

How I learned to stop
visualizing and love statistics

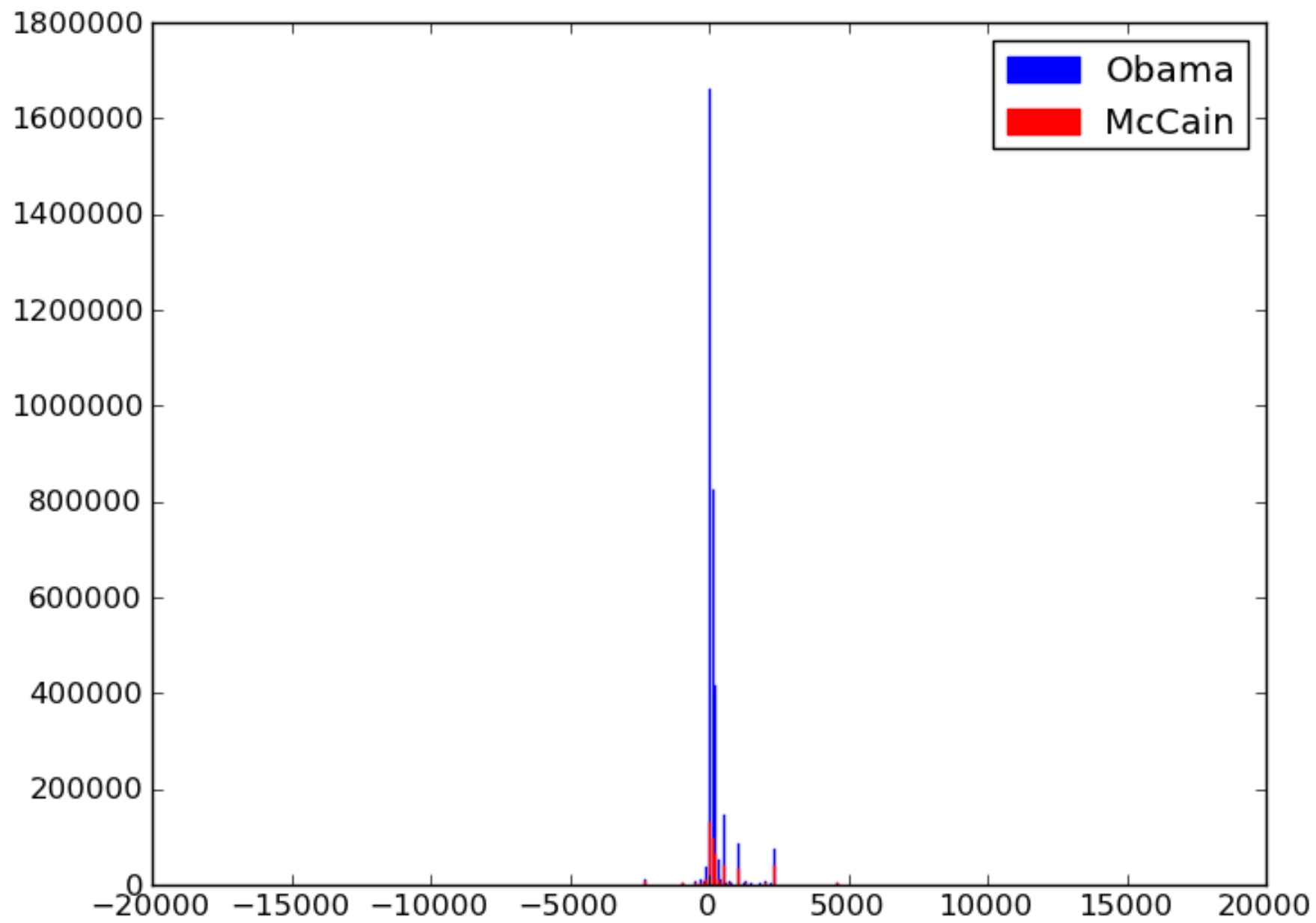
You have a hunch

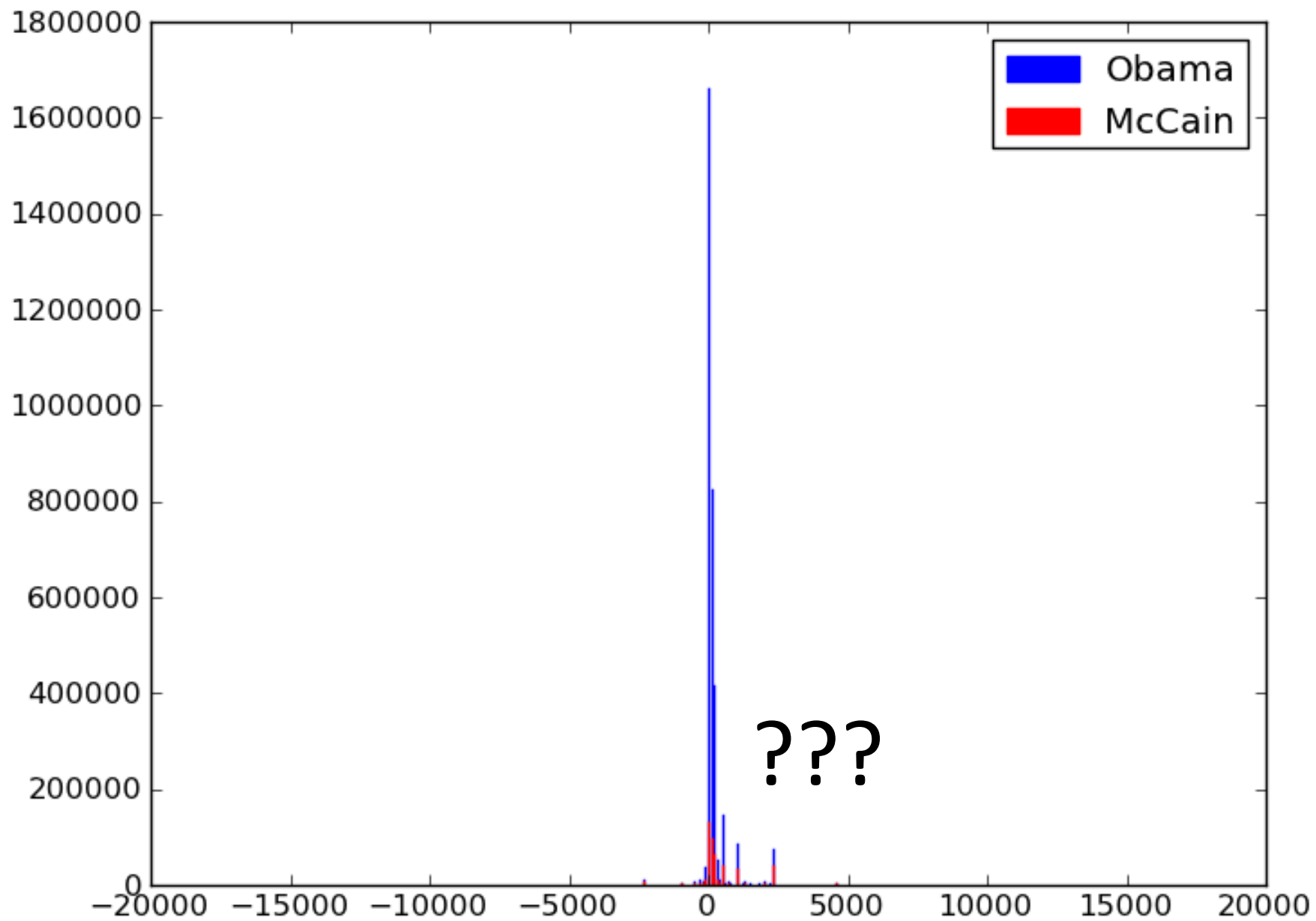
Visualizations → sanity check

