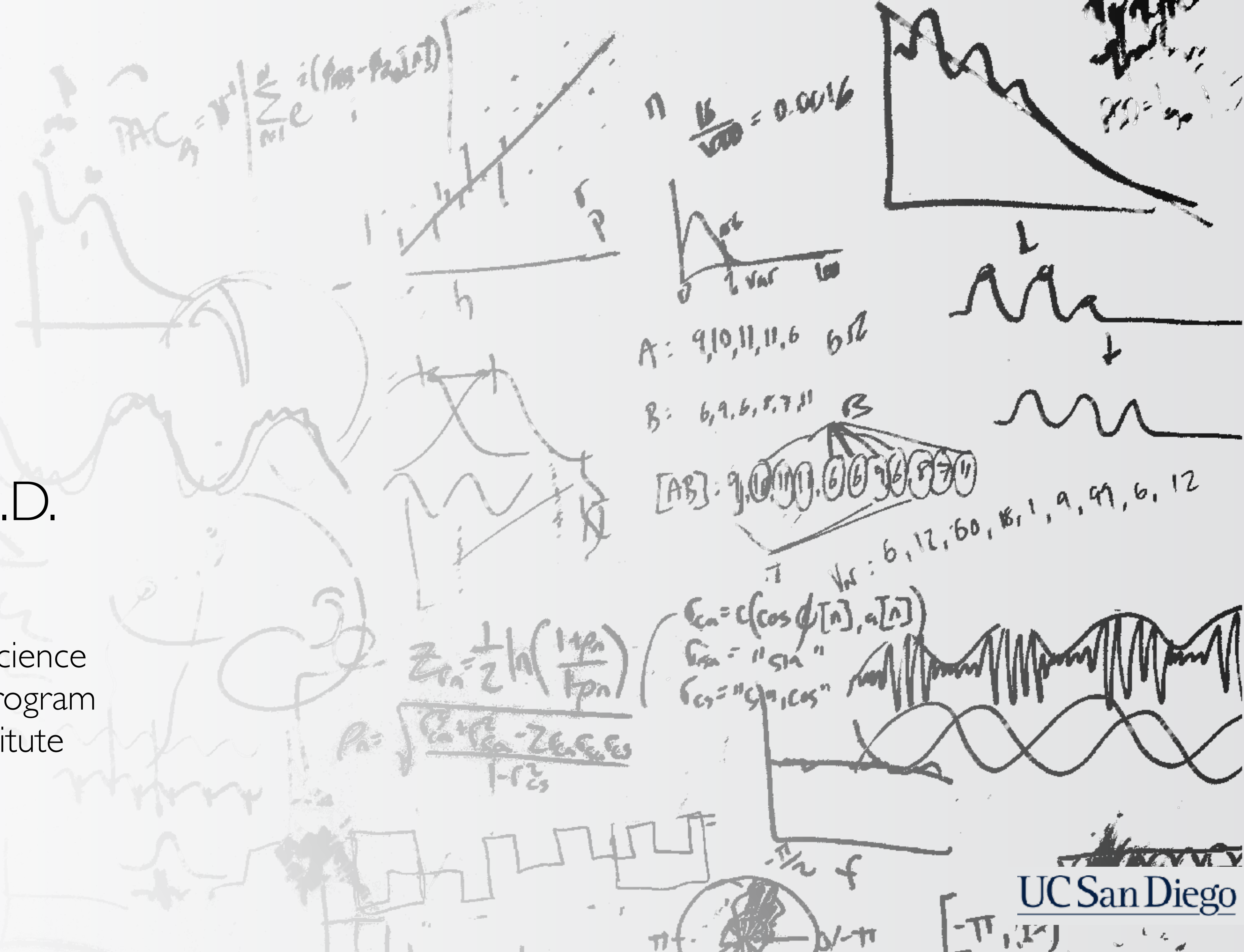


Bradley Voytek, Ph.D.
UC San Diego

Department of Cognitive Science
Neurosciences Graduate Program
Halicioglu Data Science Institute

bvoytek@ucsd.edu
@bradleyvoytek



Fermi Estimation



What is the average height of your group (cm)?

