

What is
Data ?
Why is it **Science** ?

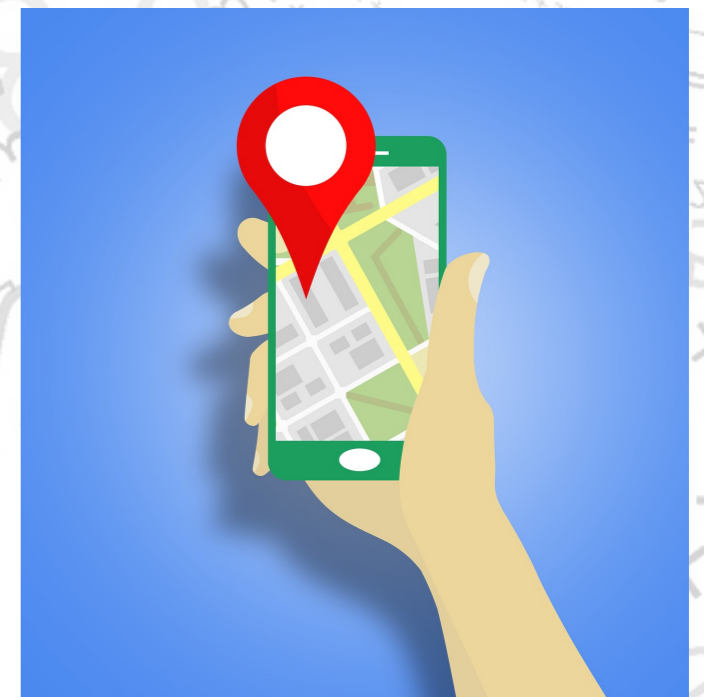
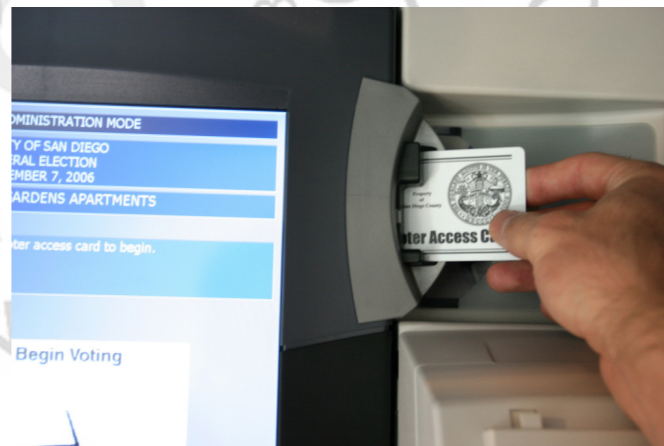
Data is everywhere...



ebay
VISA



Booking.com
amazon



Meaning from data is money...

Harvard
Business
Review

DATA

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE

DSPT
DATA SCIENCE PORTUGAL

 **DSPPA**
DATA SCIENCE
PORTUGUESE ASSOCIATION

"Data Scientist"

Percentage of matching job postings



Source: Indeed

indeed

Challenges...

DATA



SORTED



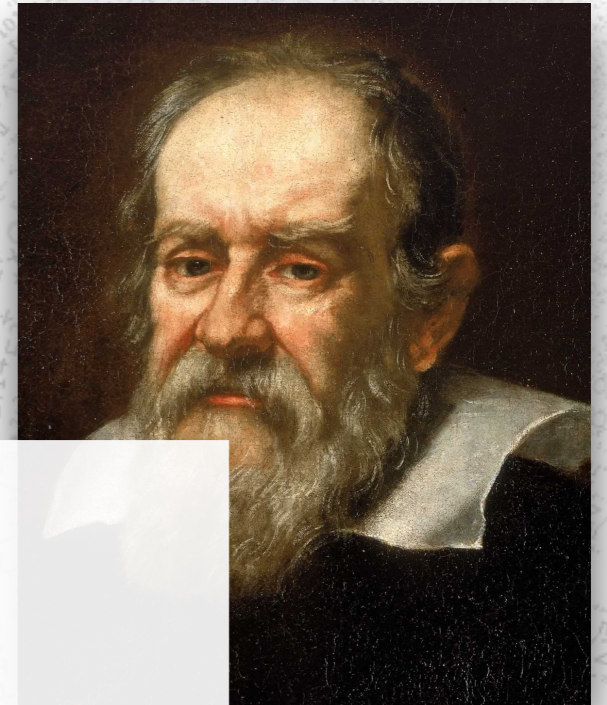
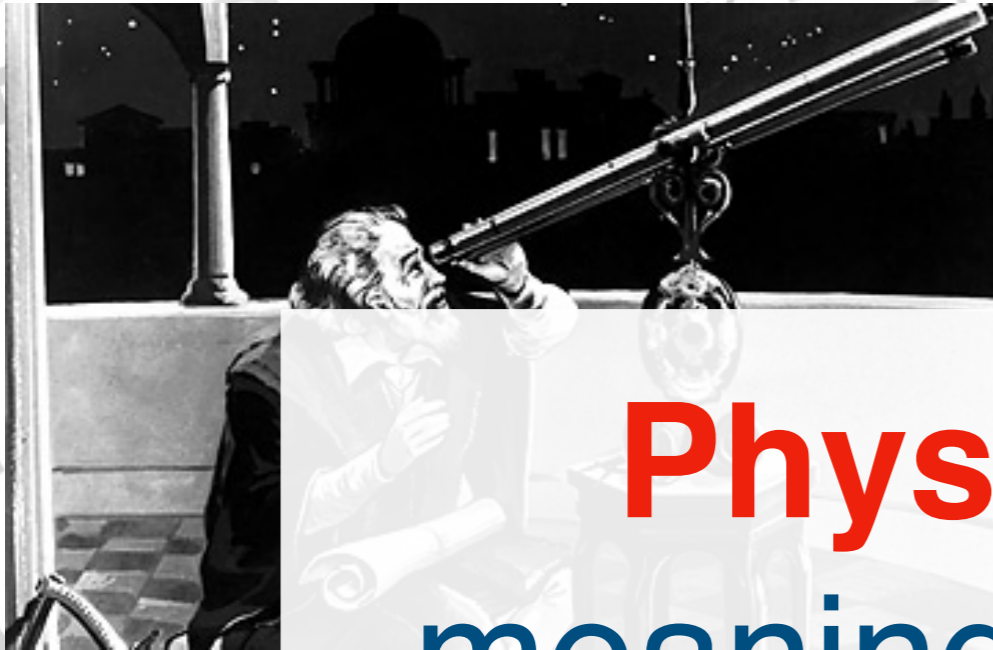
ARRANGED



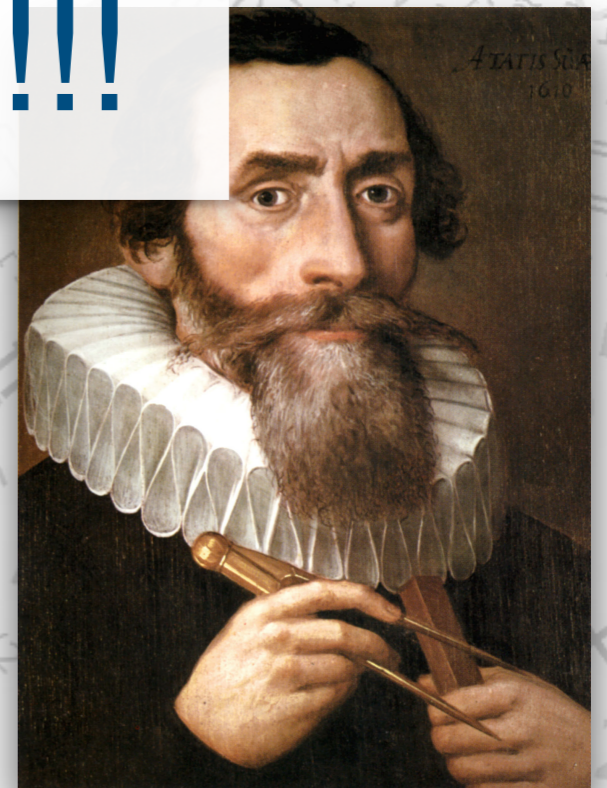
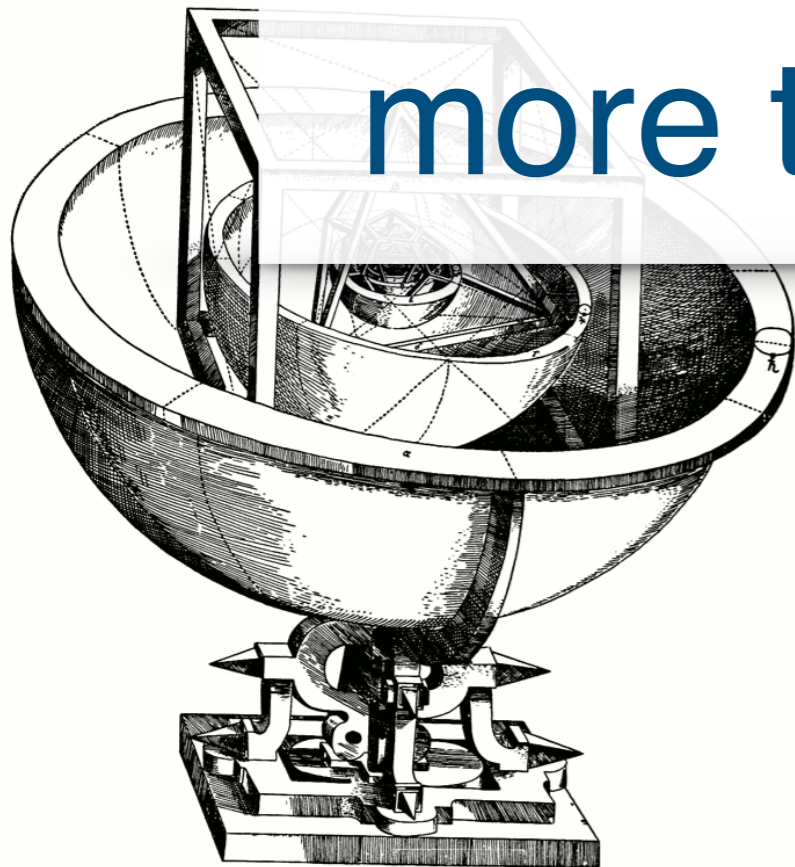
PRESENTED VISUALLY