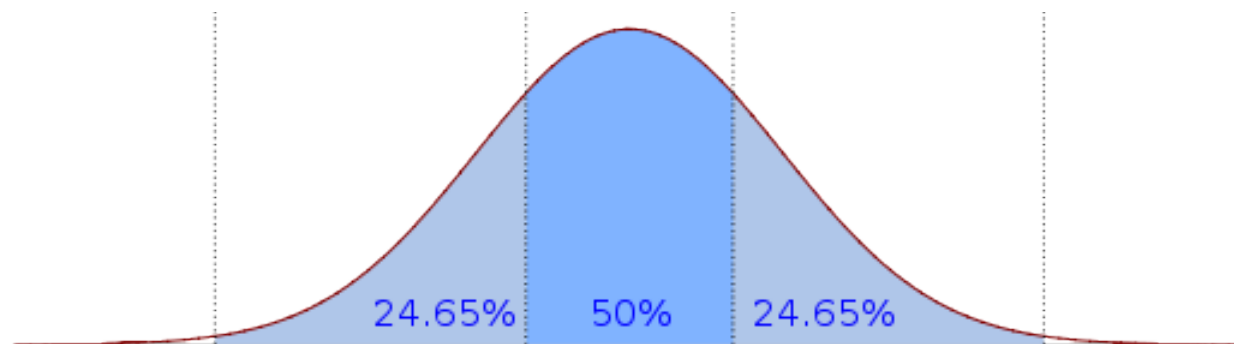
Statistics → quantify the hunch

(Visualizations → storytelling)

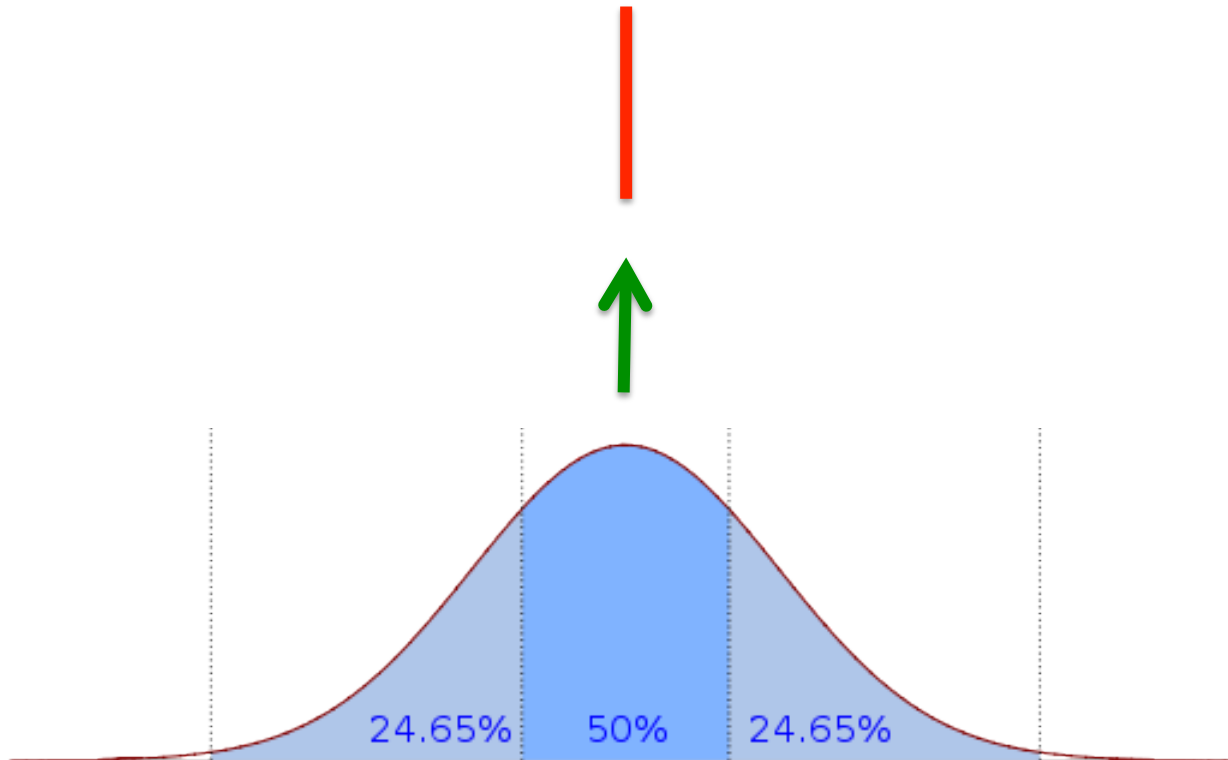
Someone says:
“Obama got more small campaign
contributions than McCain”