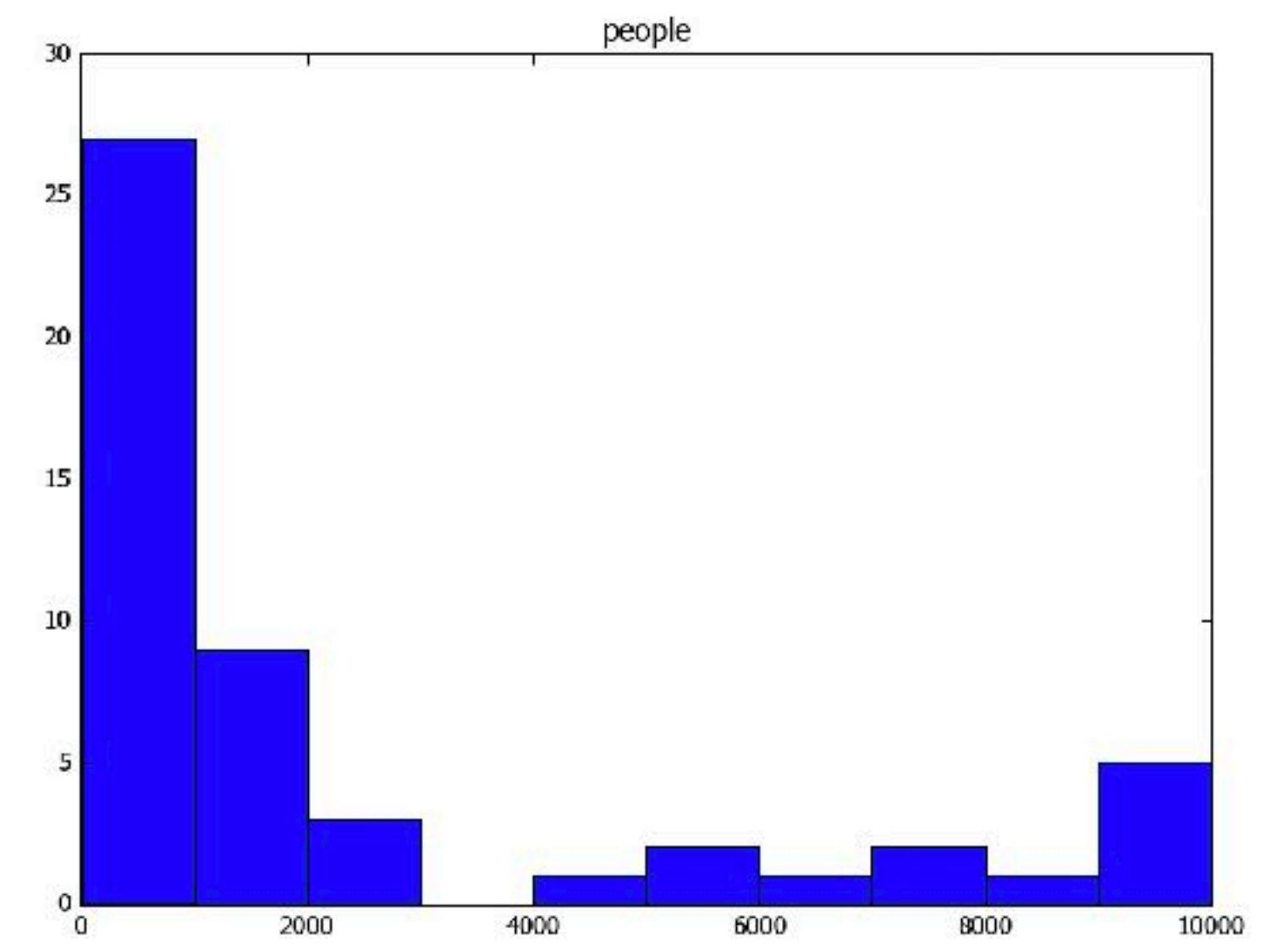
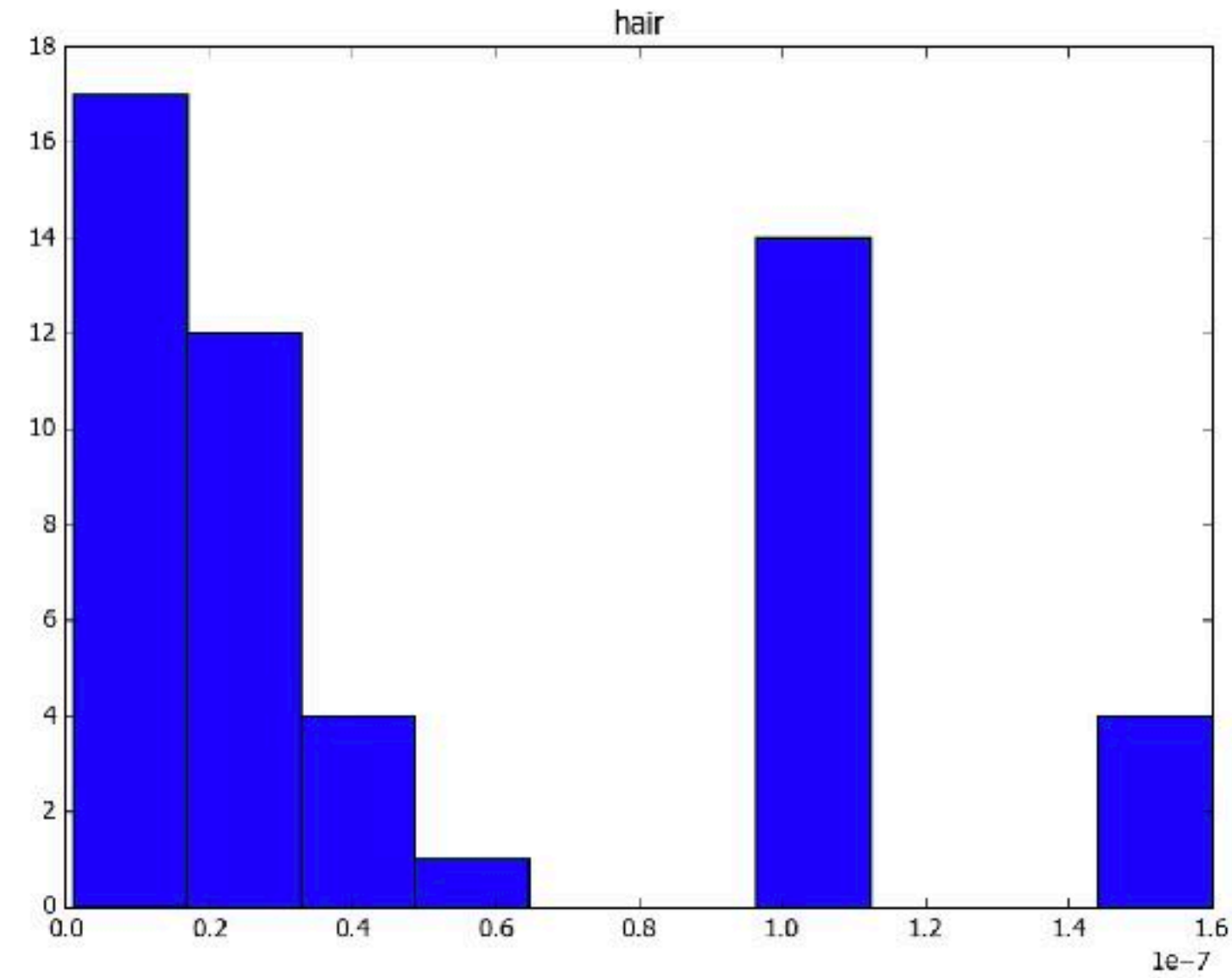