Why **us**?

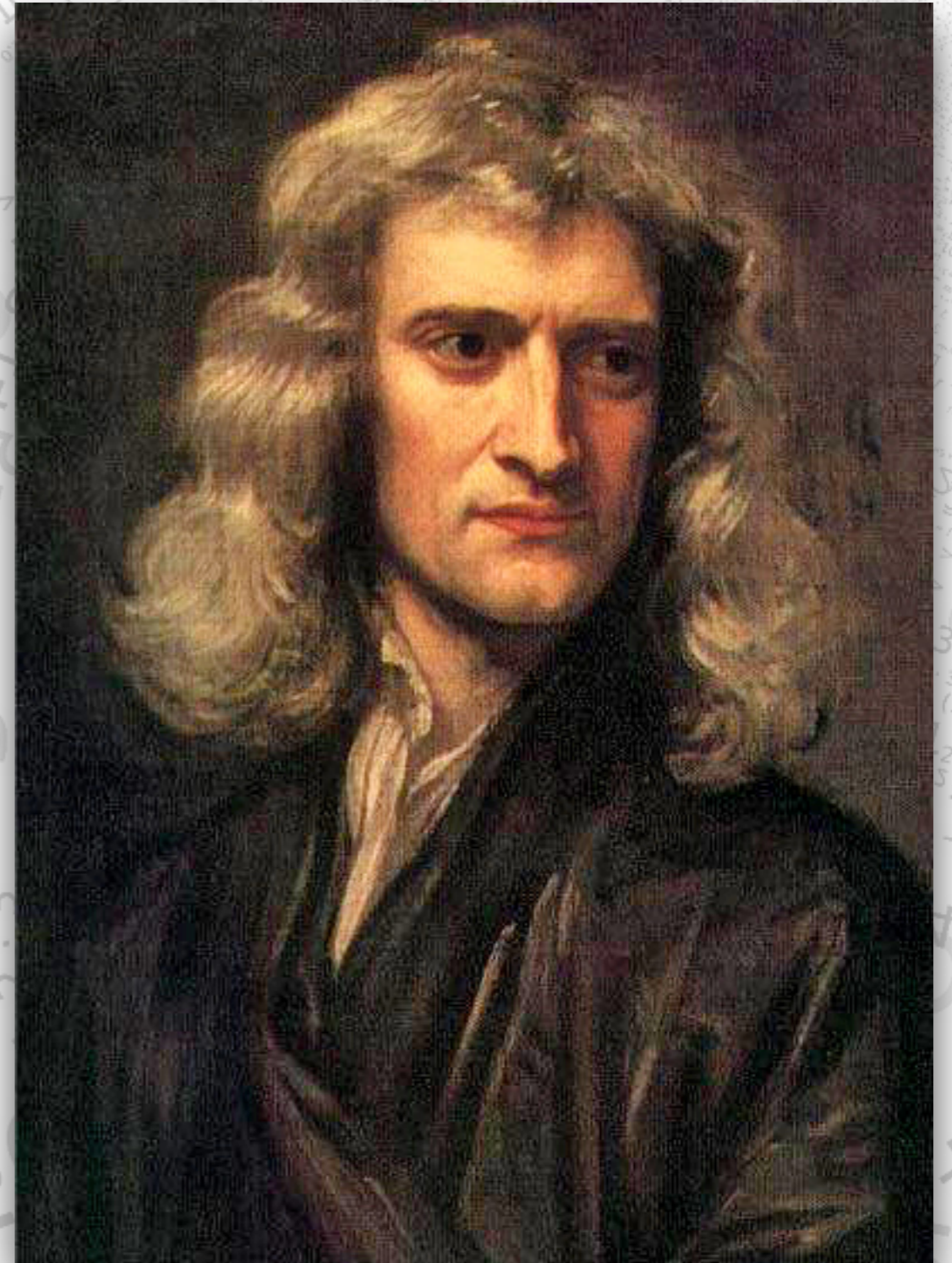


Physics: Taking
meaning from data for
more than **500 years!!!**



Johannes Kepler
1571-1630

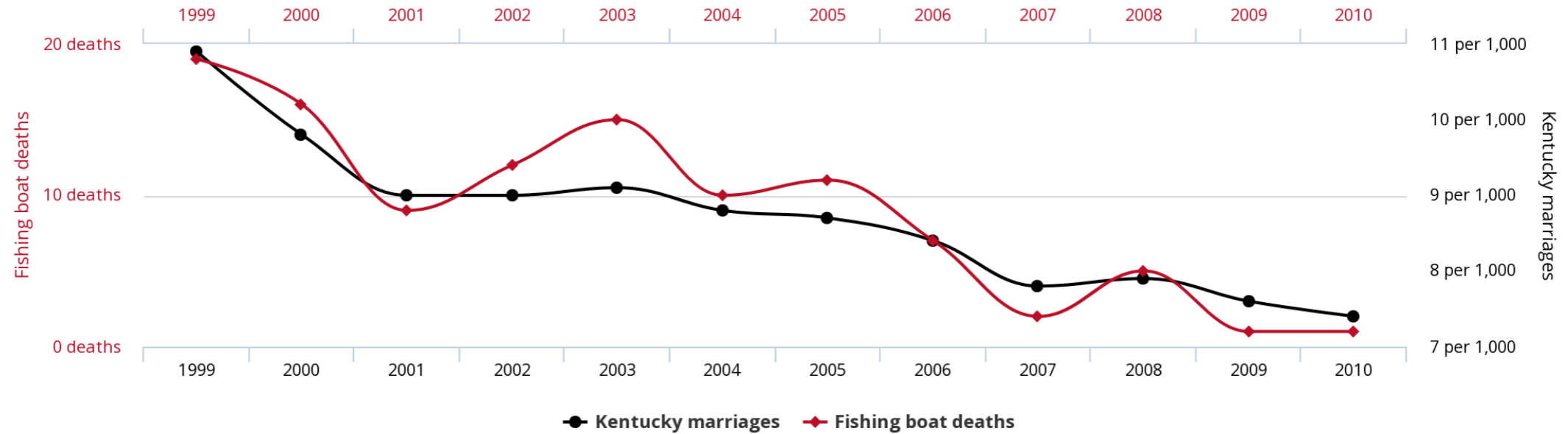
The power of the model...