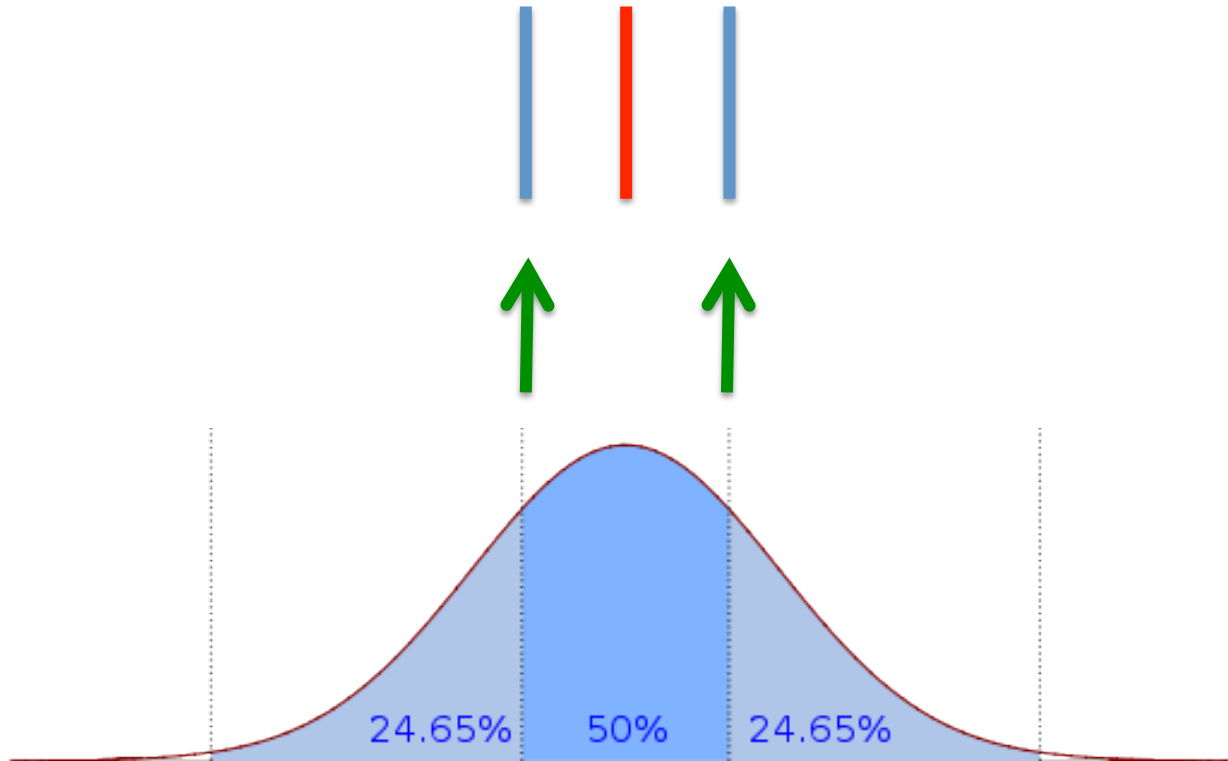




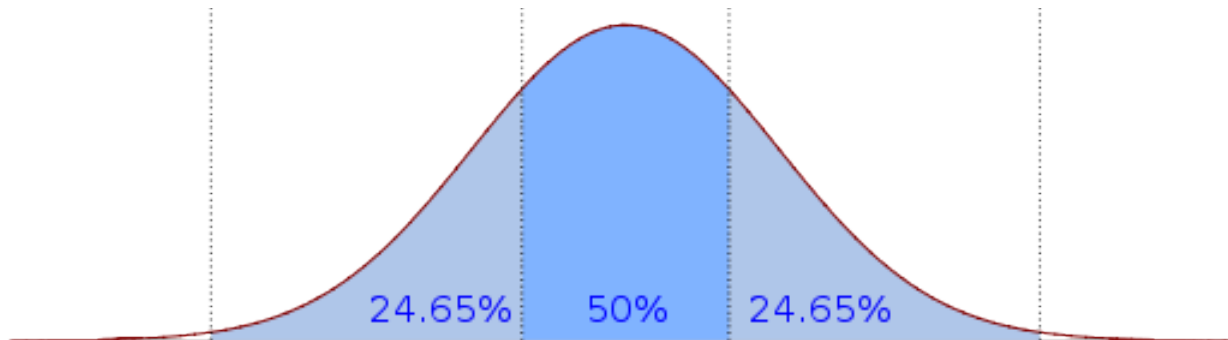
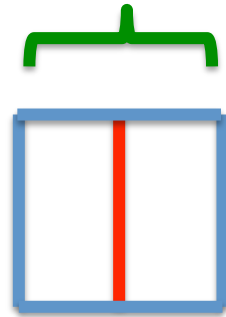
Median