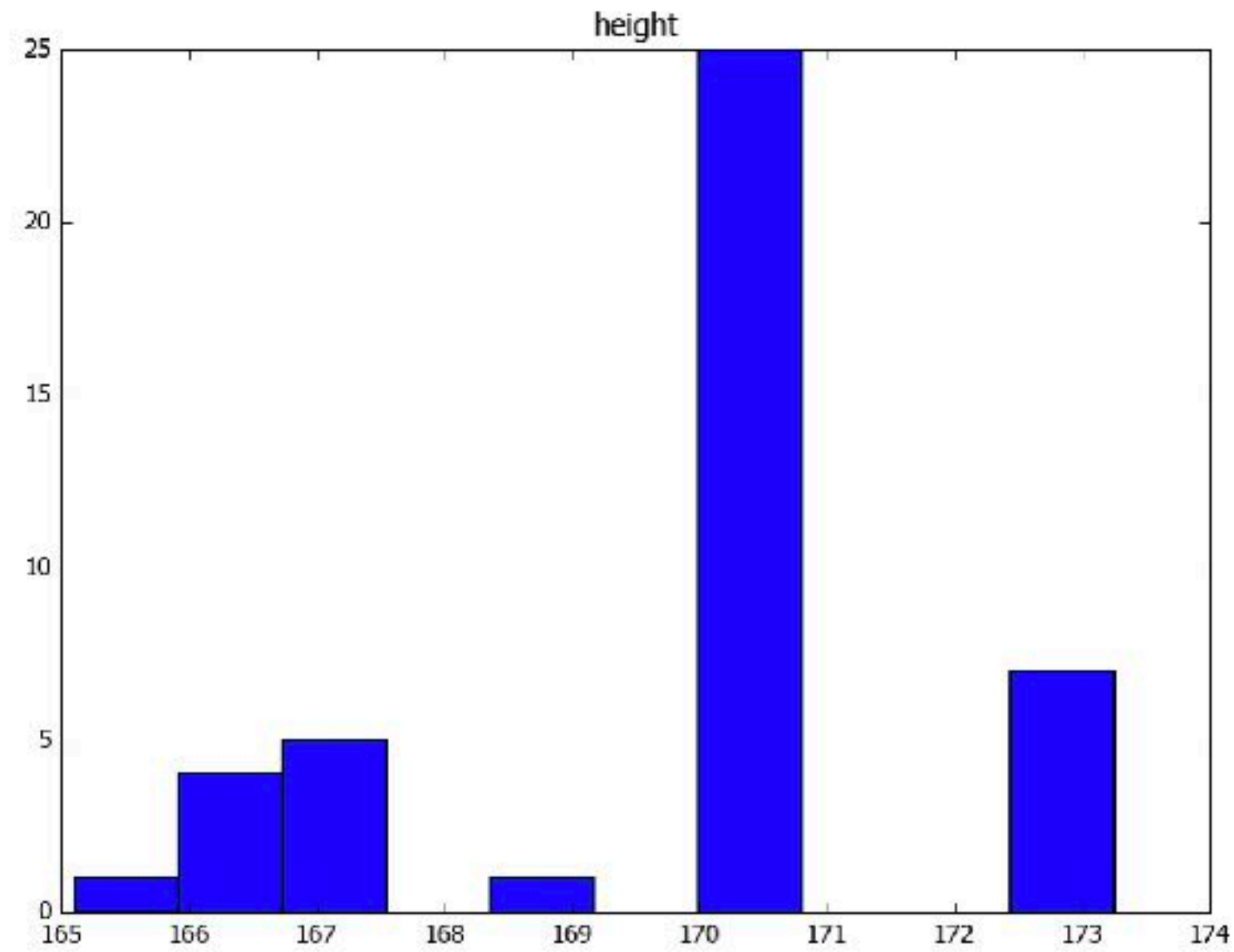
Short answer text