Isaac Newton
1643-1727

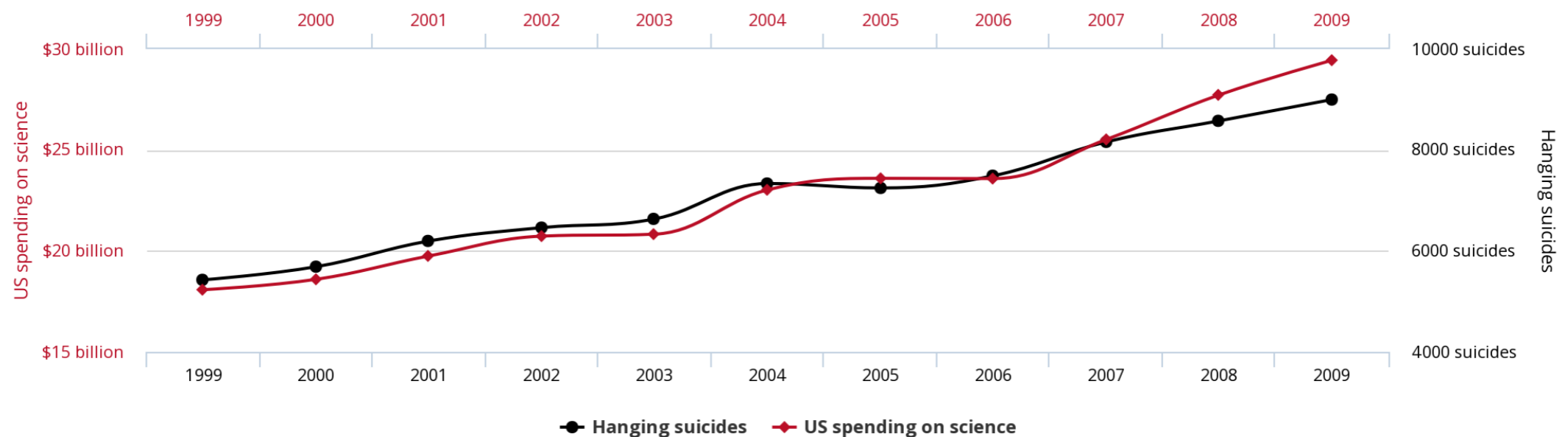
Causality

People who drowned after falling out of a fishing boat
correlates with
Marriage rate in Kentucky



tylervigen.com

US spending on science, space, and technology
correlates with
Suicides by hanging, strangulation and suffocation



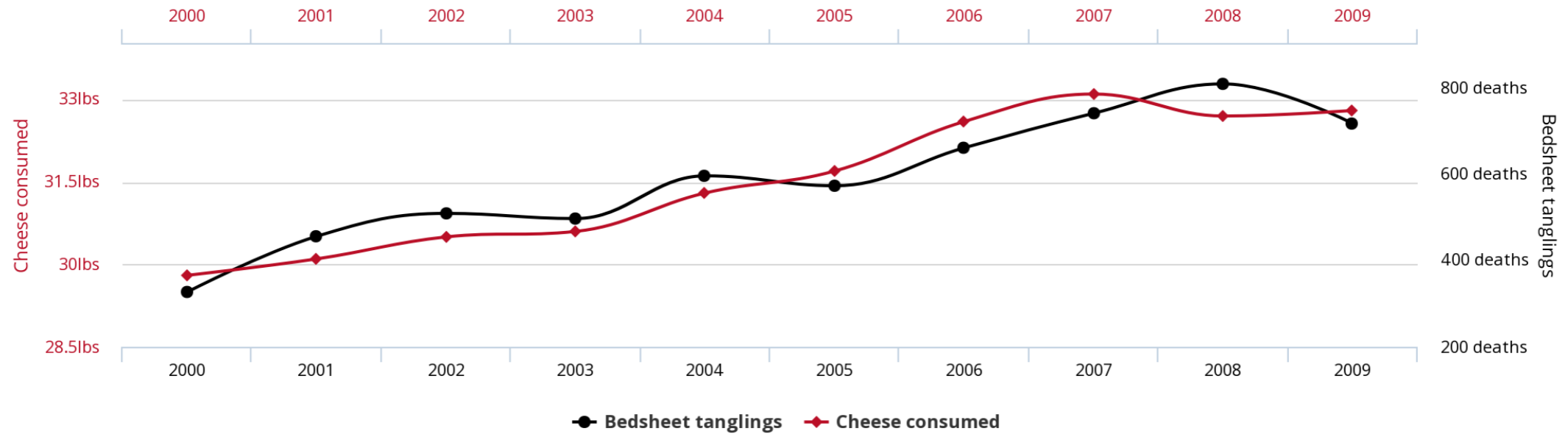
tylervigen.com

Causality

Per capita cheese consumption

correlates with

Number of people who died by becoming tangled in their bedsheets

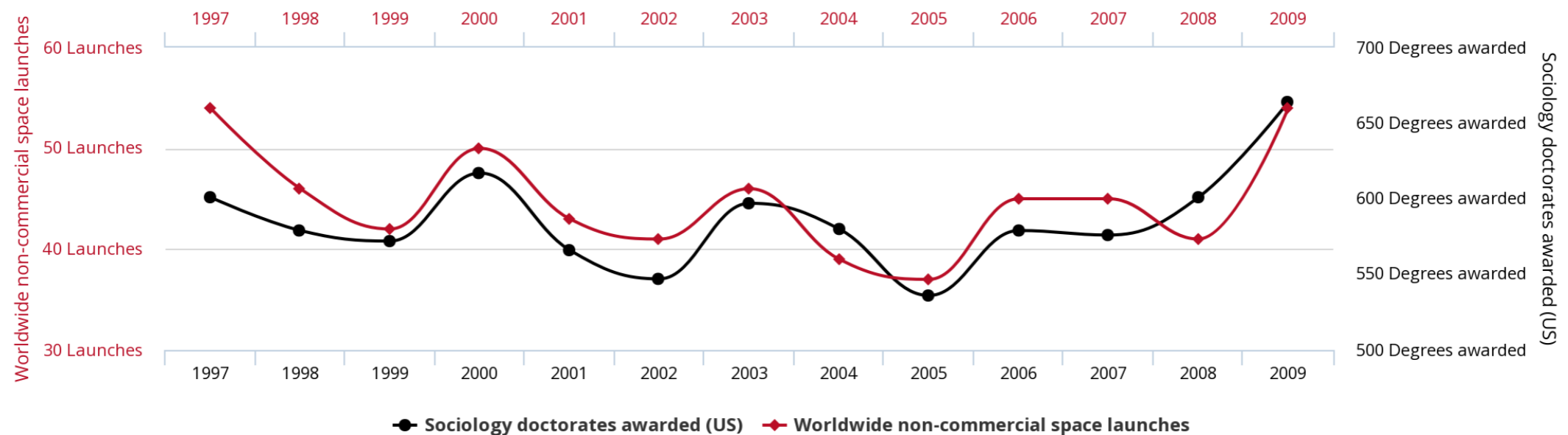


tylervigen.com

Worldwide non-commercial space launches

correlates with

Sociology doctorates awarded (US)