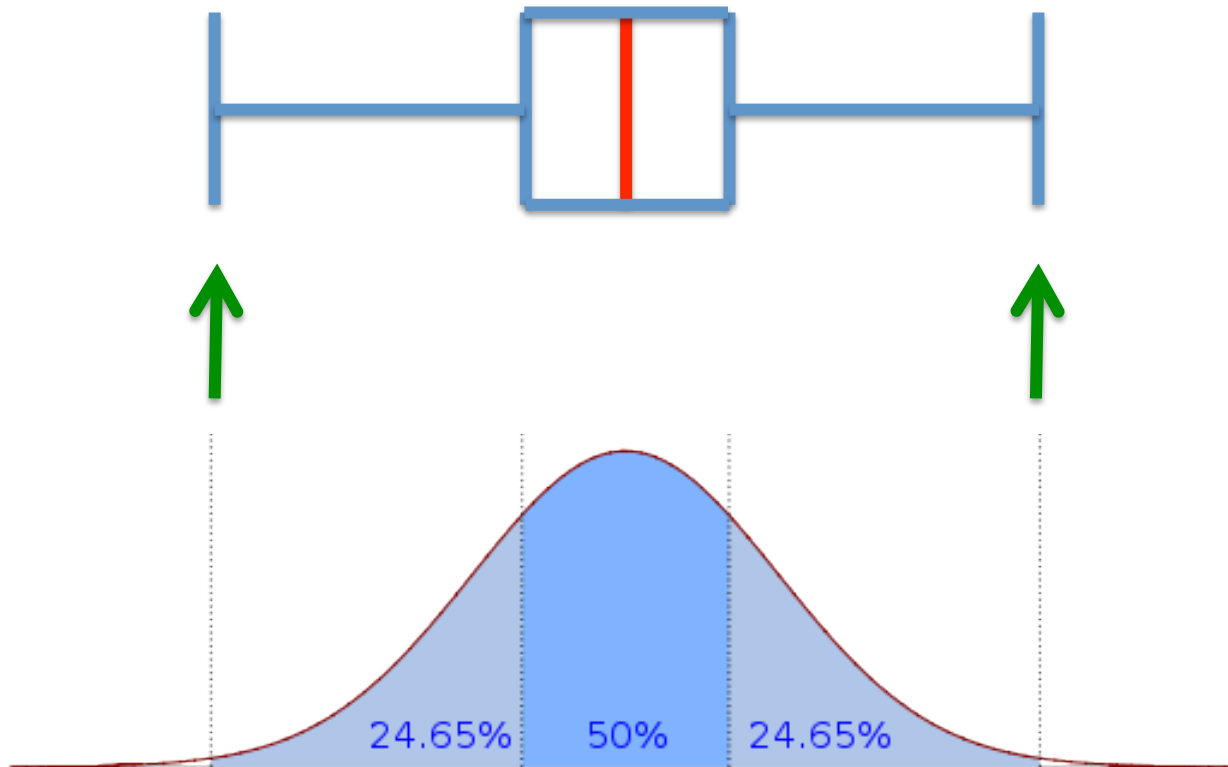
25% 75%



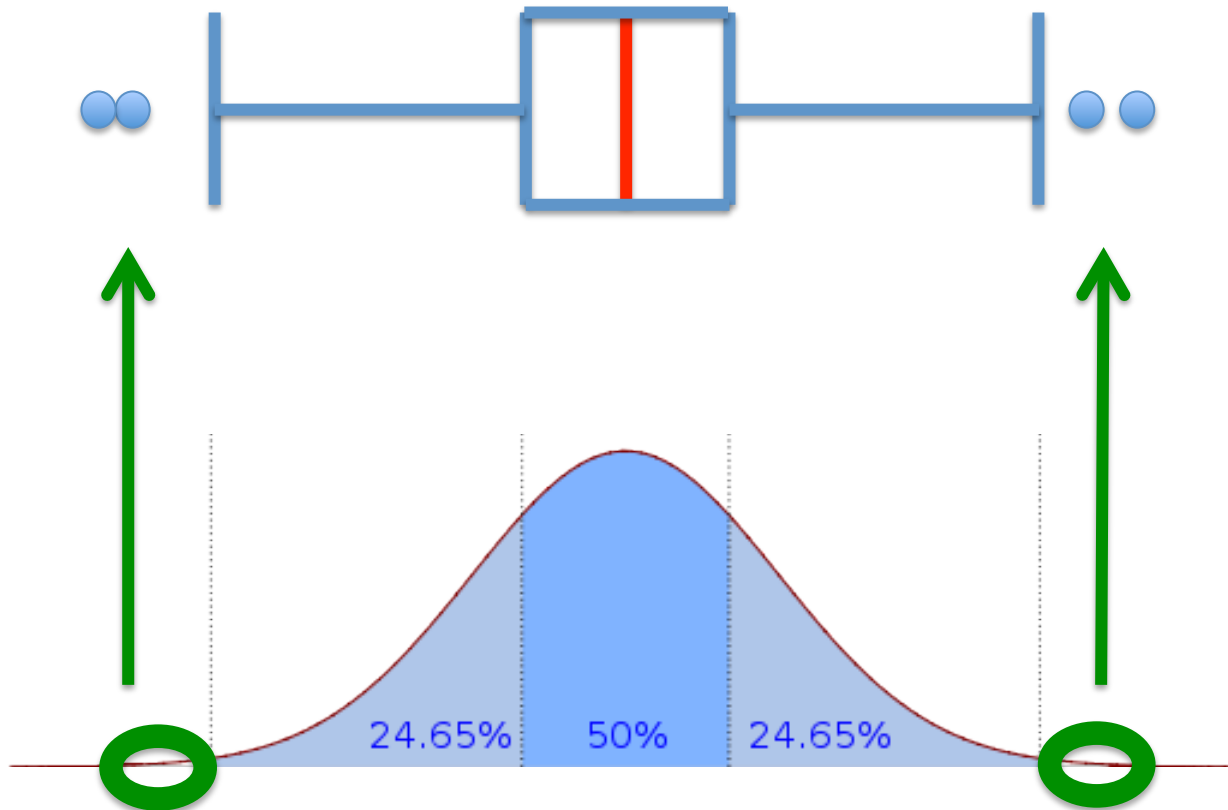
Inner Quartile Range



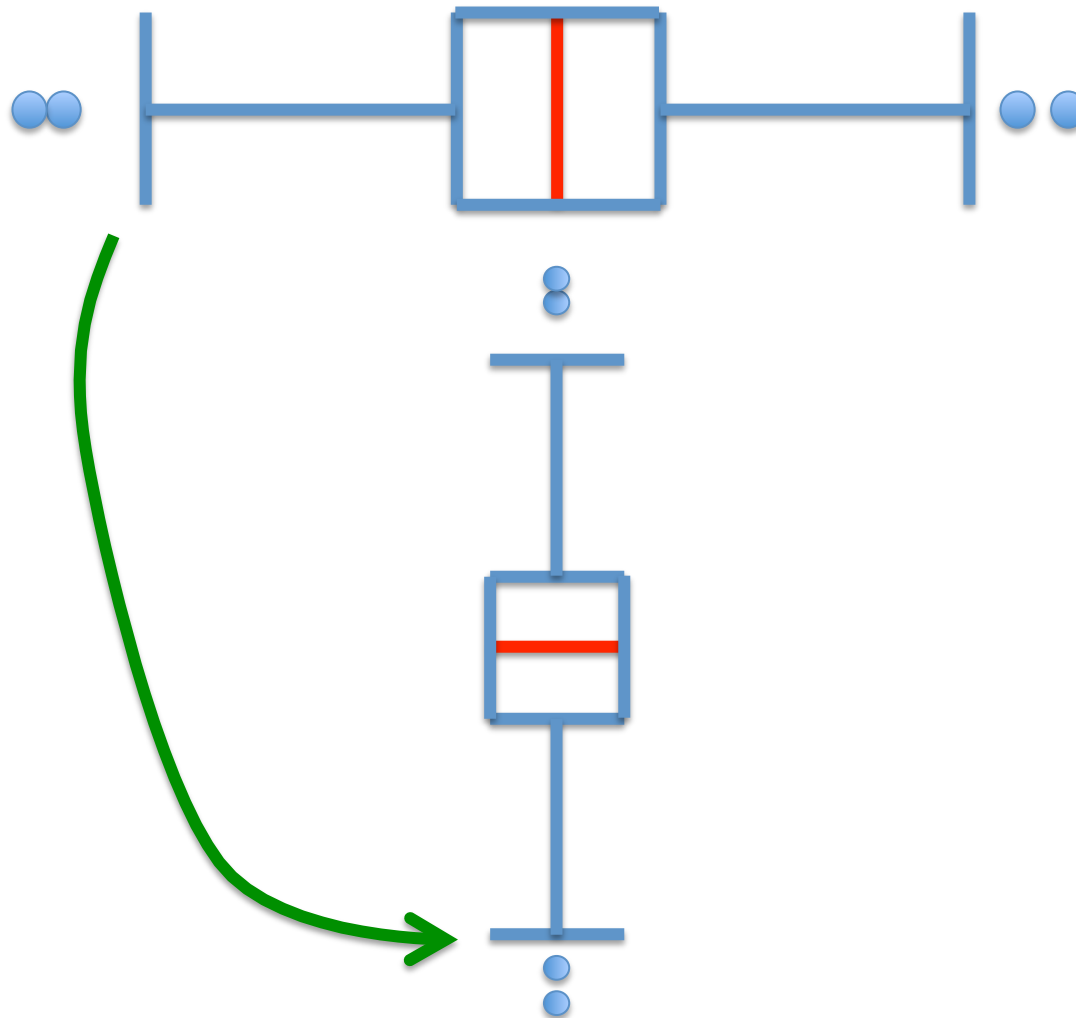
Whiskers / Extremes



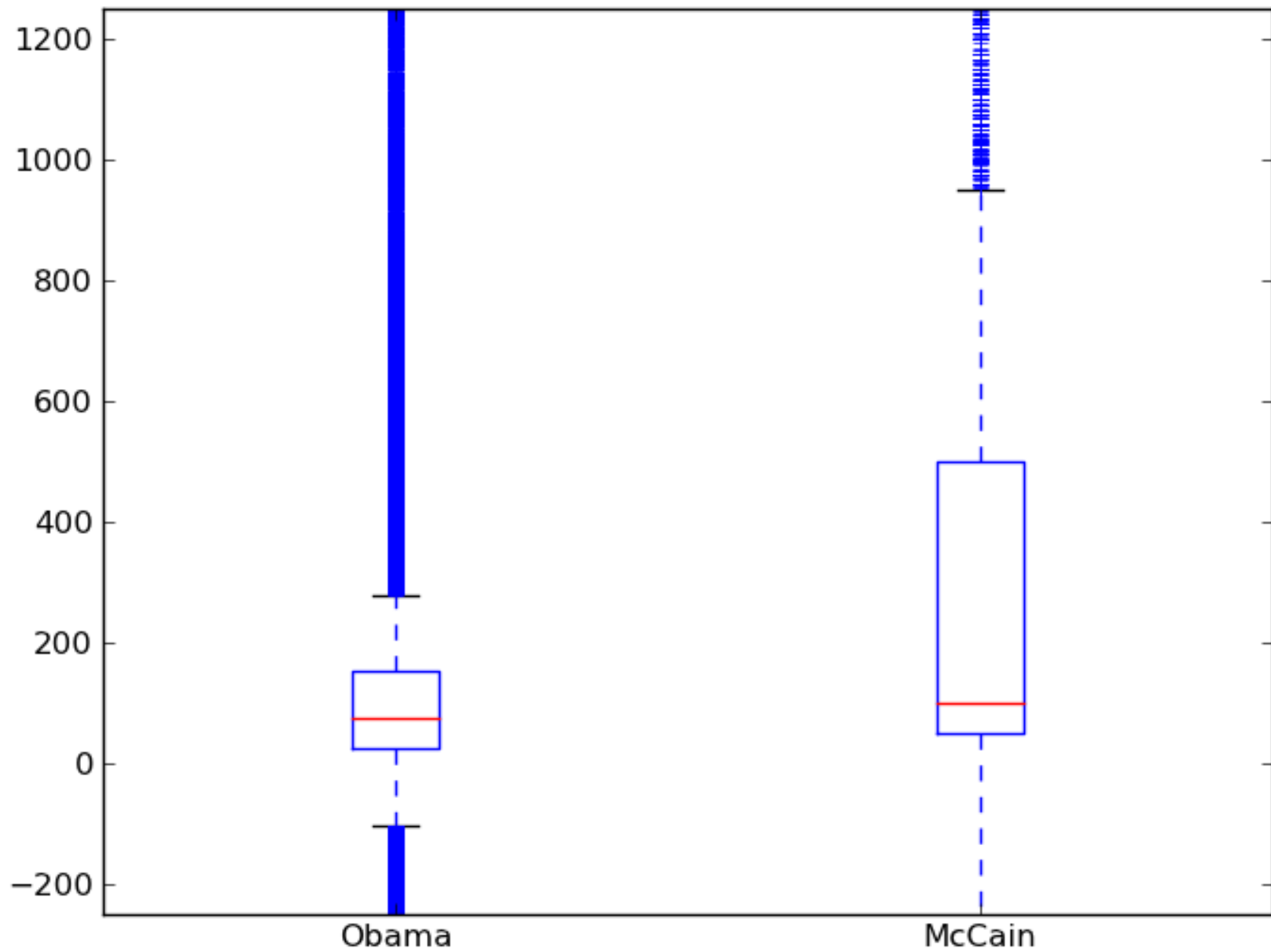