How fast does human hair grow (km/h)?

Short answer text

If every living person stood crammed together side-by-side, how large of an area would they occupy (km²)?

Fermi Estimation



Fermi Estimation



What is the average height of your group (cm)?

Short answer text

```
np.median(height), np.median(hair), np.median(people)  
(170.0, 2.4899999999999998e-08, 1000.0)
```

How fast does human hair grow (km/h)?

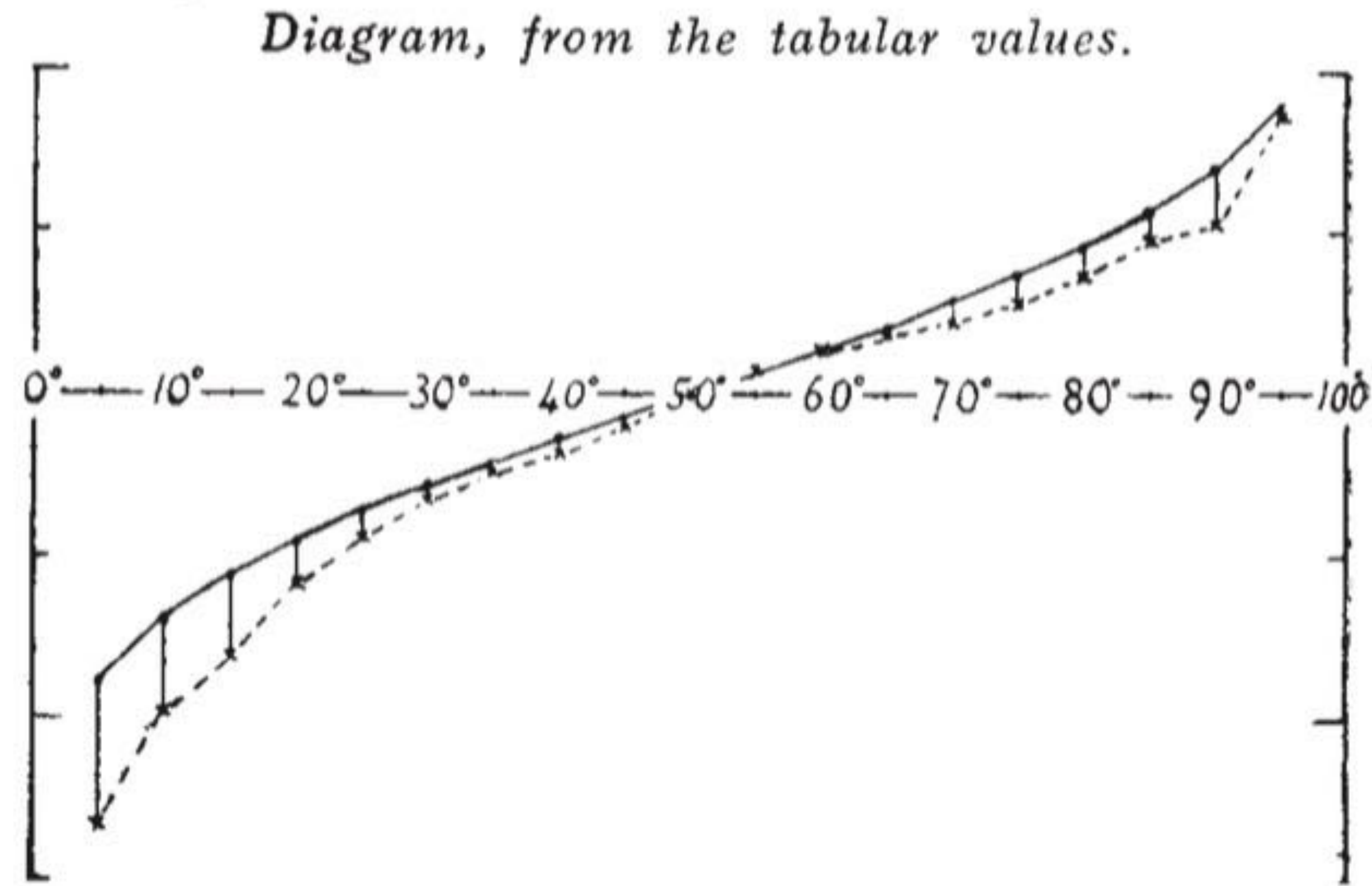
Short answer text

hair: 1.7e-08

If every living person stood crammed together side-by-side, how large of an area would they occupy (km²)?

humans: 1000 km²

Vox Populi



Middlemost estimate: 1207 lbs
True weight: 1198 lbs (0.8% under)

787 estimates

The Wisdom of Crowds

- Diversity of opinion: Each person should have private information even if it's just an eccentric interpretation of the known facts.
- Independence: People's opinions aren't determined by the opinions of those around them.
- Decentralization: People are able to specialize and draw on local knowledge.
- Aggregation: Some mechanism exists for turning private judgments into a collective decision.