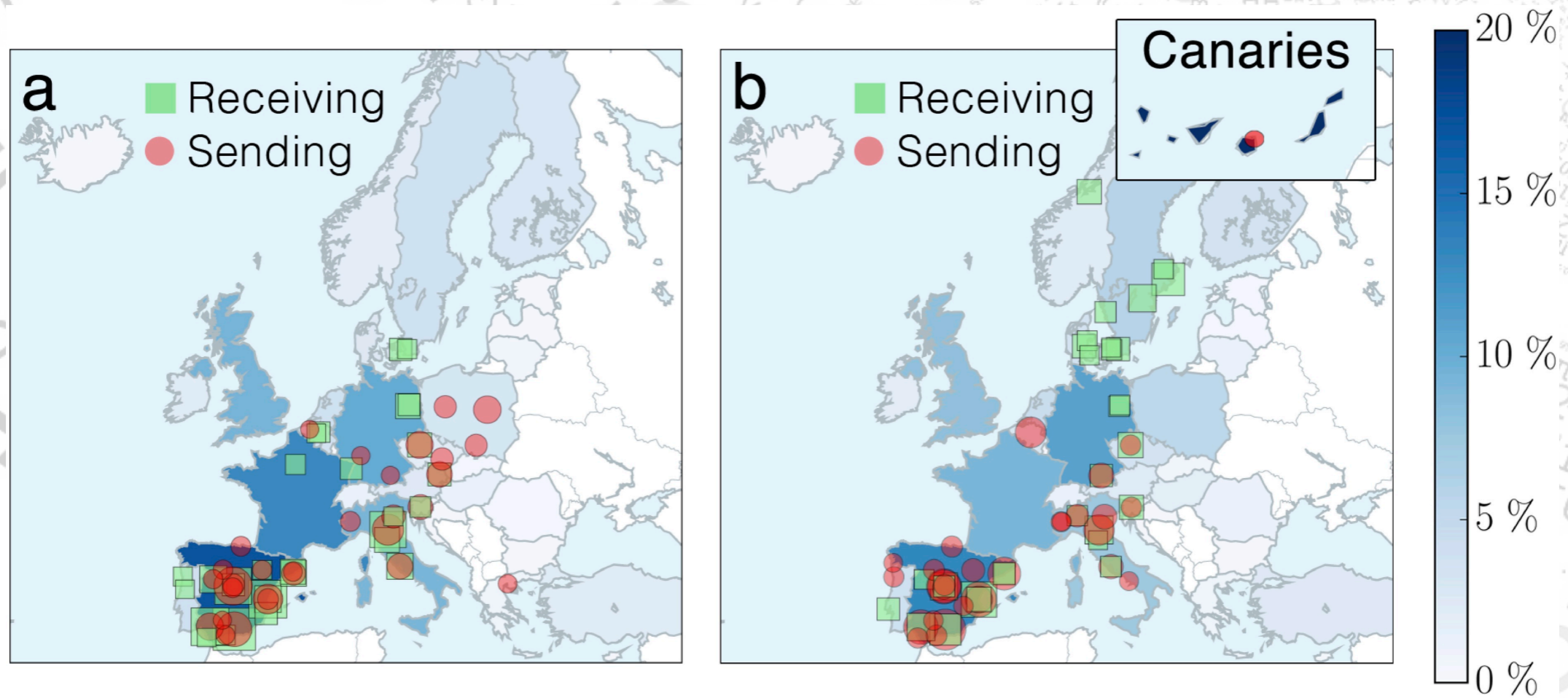


tylervigen.com

Causality: Gender bias in Erasmus

Female

Male

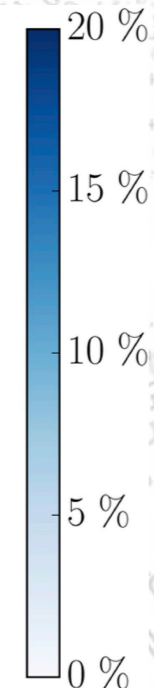
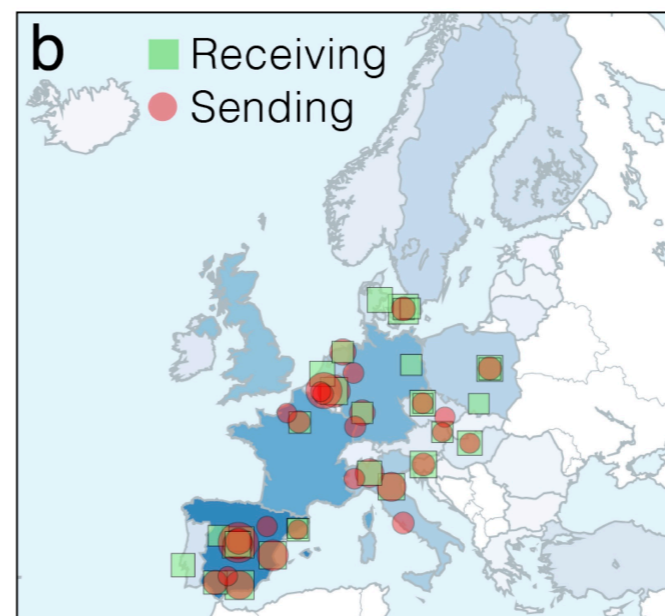
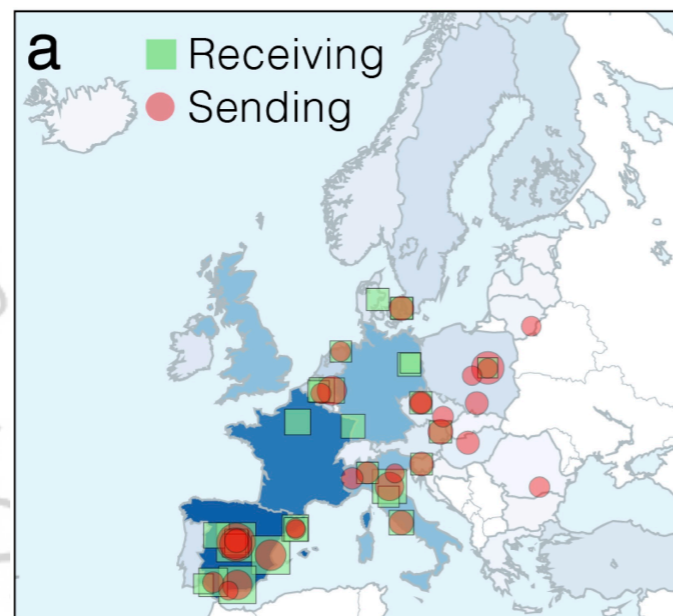


Causality: Gender bias in Erasmus

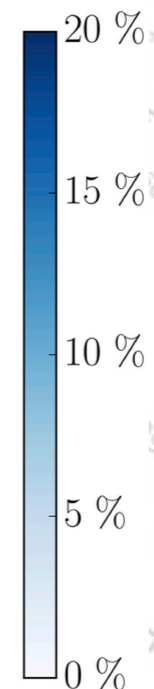
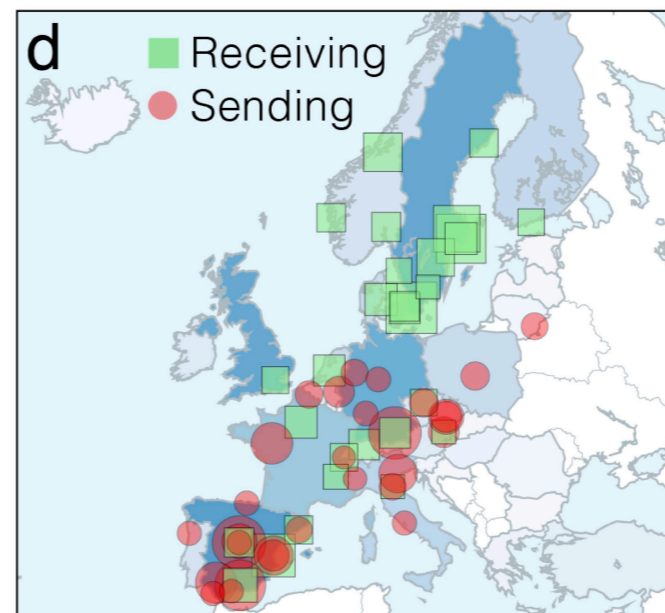
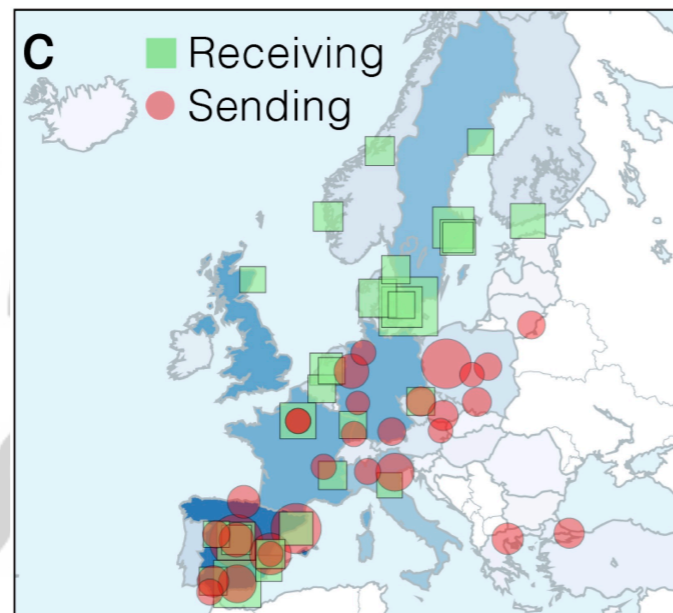
Female