Outliers



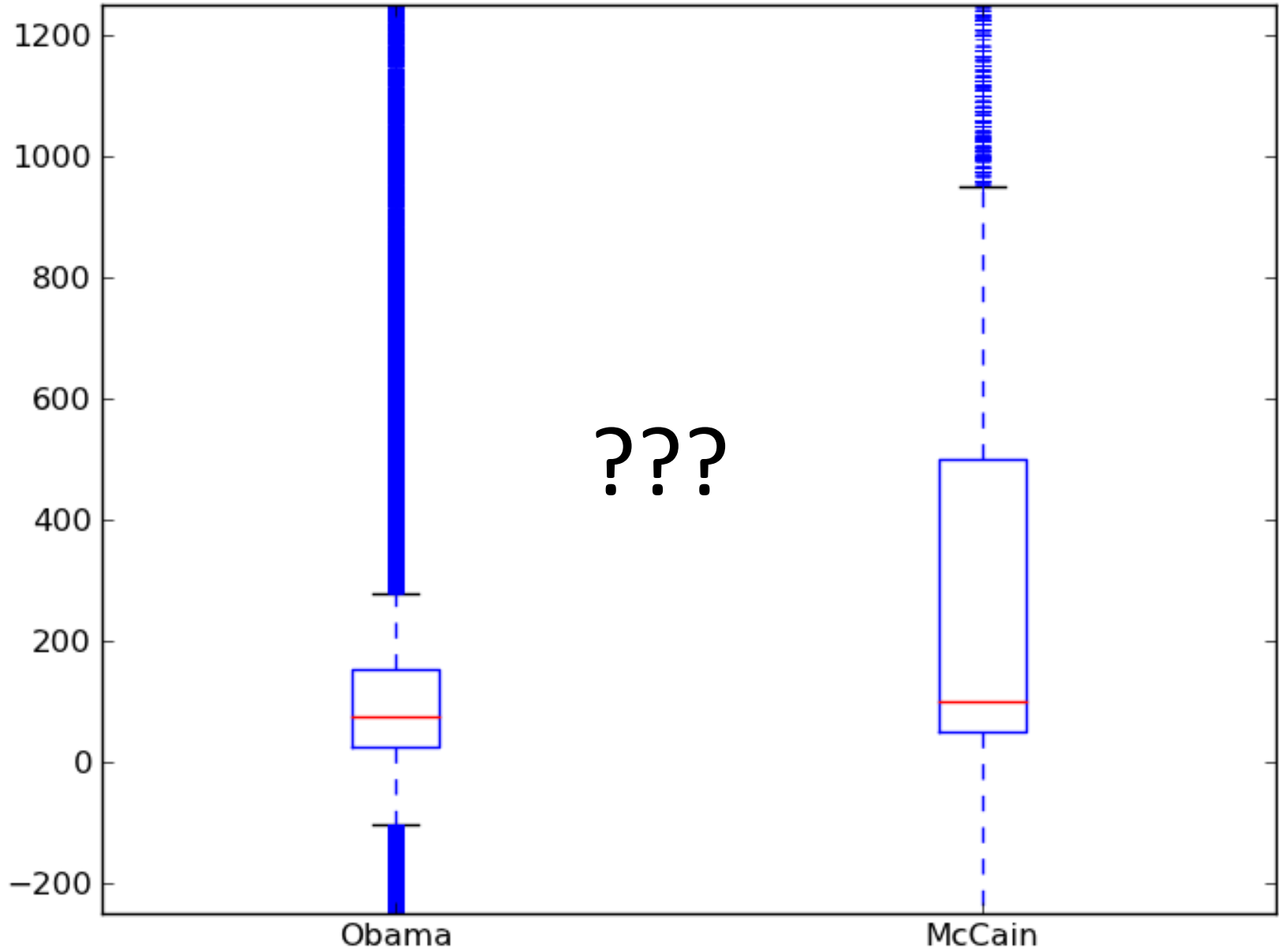
Box-and-Whiskers Plot



Obama vs. McCain Contributions



Obama vs. McCain Contributions

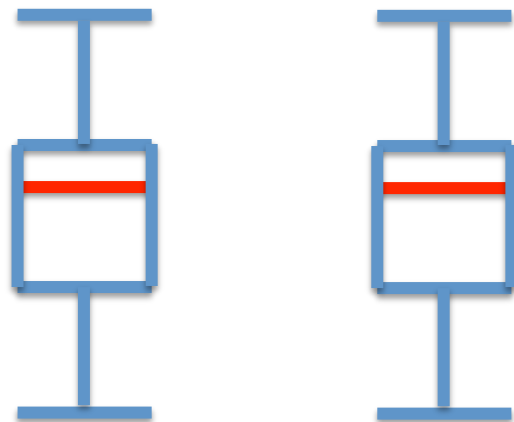


Are they actually different?



