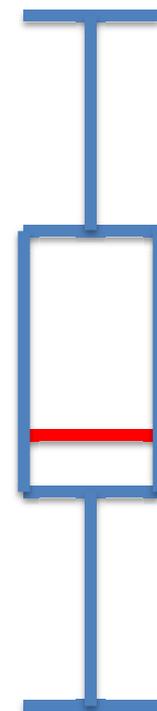
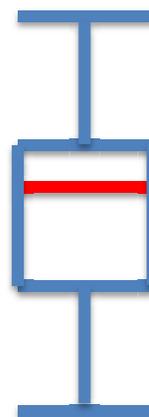
Obama McCain

Obama McCain

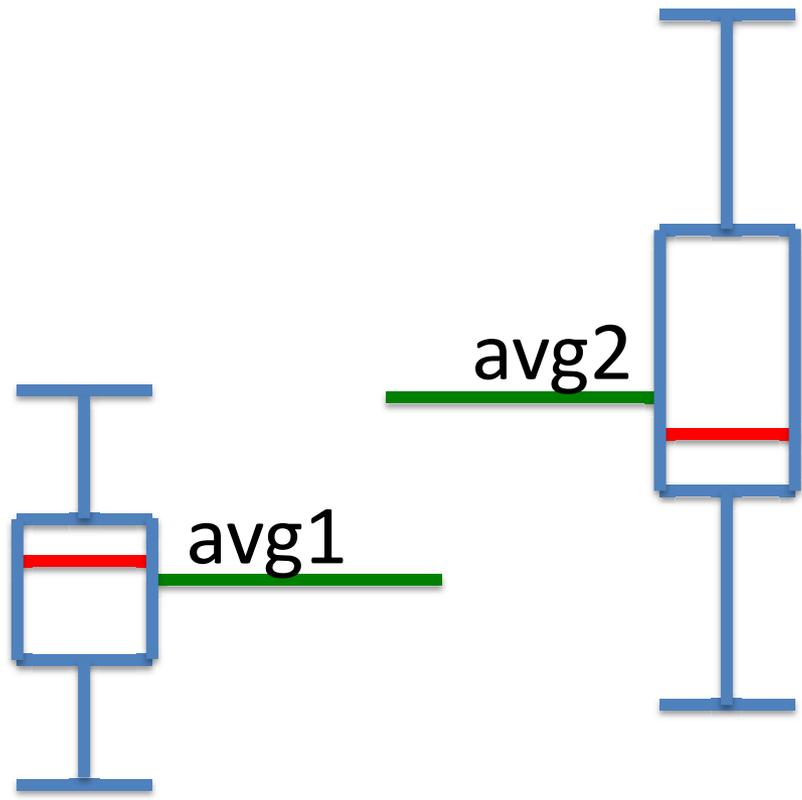
How likely is given ?