Male

*Social
Sciences*



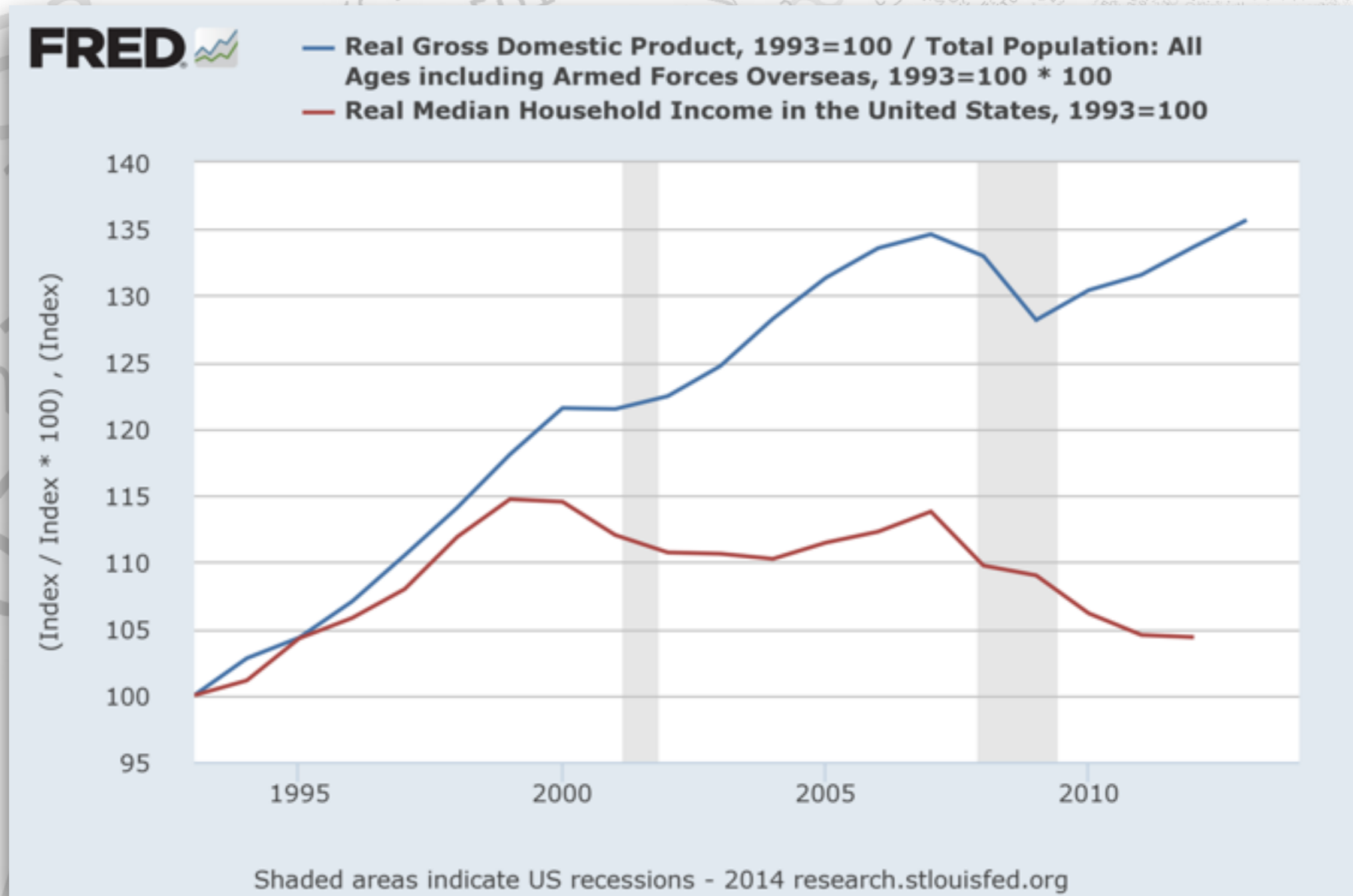
*Natural
Sciences*



Is the **average** meaningful?