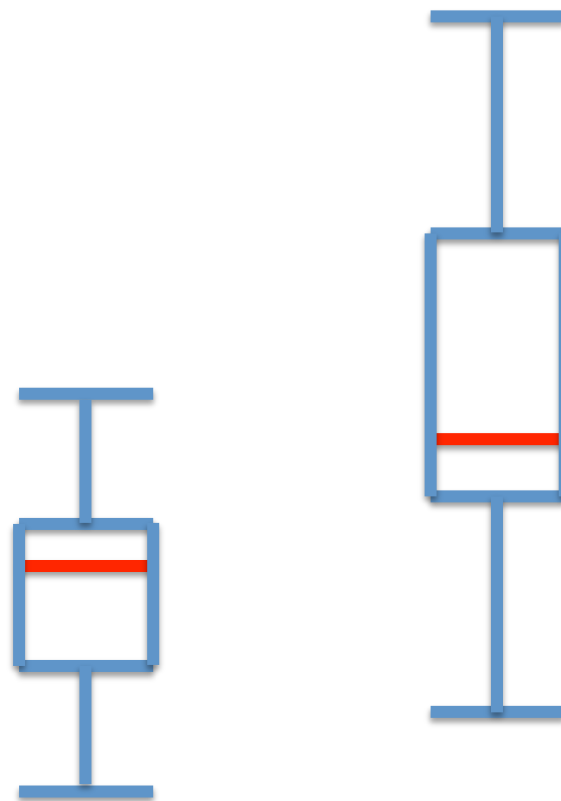
T-Test

Assume



Obama McCain

Reality

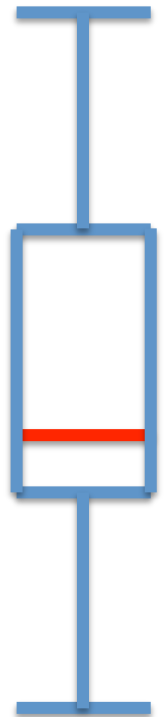


Obama McCain

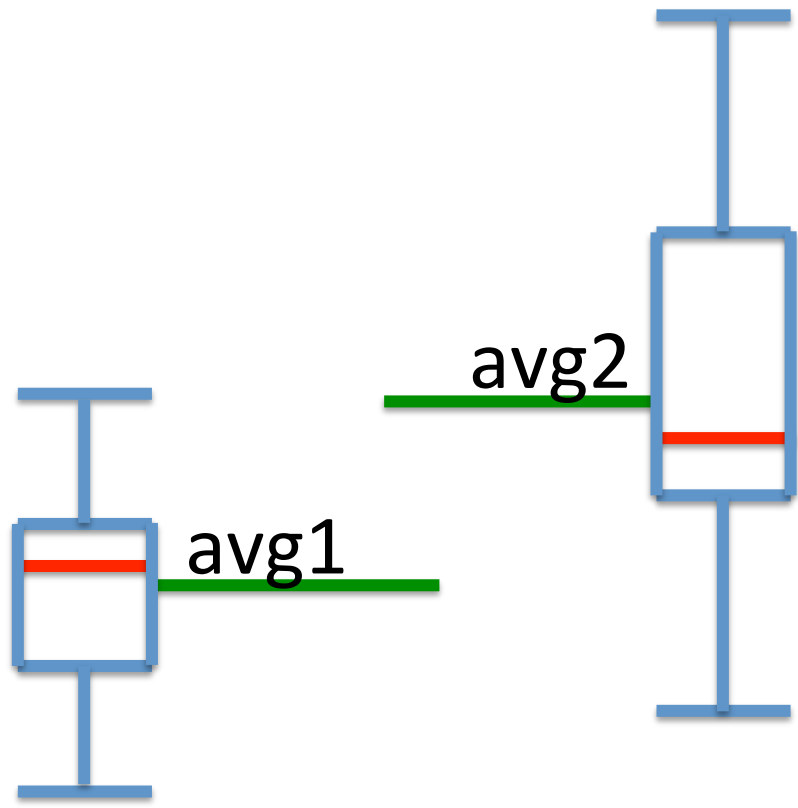
How likely is given ?