Obama McCain

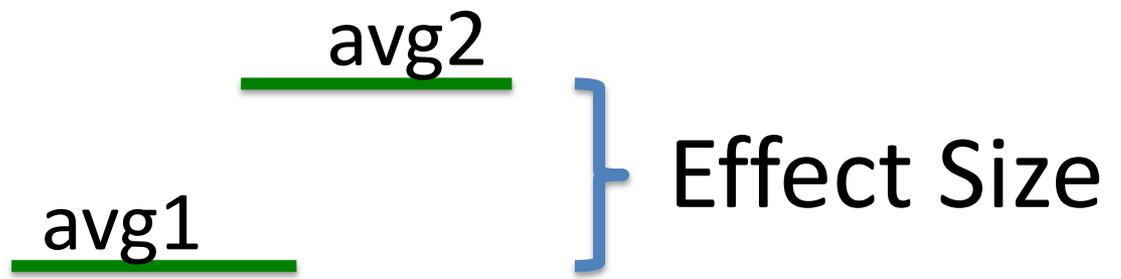


Obama McCain



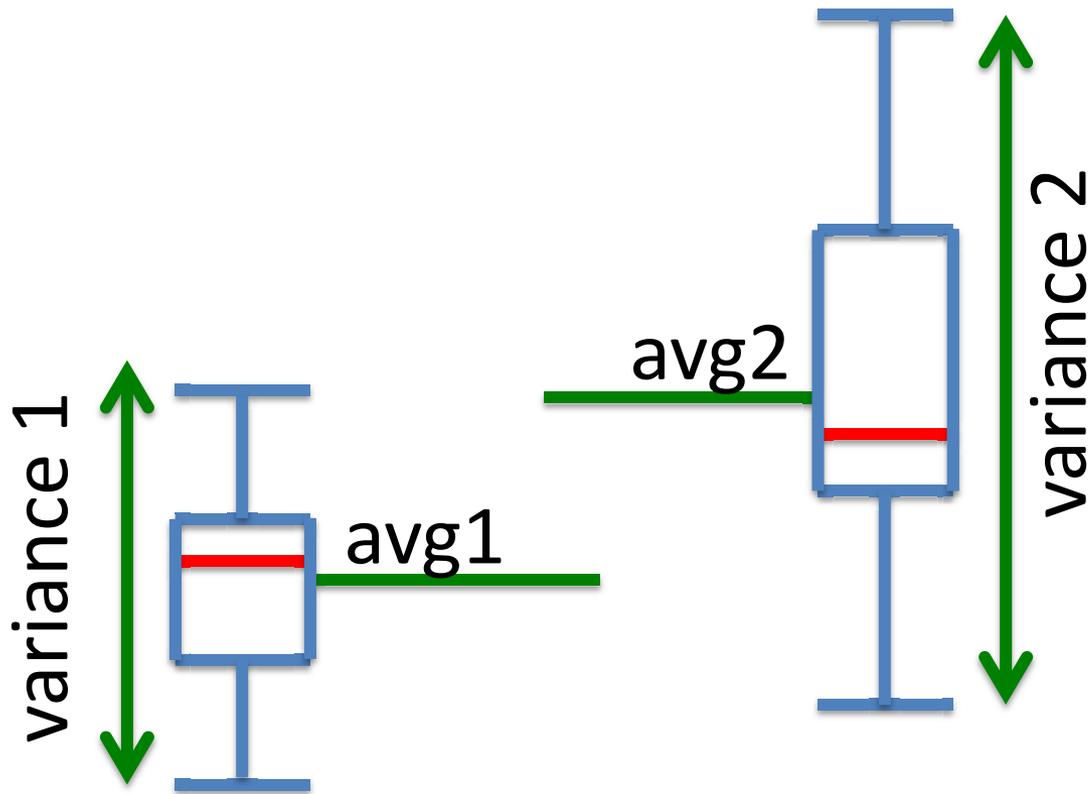
Obama

McCain



Obama

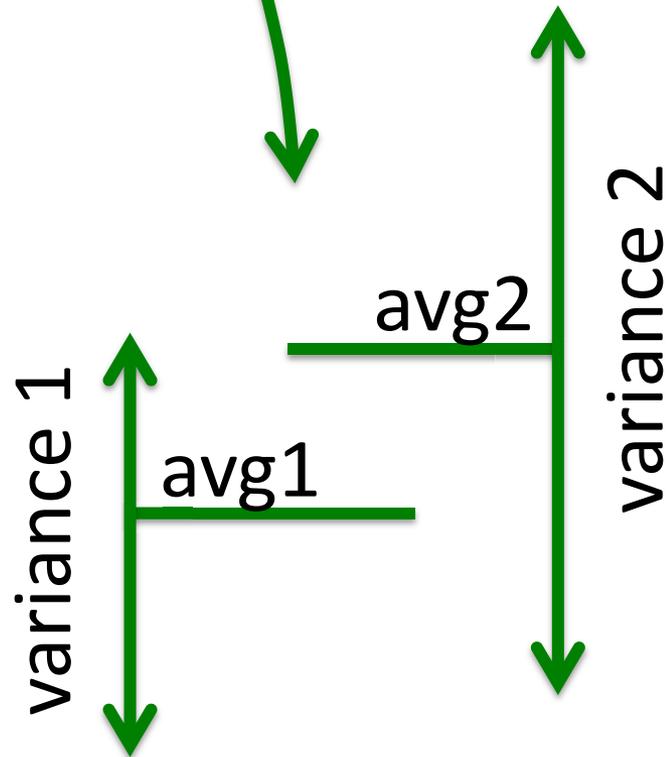
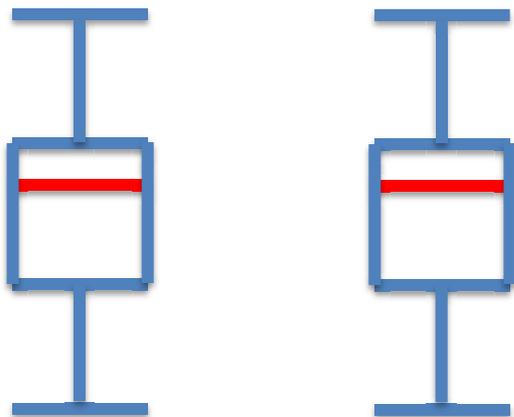
McCain