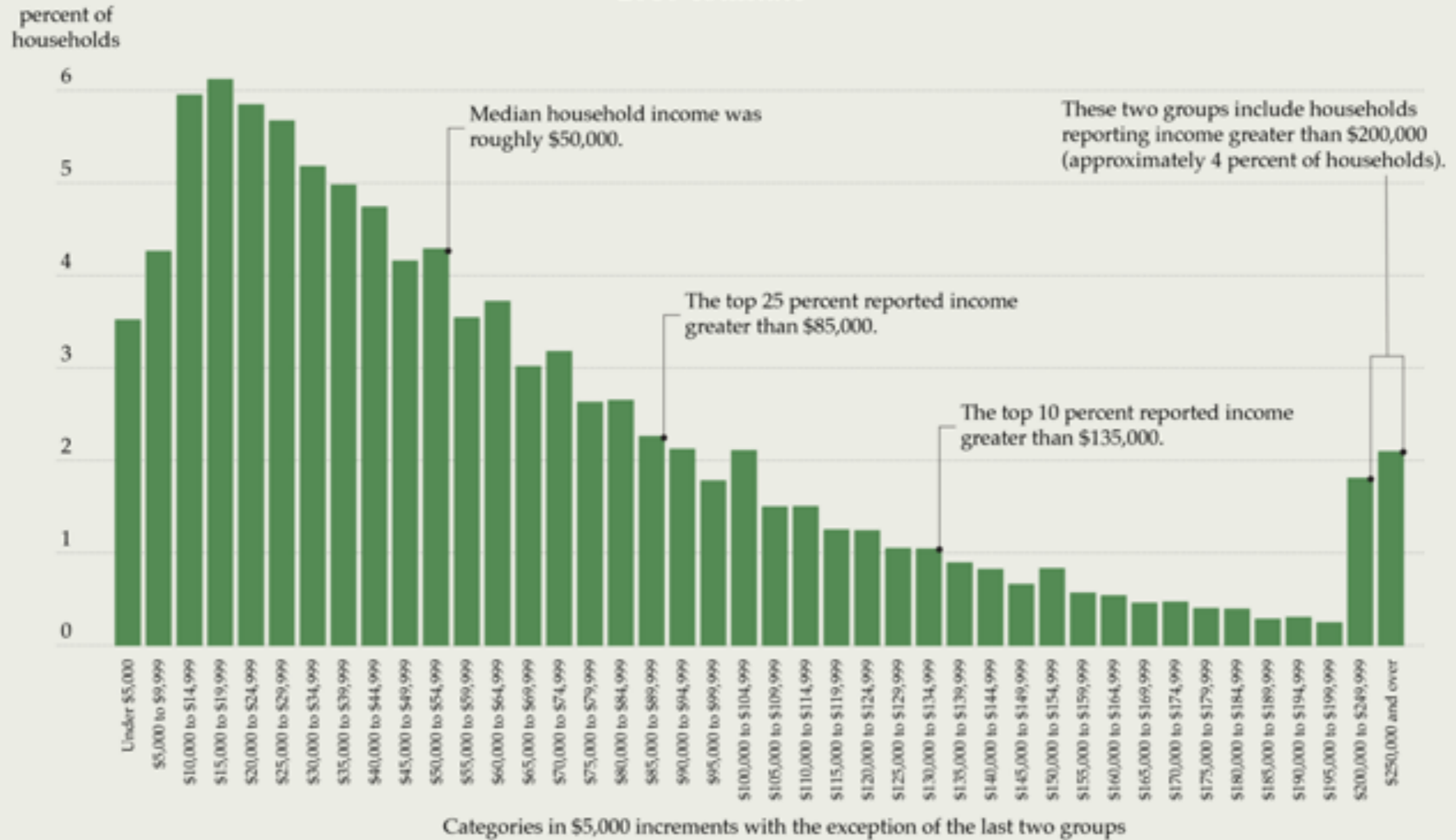


Average vs Median



Histograms

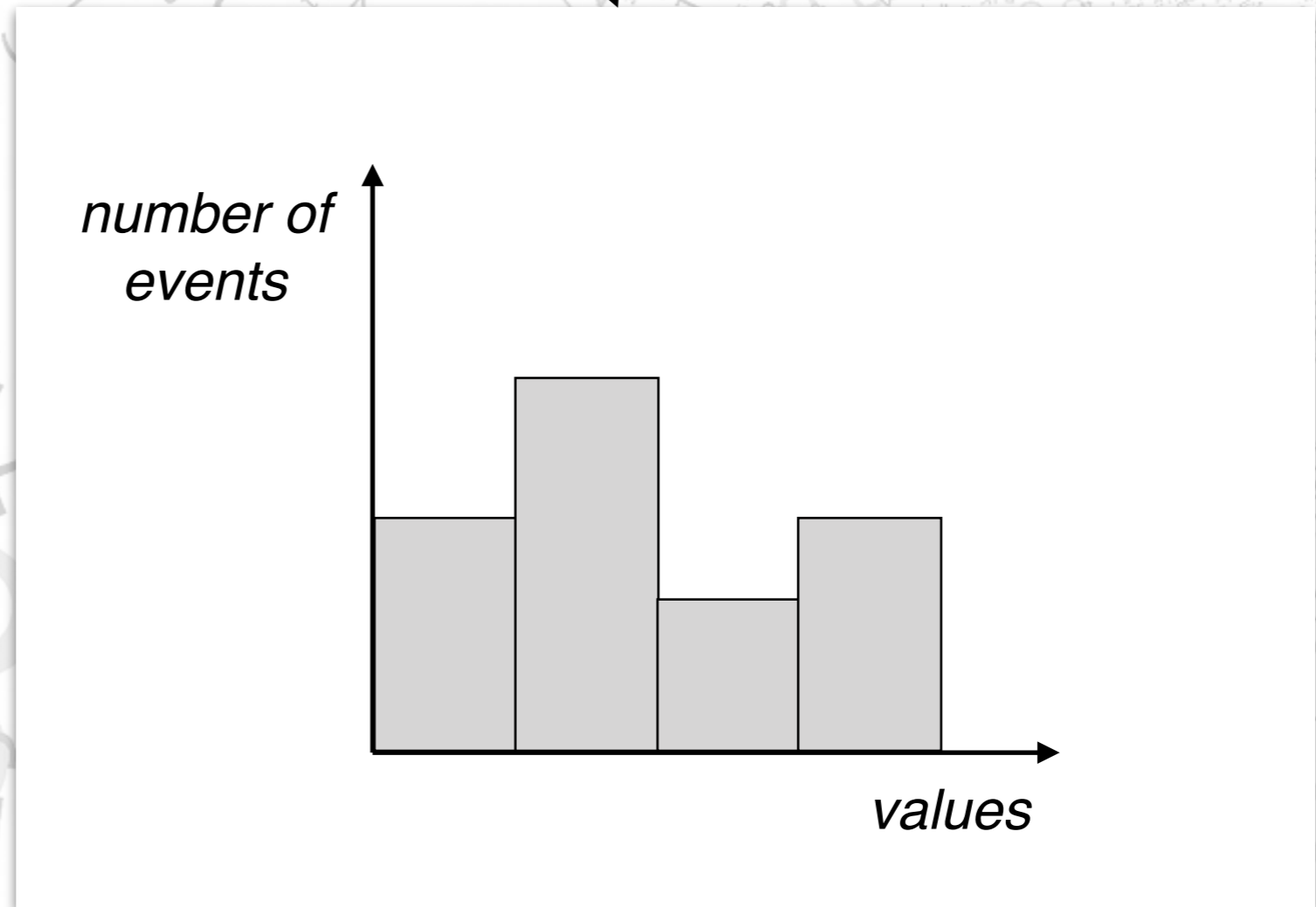
Distribution of annual household income in the United States 2010 estimate



Source: U.S. Census Bureau, Current Population Survey, 2011 Annual Social and Economic Supplement

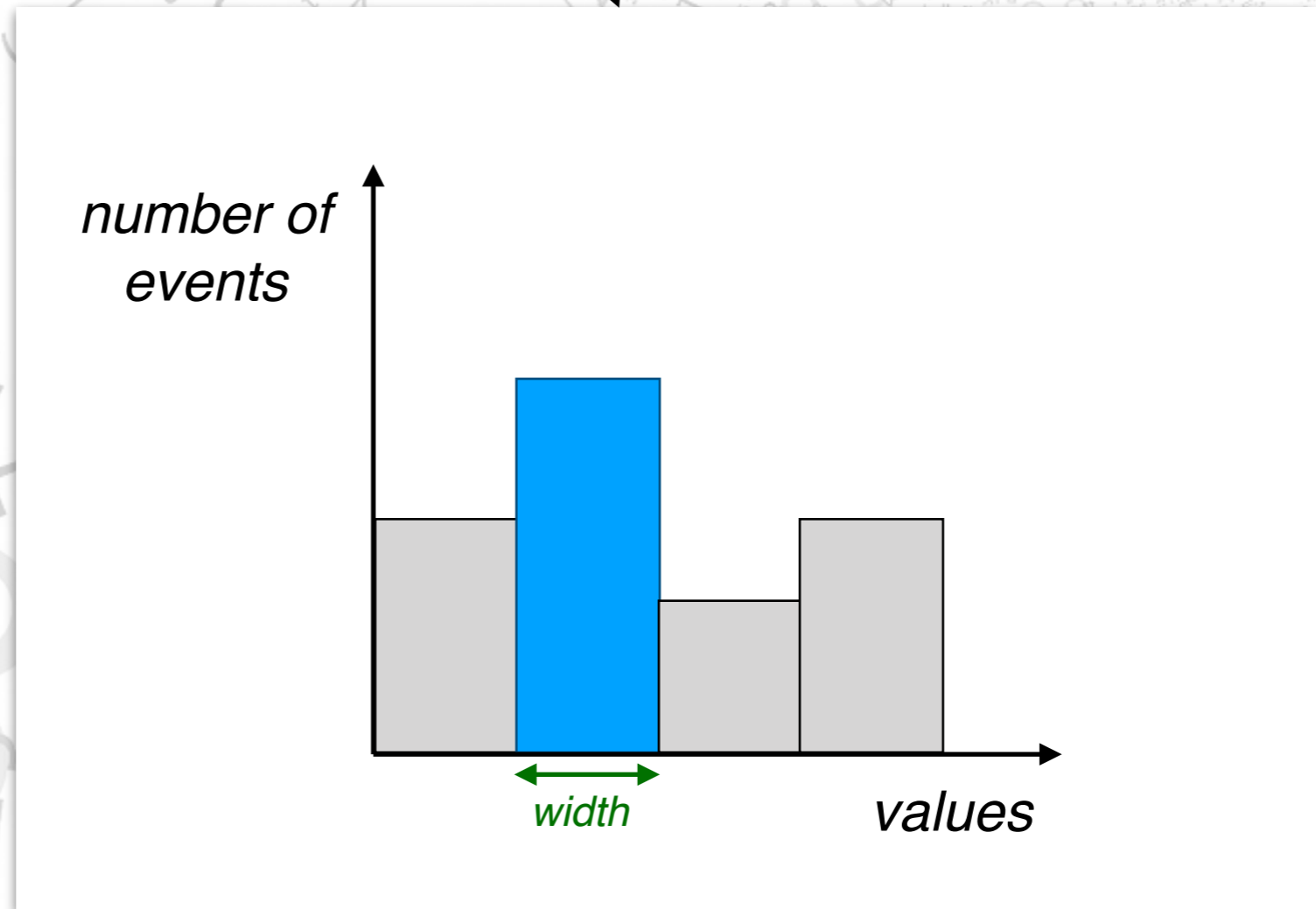
How to make a histograms...

6.39
4.62
3.13
4.96
0.09
0.19
8.88
0.59
1.74
0.24
7.58
...



How to make a histograms...

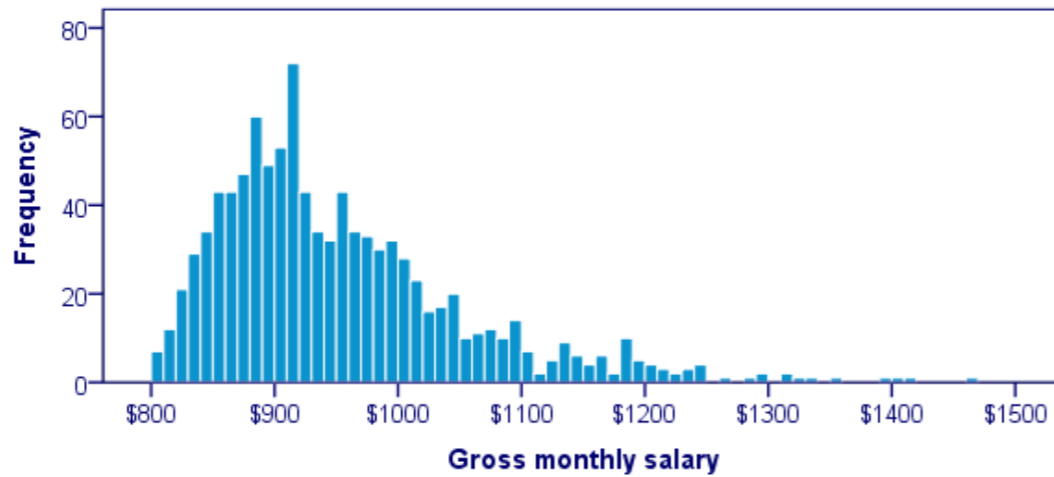
6.39
4.62
3.13
4.96
0.09
0.19
8.88
0.59
1.74
0.24
7.58
...