COGS 108

Data Science in Practice

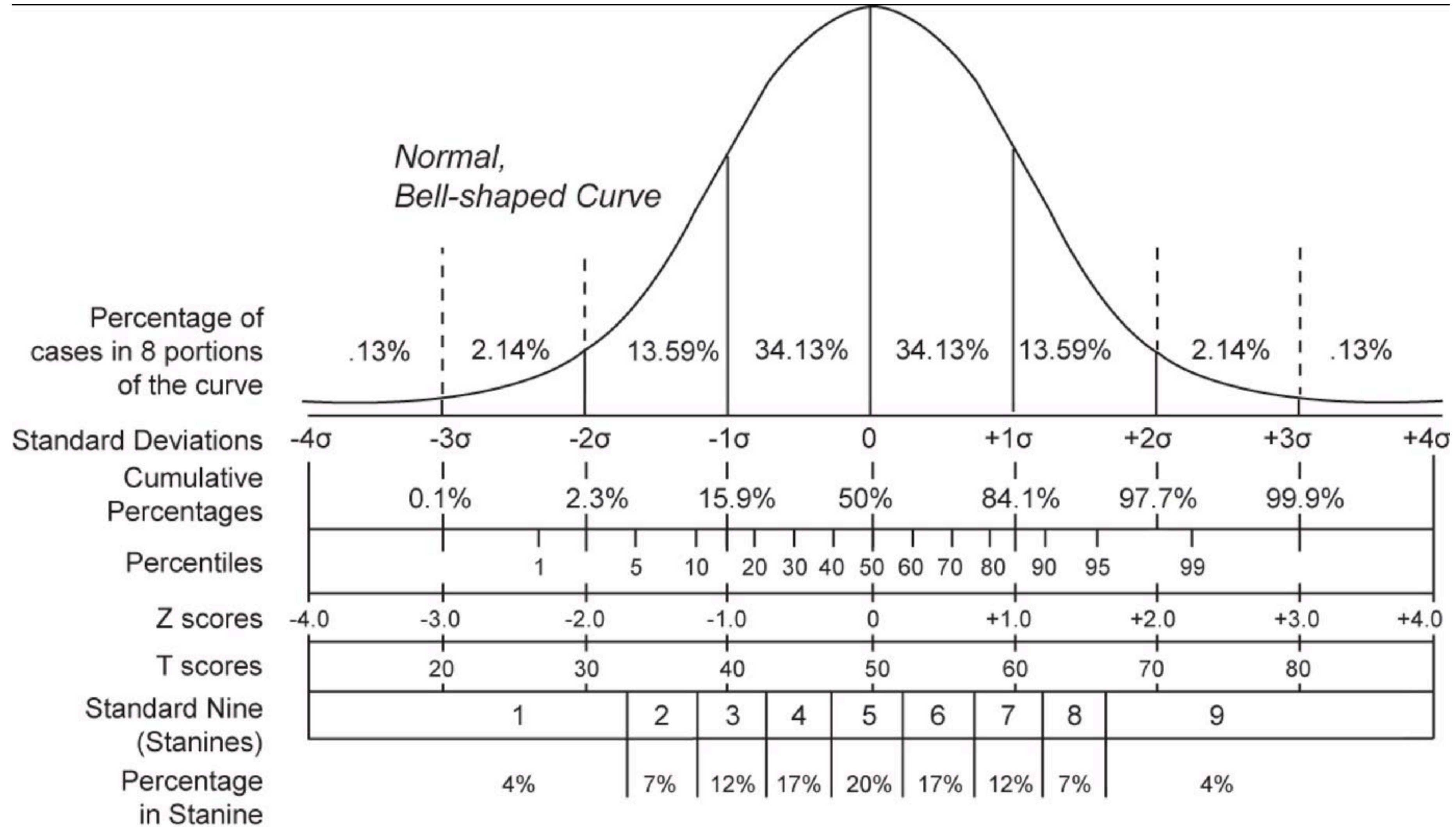
Linear modeling

COGS 108

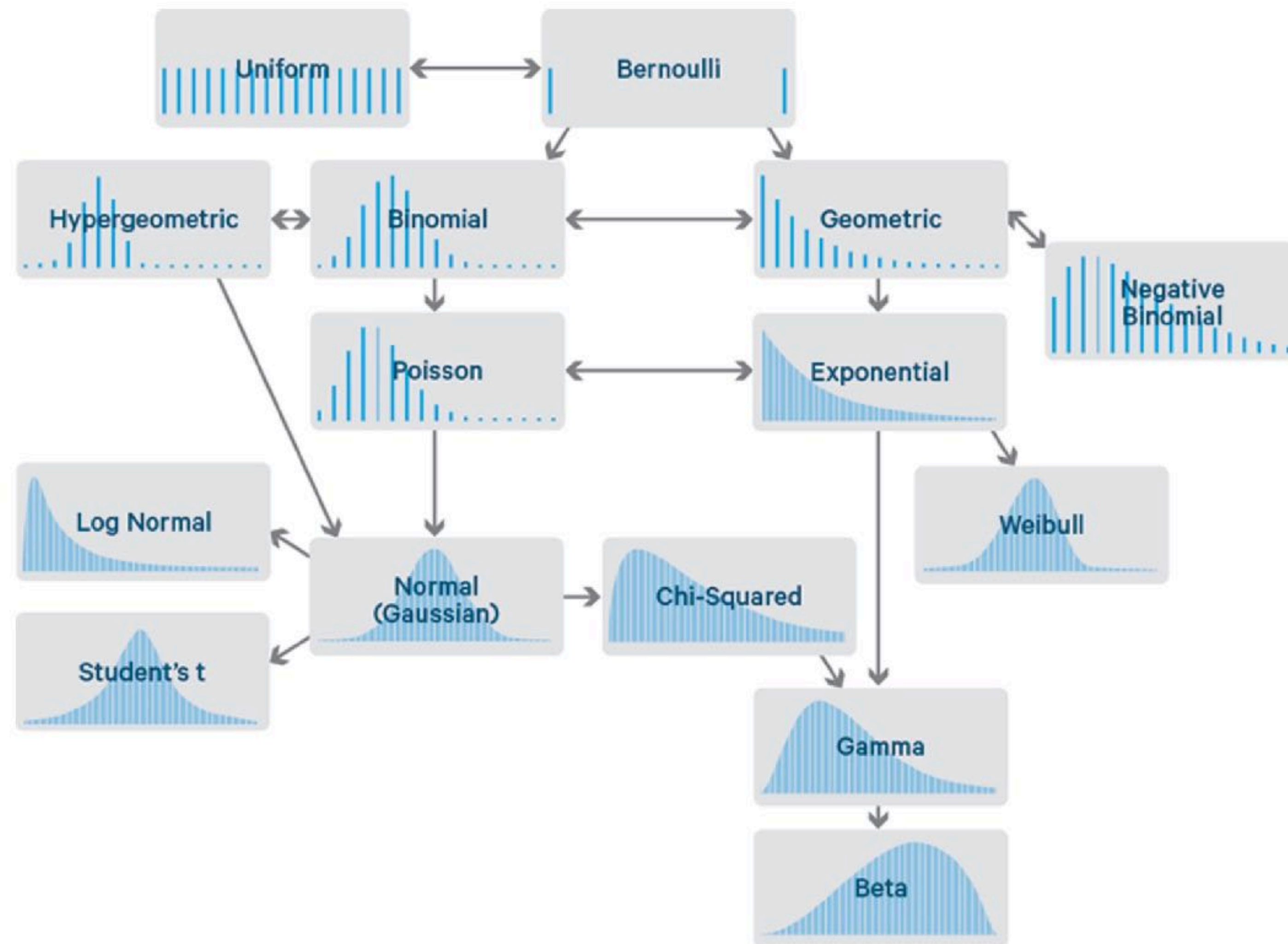
Data Science in Practice

~~*Linear modeling*~~
Intro to statistics

Distributions



Distributions



The diagram illustrates the relationships between various probability distributions, organized into a grid-like structure. The distributions are categorized by their properties and relationships.