Obama

McCain

How likely is given ?



How likely are they equal  
given avg/variance differences?



Probability  $p$



$p$  is low

Obama, McCain  
are different  
(significant)



$p$  is high

Don't trust  
the difference  
(not significant)

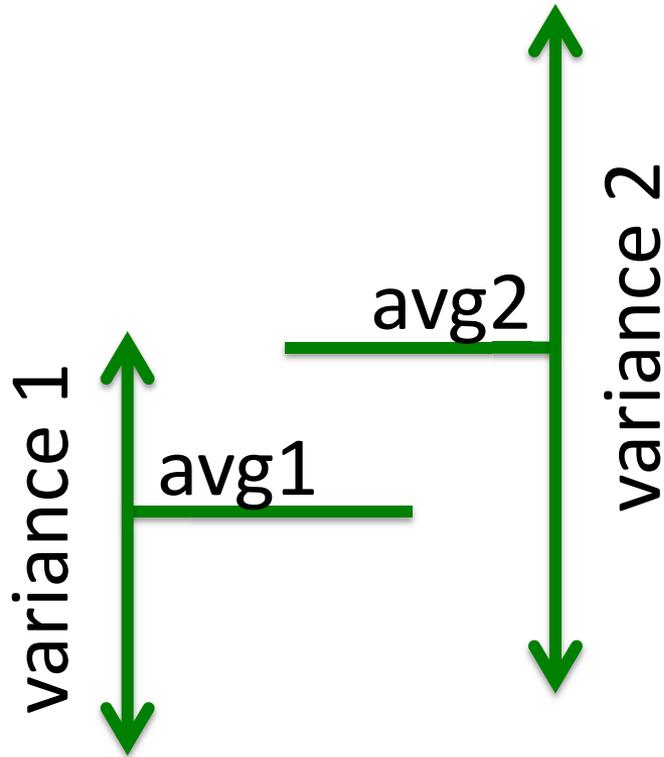
# Significance is binary

- Pick a threshold: .01? .05?
- Is  $p >$  threshold, or  $\leq$  threshold?

$p \leq .05?$  significant

$p > .05?$  don't trust the difference

# T-Test Significance



+

# Samples

Obama: >1M

McCain: >1M

# Correlation, Linear Regression

# County Health Rankings

- Every county in USA
- Years of Potential Life Lost (YPLL): early morbidity
  - less is good
  - more is bad
- Median income, % population w/ diabetes, % population under 18, ...

What is correlated with early death in a community?