Obama McCain

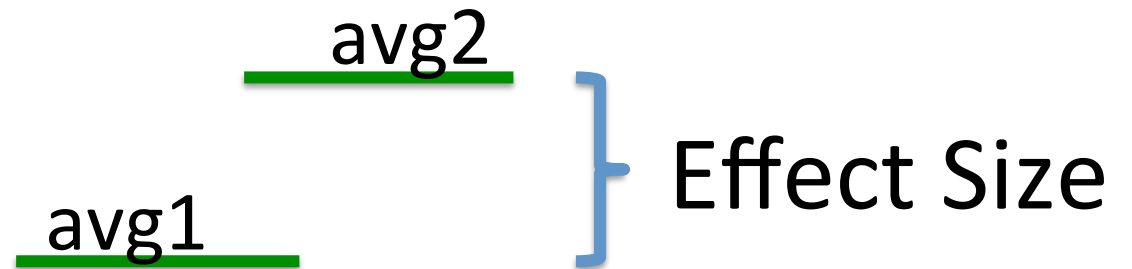


Obama McCain



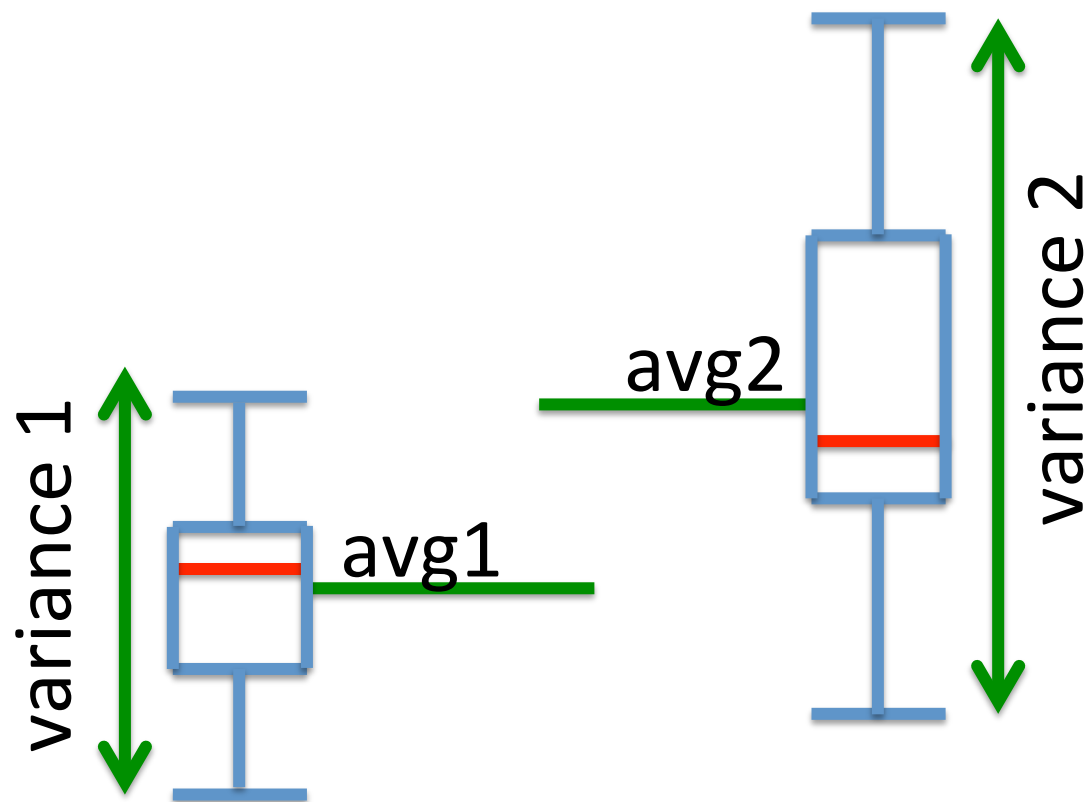
Obama

McCain



Obama

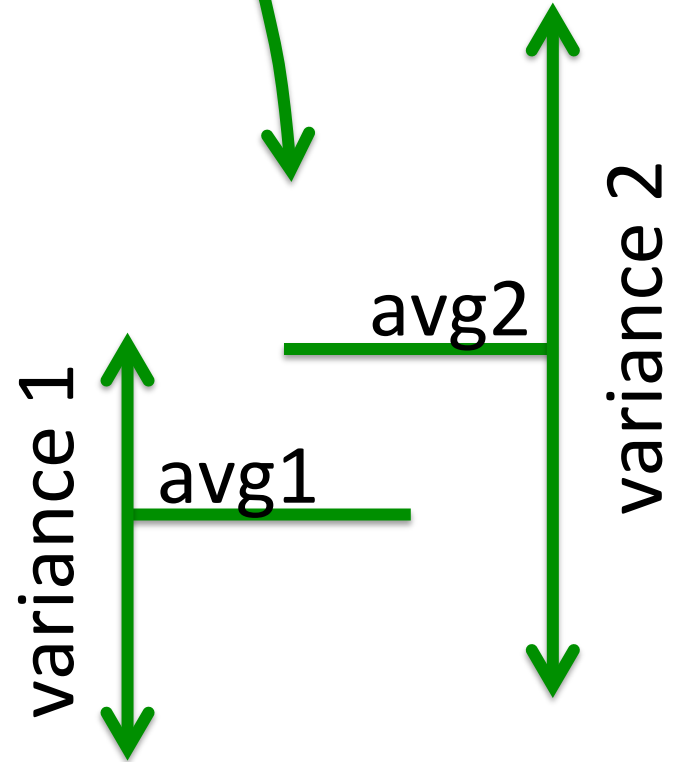
McCain



Obama

McCain

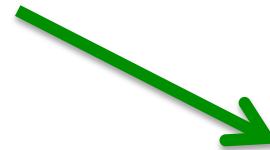
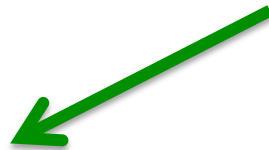
How likely is given ?



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)