Properties:

- C: Convolution
- F: Forgetfulness
- I: Inverse
- L: Linear combination
- M: Minimum
- P: Product
- R: Residual
- S: Scaling
- V: Variate generation
- X: Maximum

Relationships:

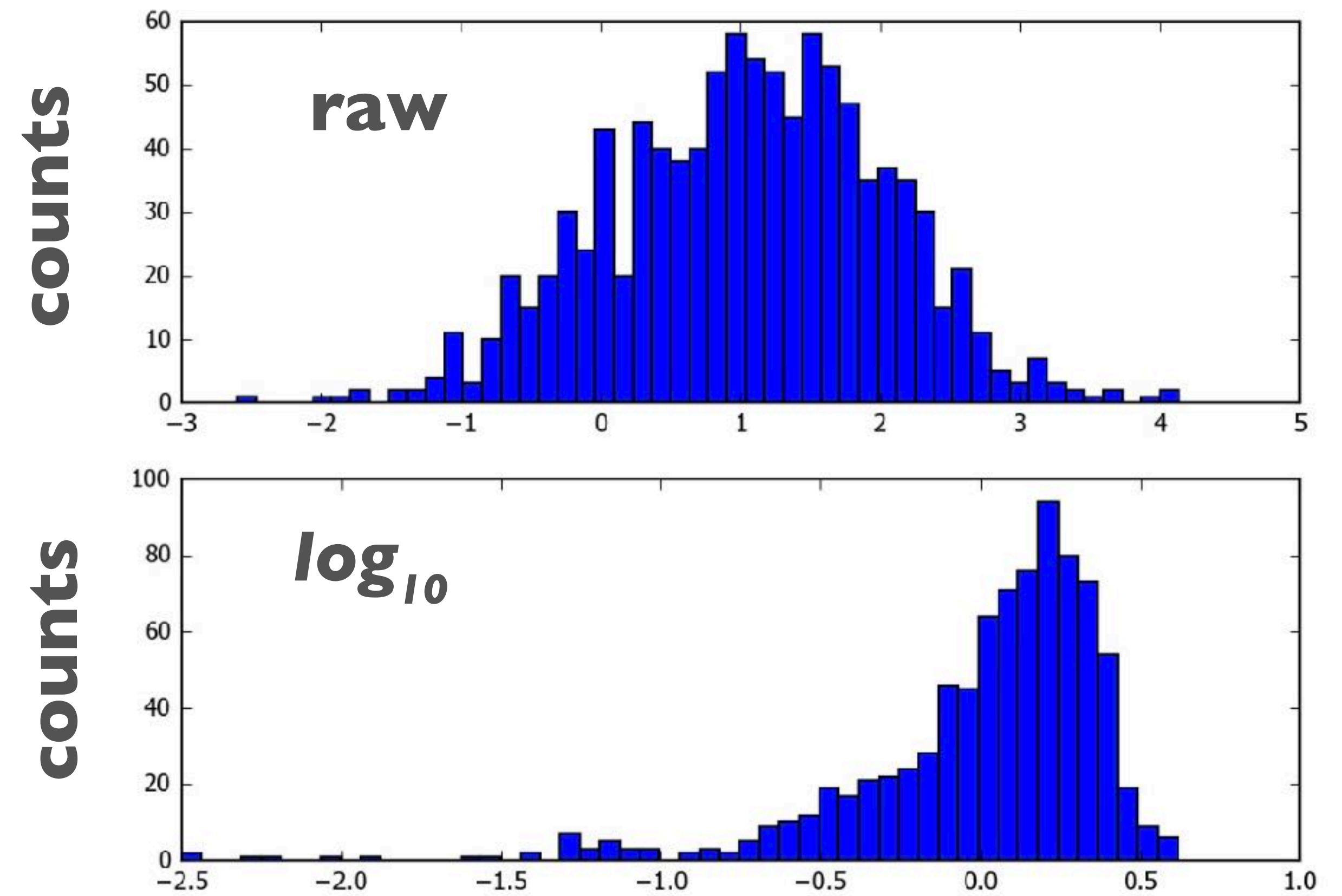
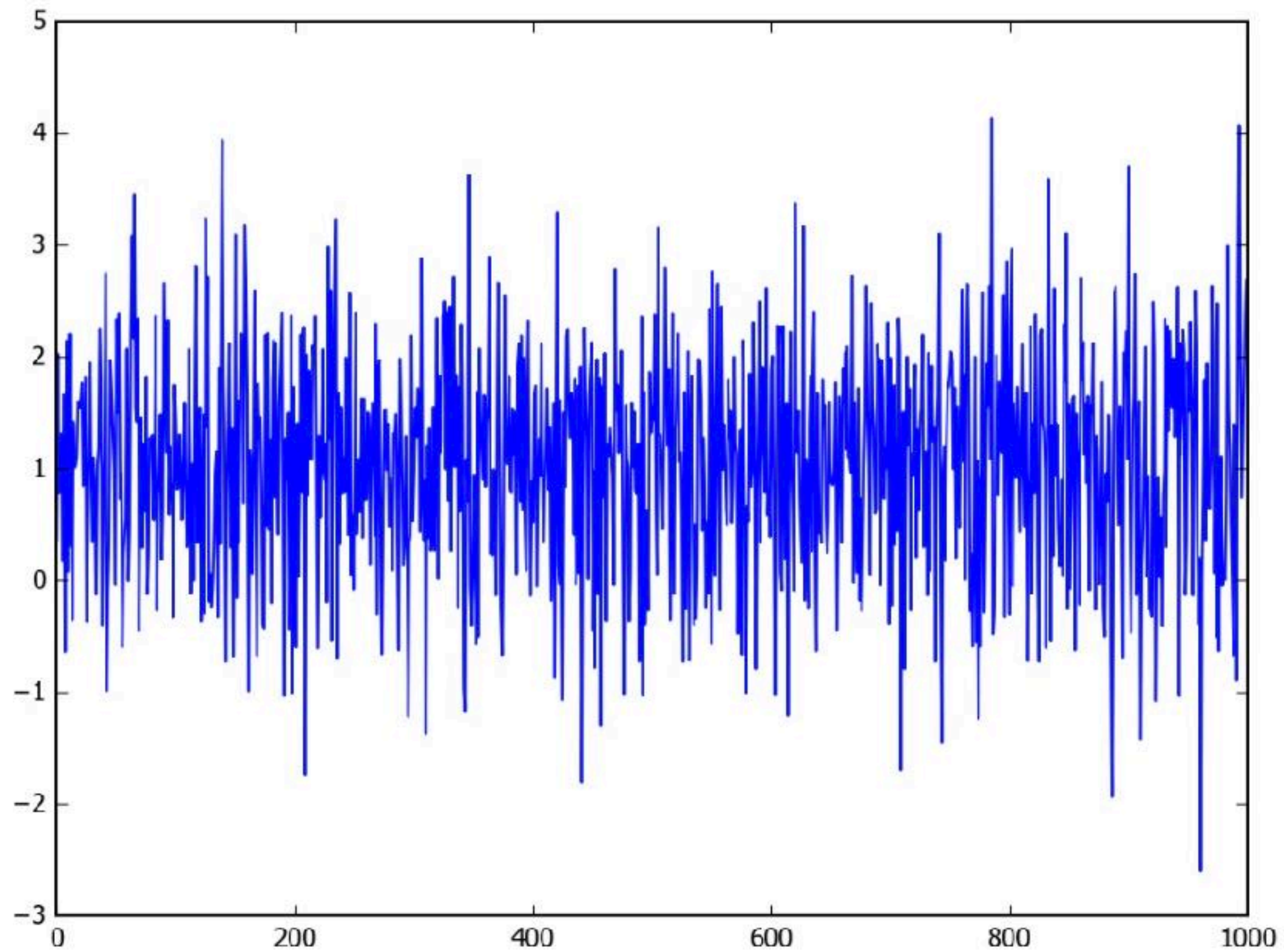
- Special cases
- Transformations
- Limiting
- Bayesian