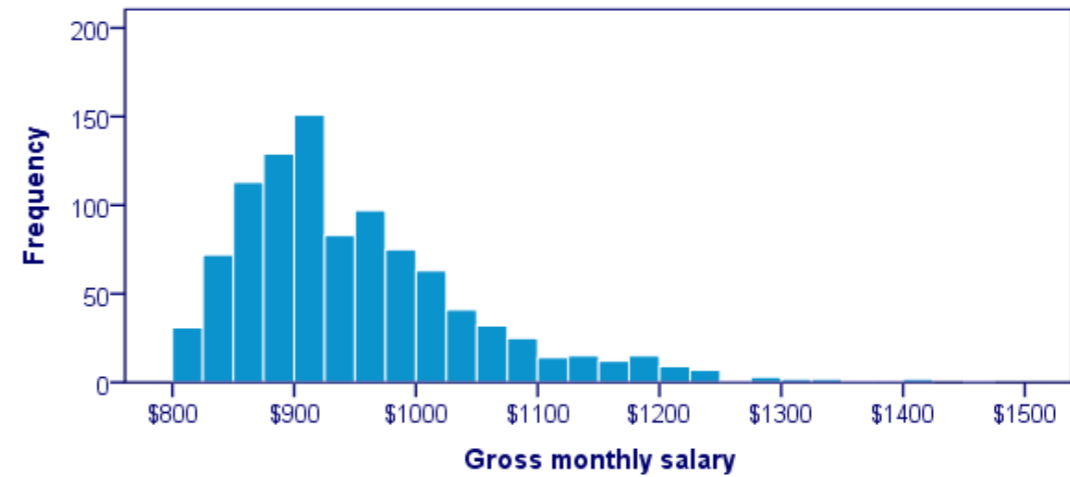
$$\text{number of the bin} = (\text{int}) (\text{value}/\text{width})$$

How to make a histograms...