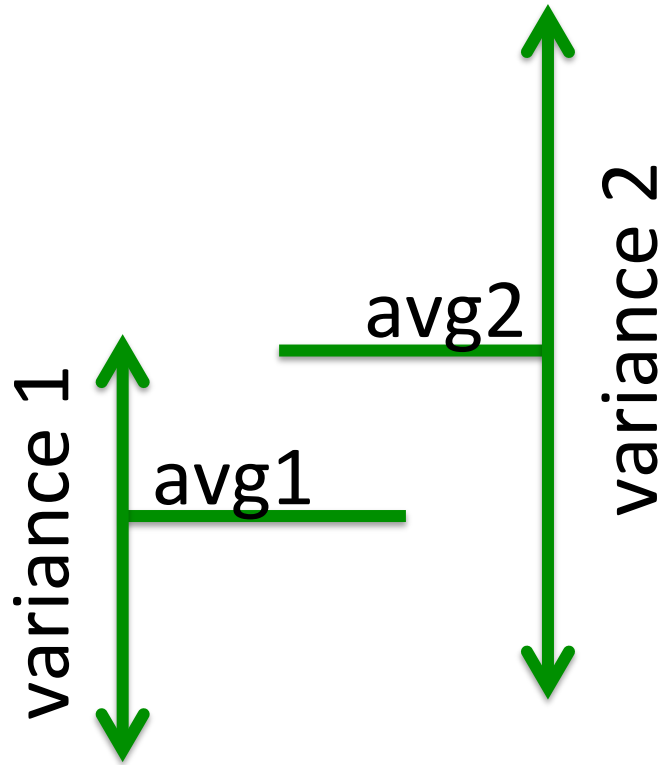
Significance is binary

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

$p \leq .05$? 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

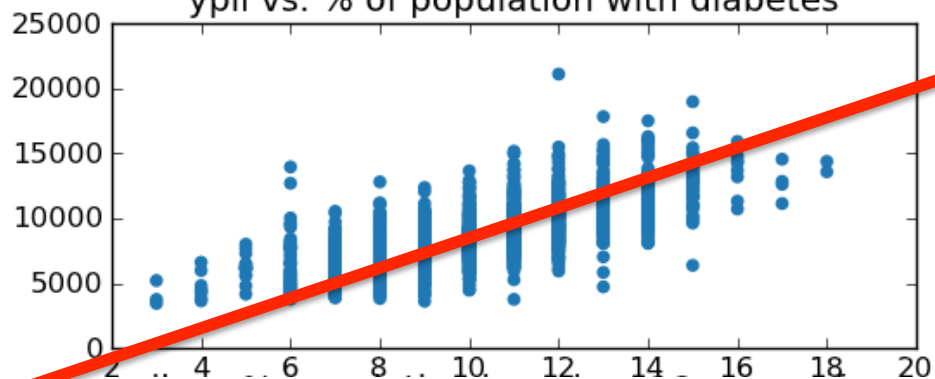
Education

Exercise

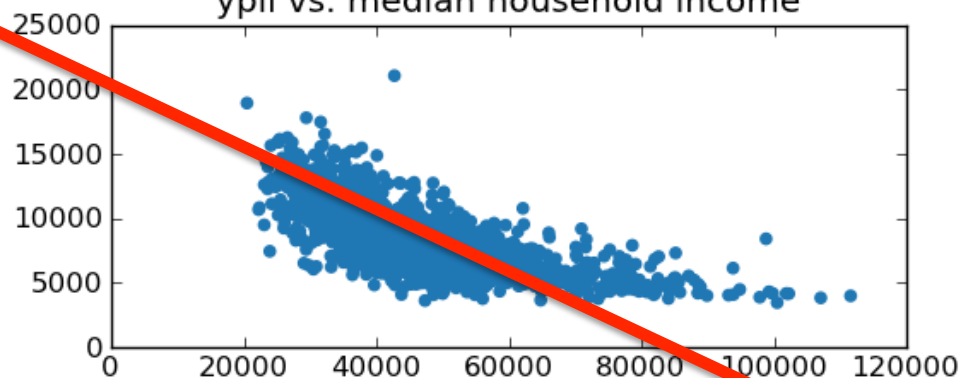
Rappers

Your theory here

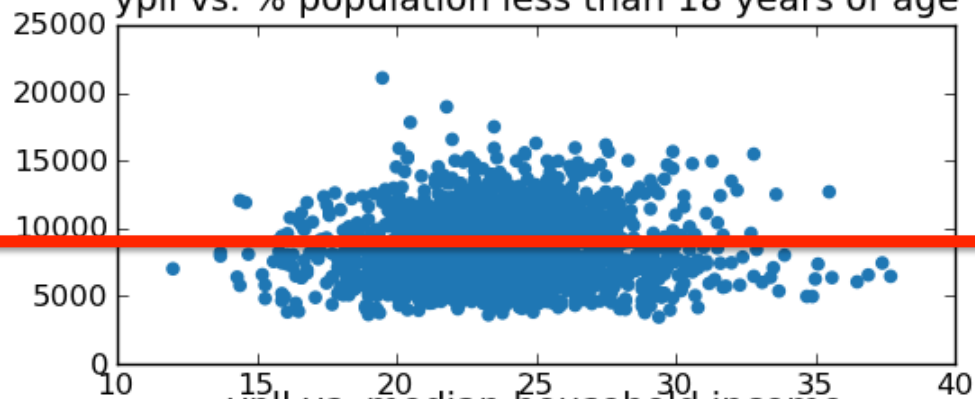
yp11 vs. % of population with diabetes



yp11 vs. median household income



yp11 vs. % population less than 18 years of age

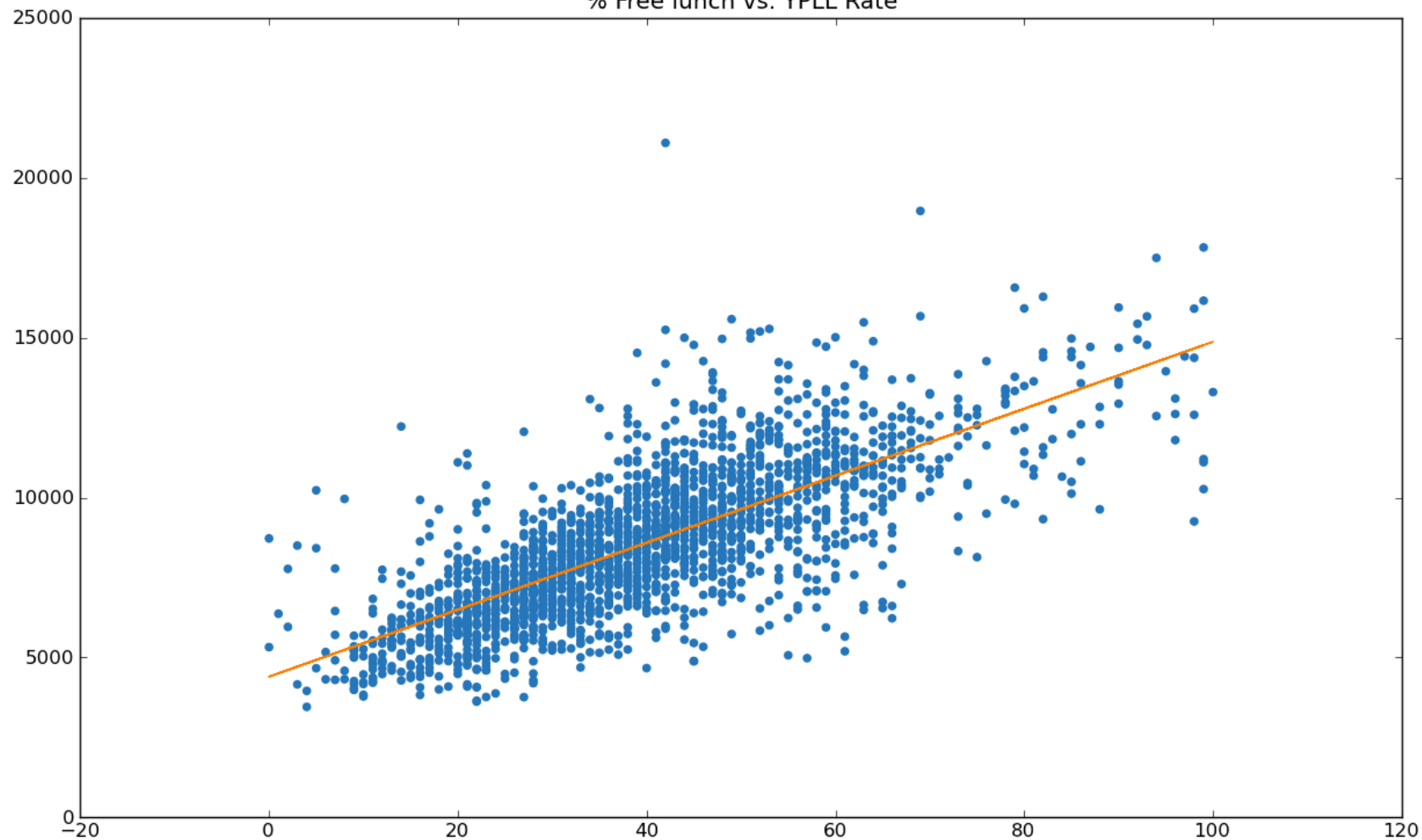


Line coefficients: $y = mx + b$

Correlation amount: R^2 (0 to 1)

Significance: $p < .05?$

% Free lunch vs. YPLL Rate



Decrease amount of free lunch



Reduce early morbidity!



WARNING

**ZOMBIES
AHEAD**



Correlation \neq Causation

Correlation



Causal Hunch



Randomized Trial



T-Test!

Remember to **git pull**

<http://dataiap.github.com/dataiap/day3/>