The diagram shows a wide range of distributions, including:

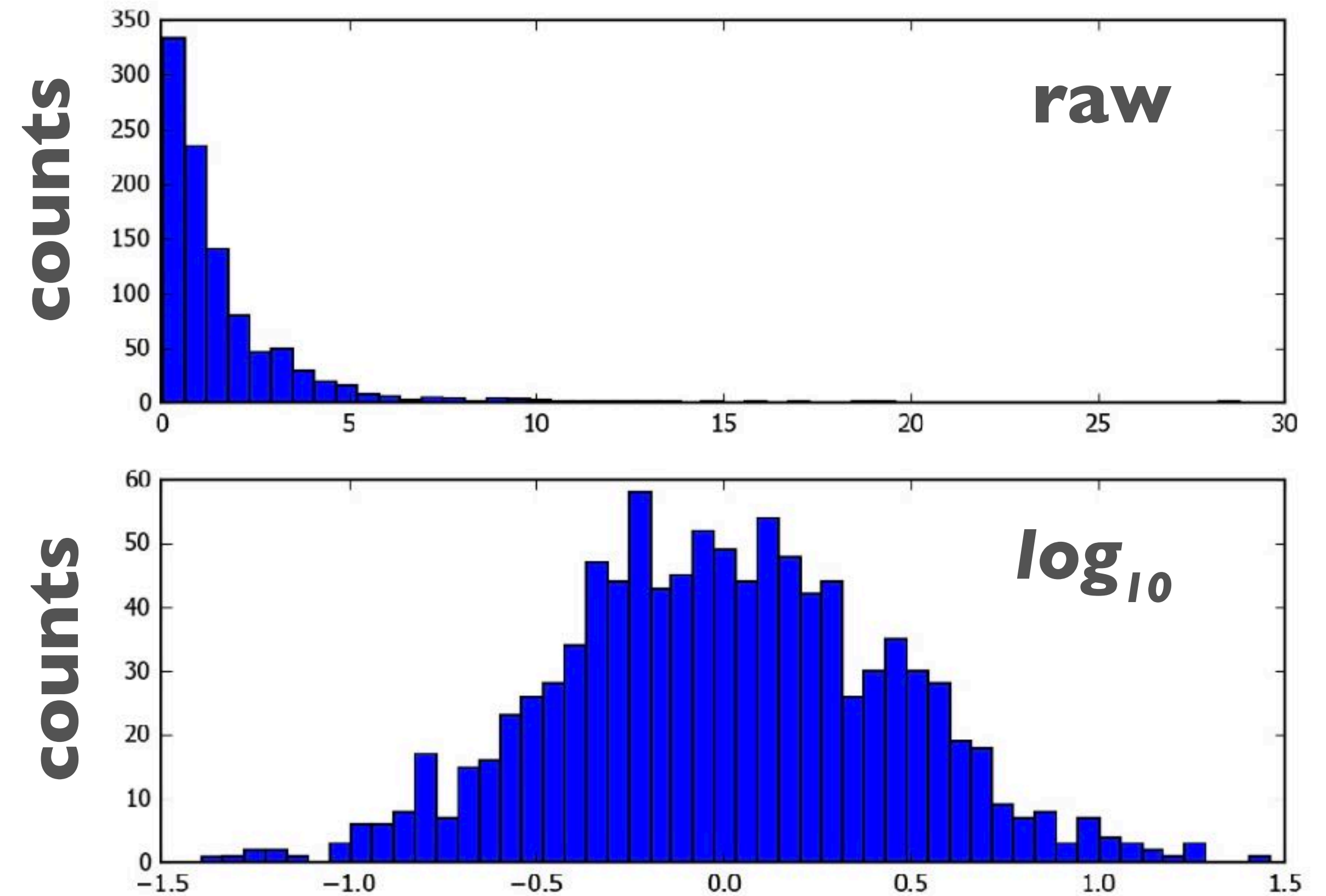
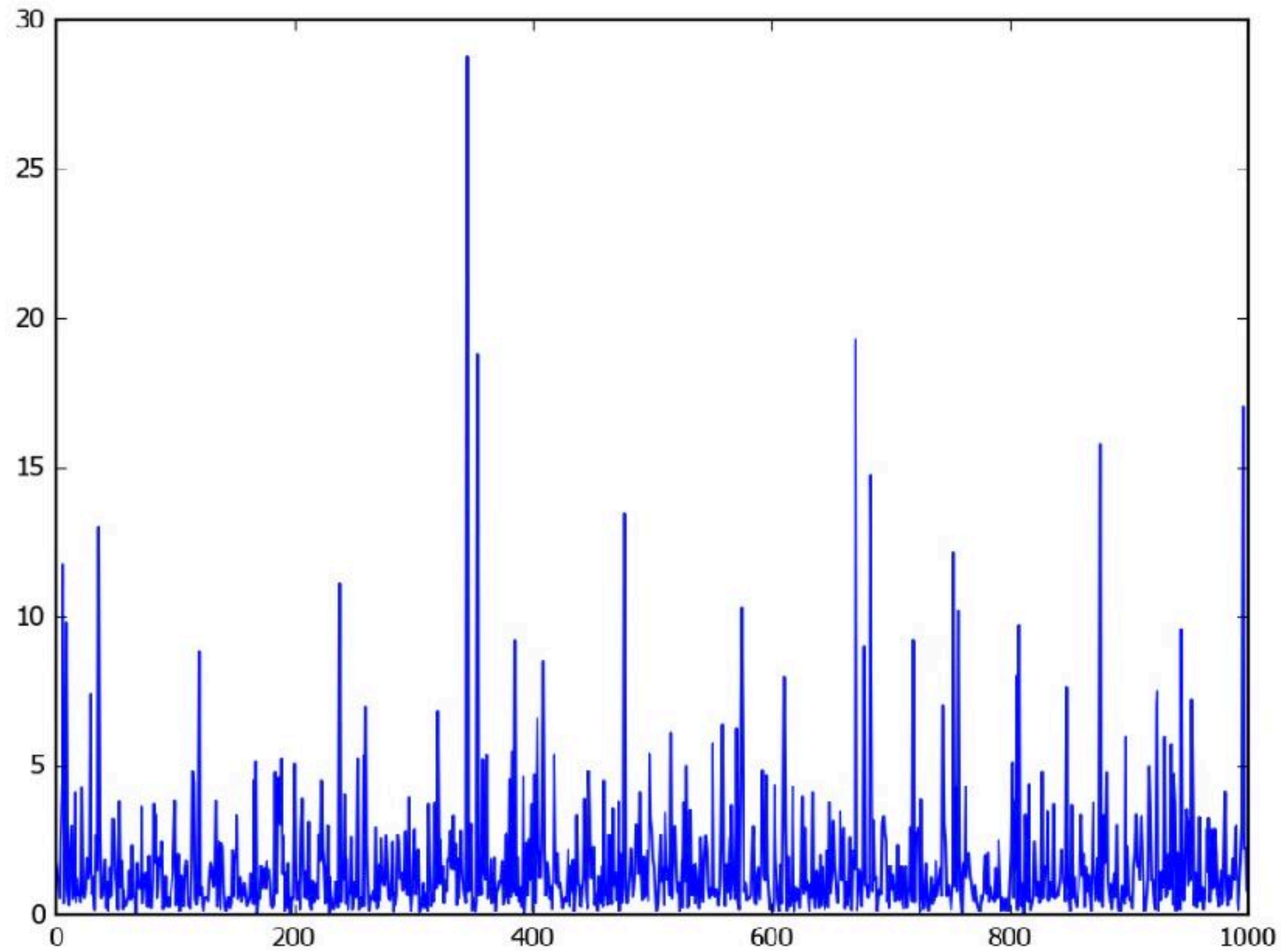
- Discrete distributions:** Ziplf, Discrete uniform, Rectangular, Beta-binomial, Negative hypergeometric, Poisson, Binomial, Pascal, Geometric, Discrete Weibull, Hypergeometric, Bernoulli, Gamma-Poisson, Normal, Noncentral chi-square, Chi, Inverse Gaussian, Standard Wald, t, F, Noncentral t, Doubly noncentral t, Noncentral F, Doubly noncentral F, Rayleigh, Weibull, Extreme value, Generalized Pareto, Logistic, Triangular, Standard triangular, Uniform, von Mises, Kolmogorov-Esmirnov.
- Continuous distributions:** Beta-Pascal, Gamma, Lognormal, Log gamma, Generalized gamma, Gamma, Inverted gamma, Gamma, Inverted beta, Beta, Dirichlet, Logistic-exponential, Exponential power, Exponential, Gompertz, Makhsoum, Standard power, Power, TSP, Benford, Lomax, Logistic, Triangular, Standard triangular, Uniform, von Mises, Kolmogorov-Esmirnov.
- Other distributions:** Zeta, Logarithm, Power series, Poisson, Binomial, Pascal, Geometric, Discrete Weibull, Hypergeometric, Bernoulli, Gamma-Poisson, Normal, Noncentral chi-square, Chi, Inverse Gaussian, Standard Wald, t, F, Noncentral t, Doubly noncentral t, Noncentral F, Doubly noncentral F, Rayleigh, Weibull, Extreme value, Generalized Pareto, Logistic, Triangular, Standard triangular, Uniform, von Mises, Kolmogorov-Esmirnov.

The diagram is a comprehensive reference for understanding the connections between different probability distributions.