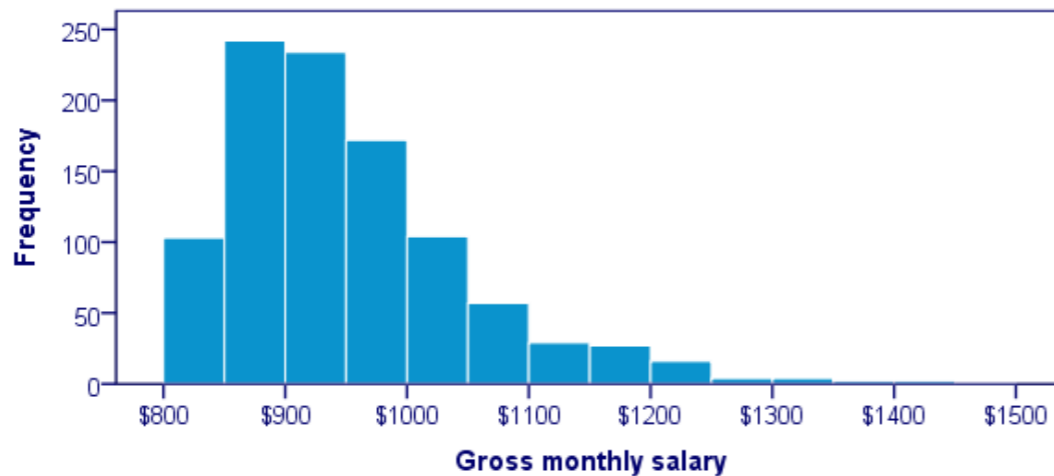
Histogram Bin Width = \$10,-



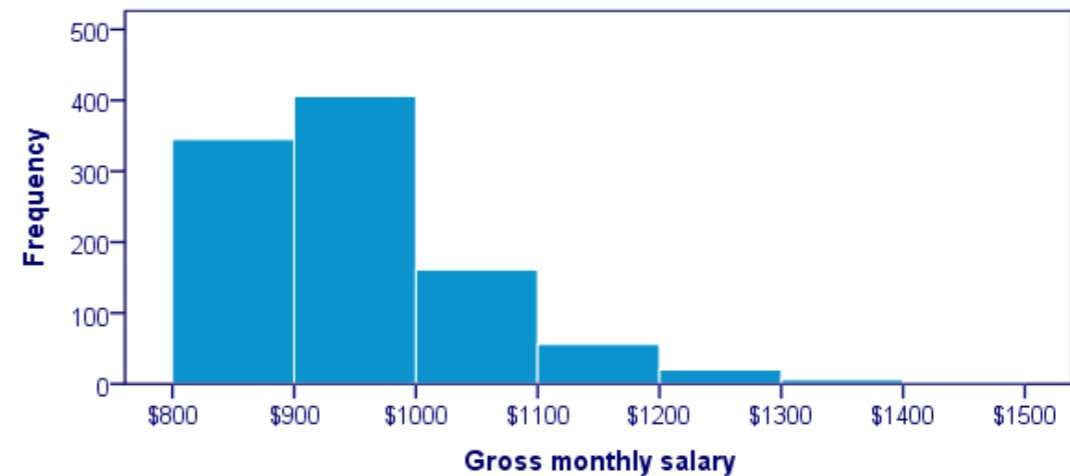
Histogram Bin Width = \$25,-



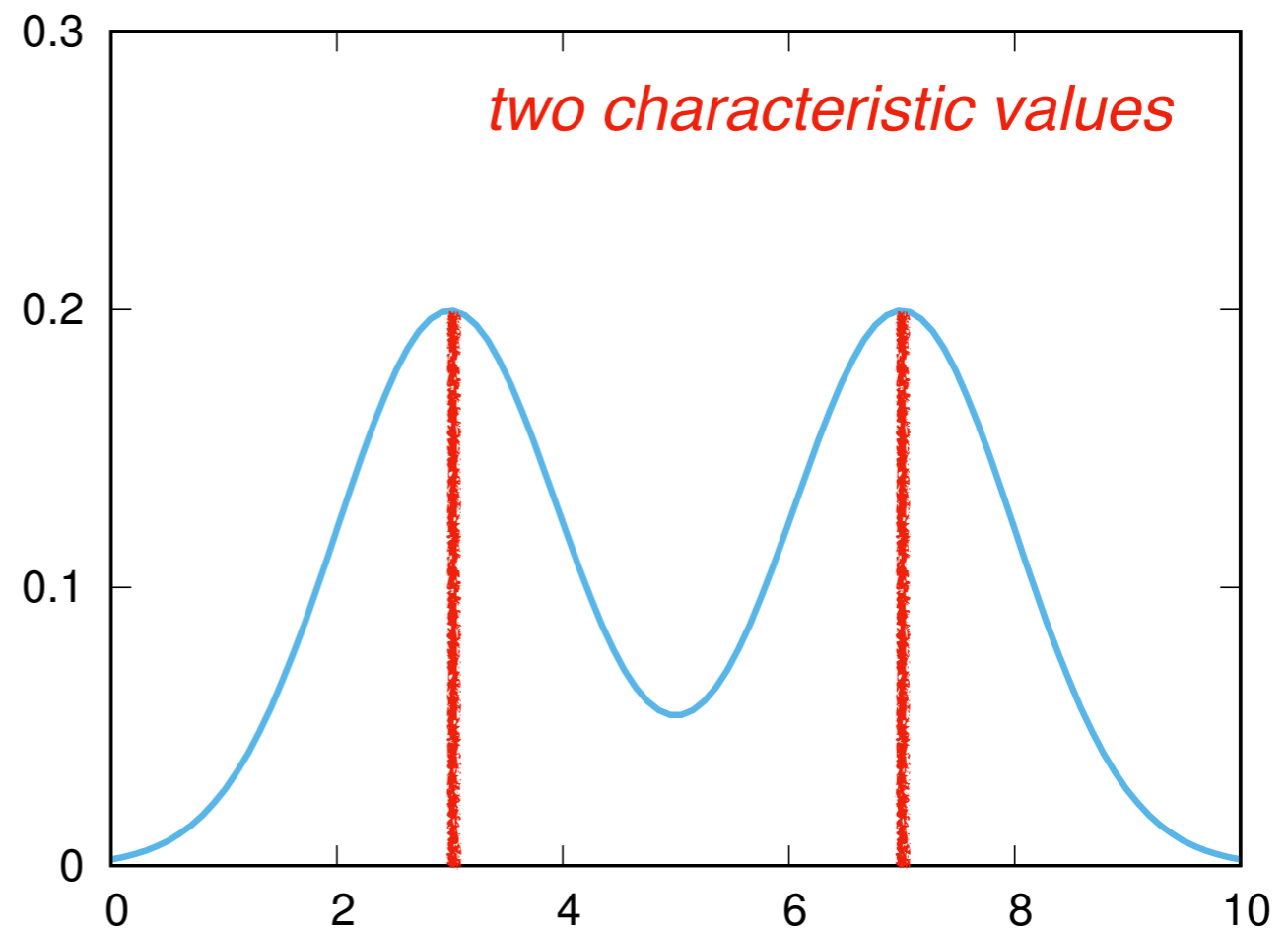
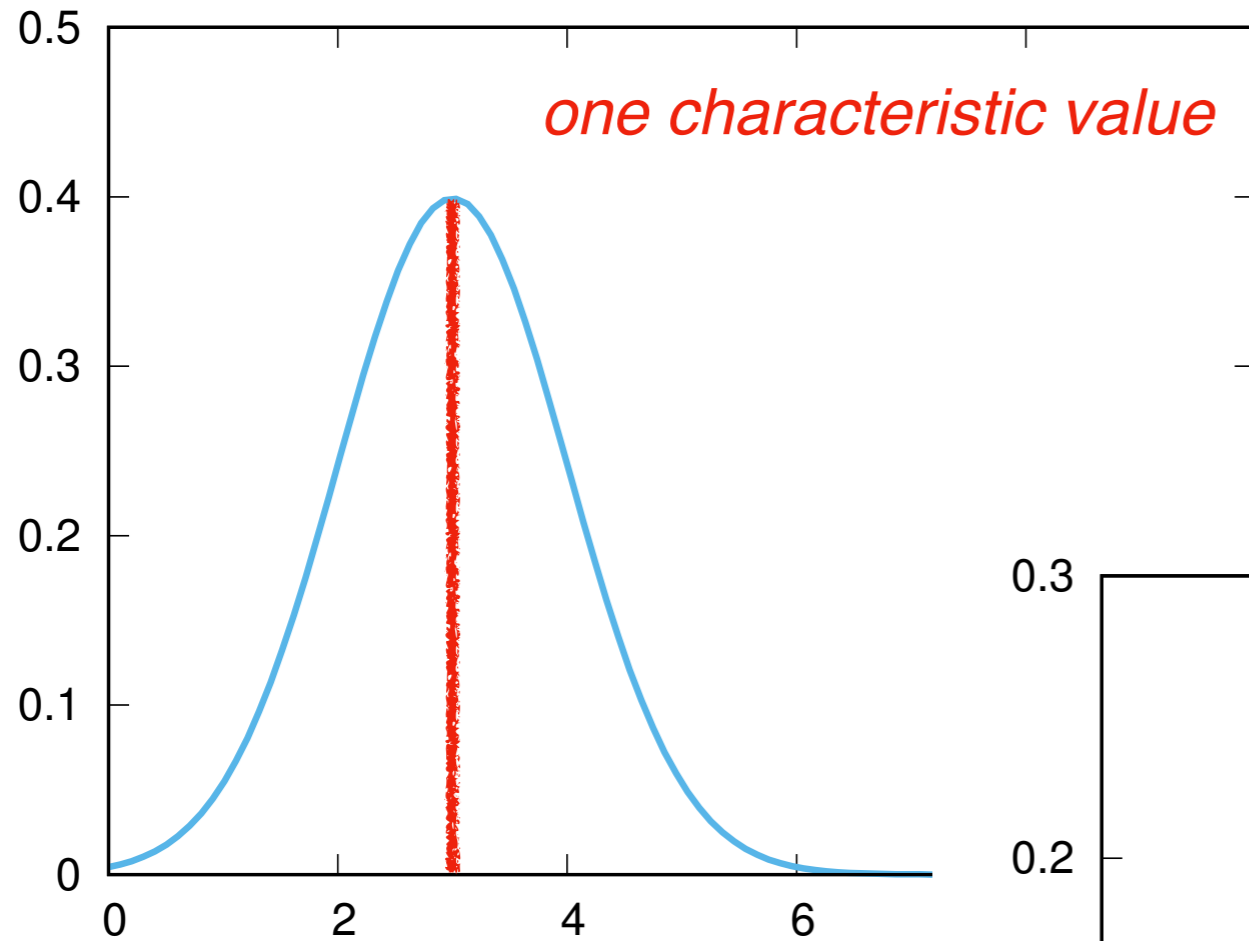
Histogram Bin Width = \$50,-



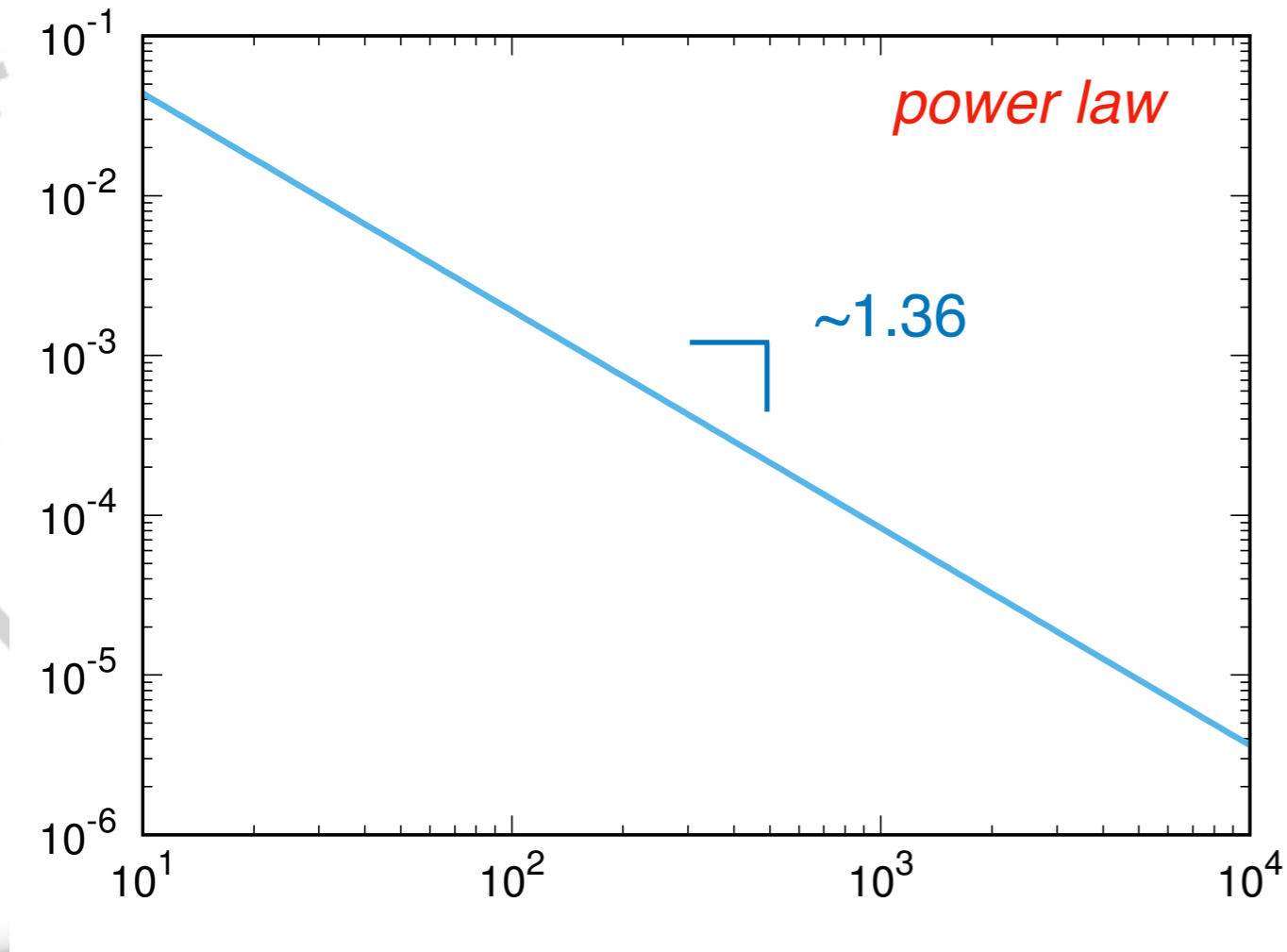
Histogram Bin Width = \$100,-



When does the **average** has a meaning?



When does the **average** has a meaning?



$$P(x) = Ax^{-\alpha}$$
$$\langle x \rangle = \int_0^{x_{\max}} xP(x)dx =$$
$$= A \int_0^{x_{\max}} x^{-\alpha+1} dx$$

$$\langle x \rangle = \frac{A}{2-\alpha} [x_{\max}^{2-\alpha}]$$

$\alpha < 2 :$

$\langle x \rangle$ does not converge...

$\alpha < 3 :$

$\langle x^2 \rangle$ does not converge...