Burgers

Sleep

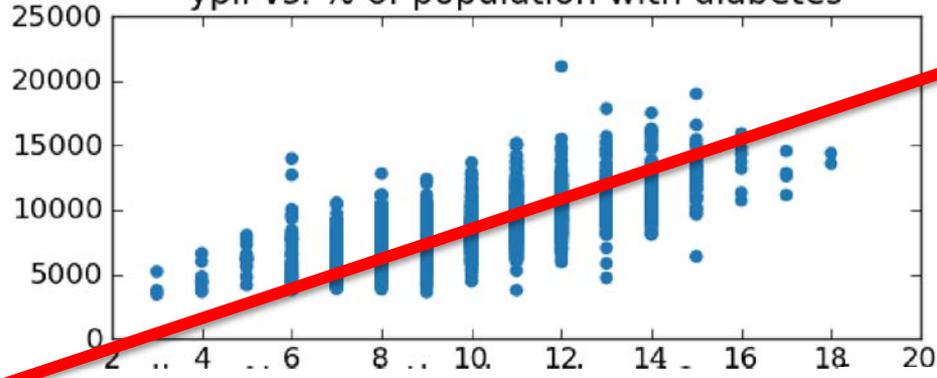
Education

Exercise

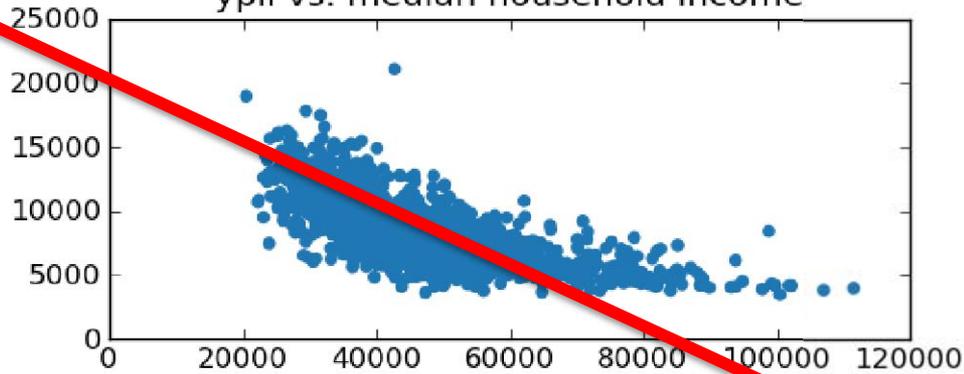
# Rappers

Your theory here

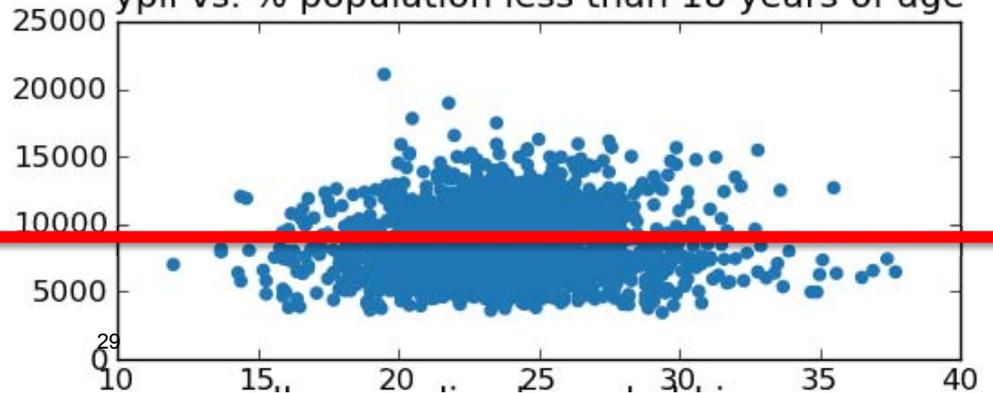
ypll vs. % of population with diabetes



ypll vs. median household income



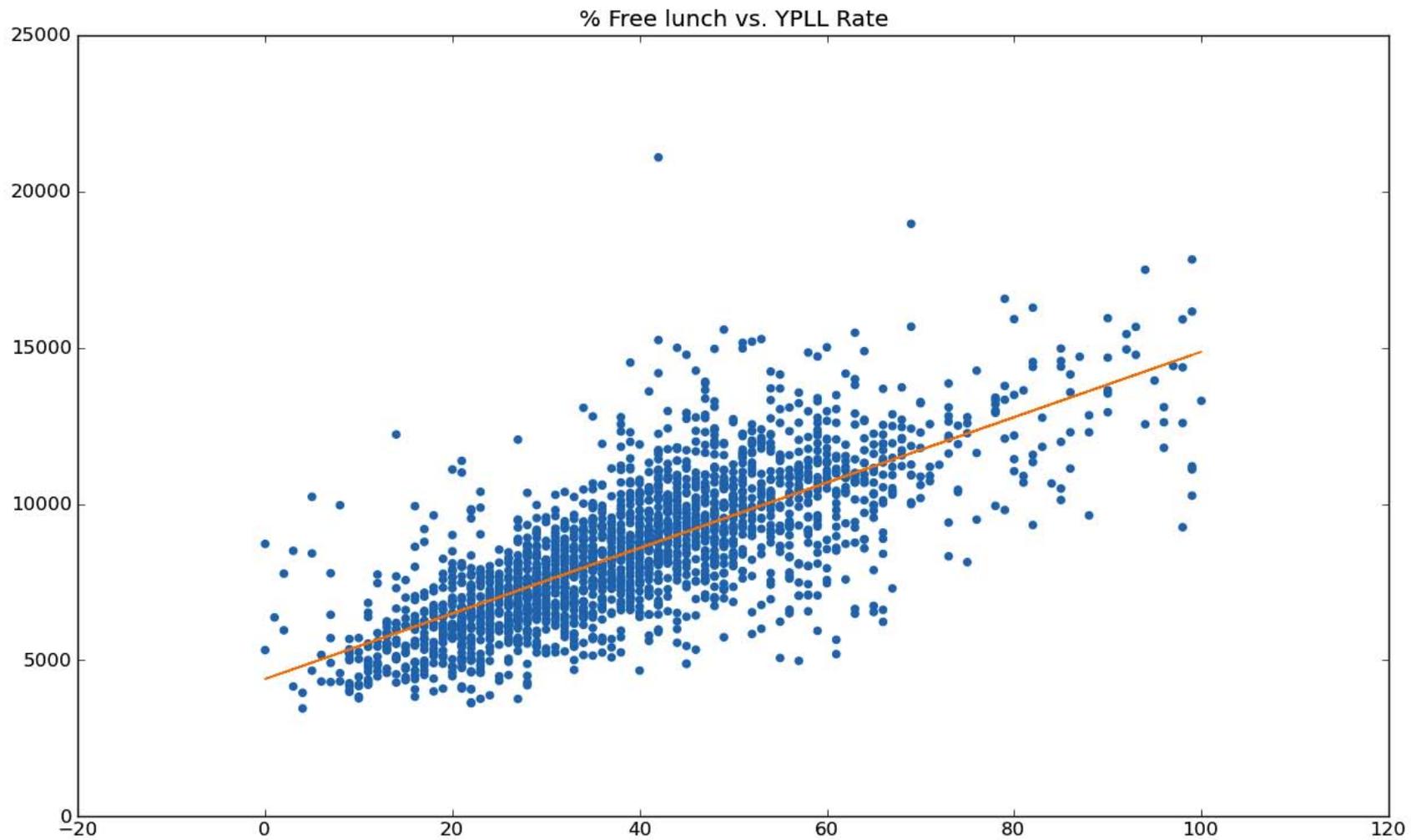
ypll vs. % population less than 18 years or age



Line coefficients:  $y = mx + b$

Correlation amount:  $R^2$  (0 to 1)

Significance:  $p < .05?$



# Correlation != Causation

Correlation



Causal Hunch



Randomized Trial



T-Test!

MIT OpenCourseWare  
<http://ocw.mit.edu>

Resource: How to Process, Analyze and Visualize Data  
Adam Marcus and Eugene Wu

The following may not correspond to a particular course on MIT OpenCourseWare, but has been provided by the author as an individual learning resource.

For information about citing these materials or our Terms of Use, visit: <http://ocw.mit.edu/terms>.