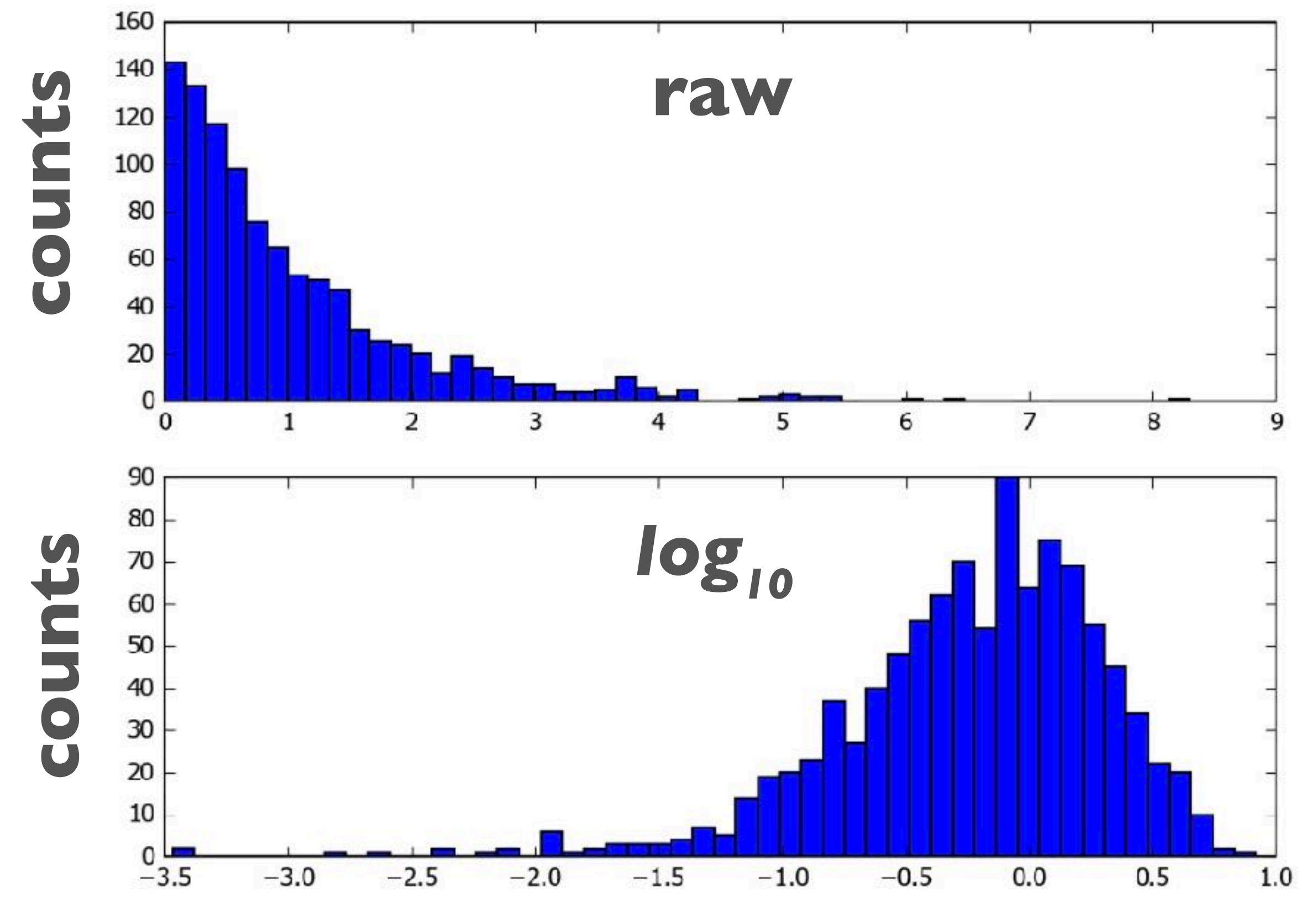
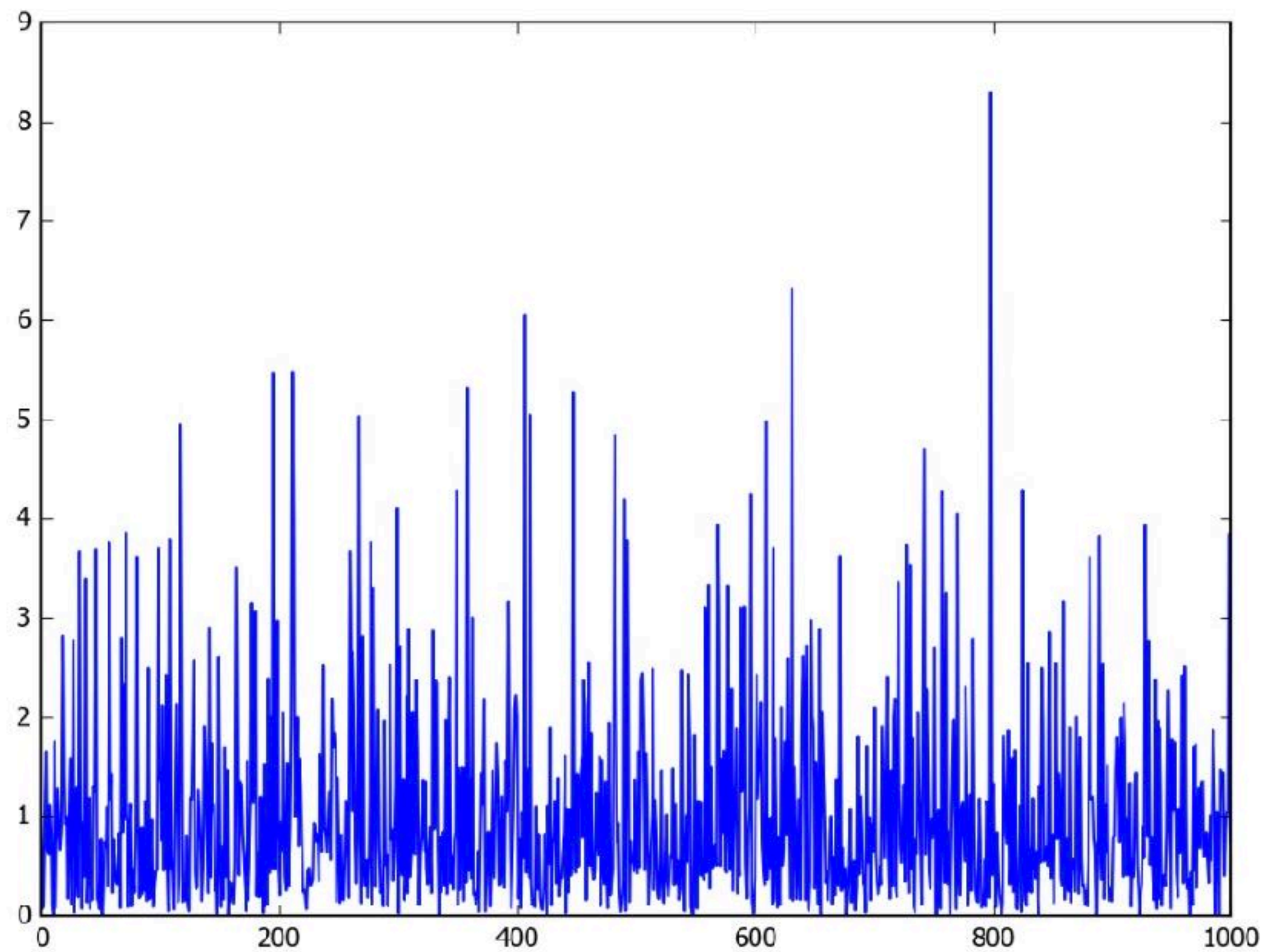
Normal



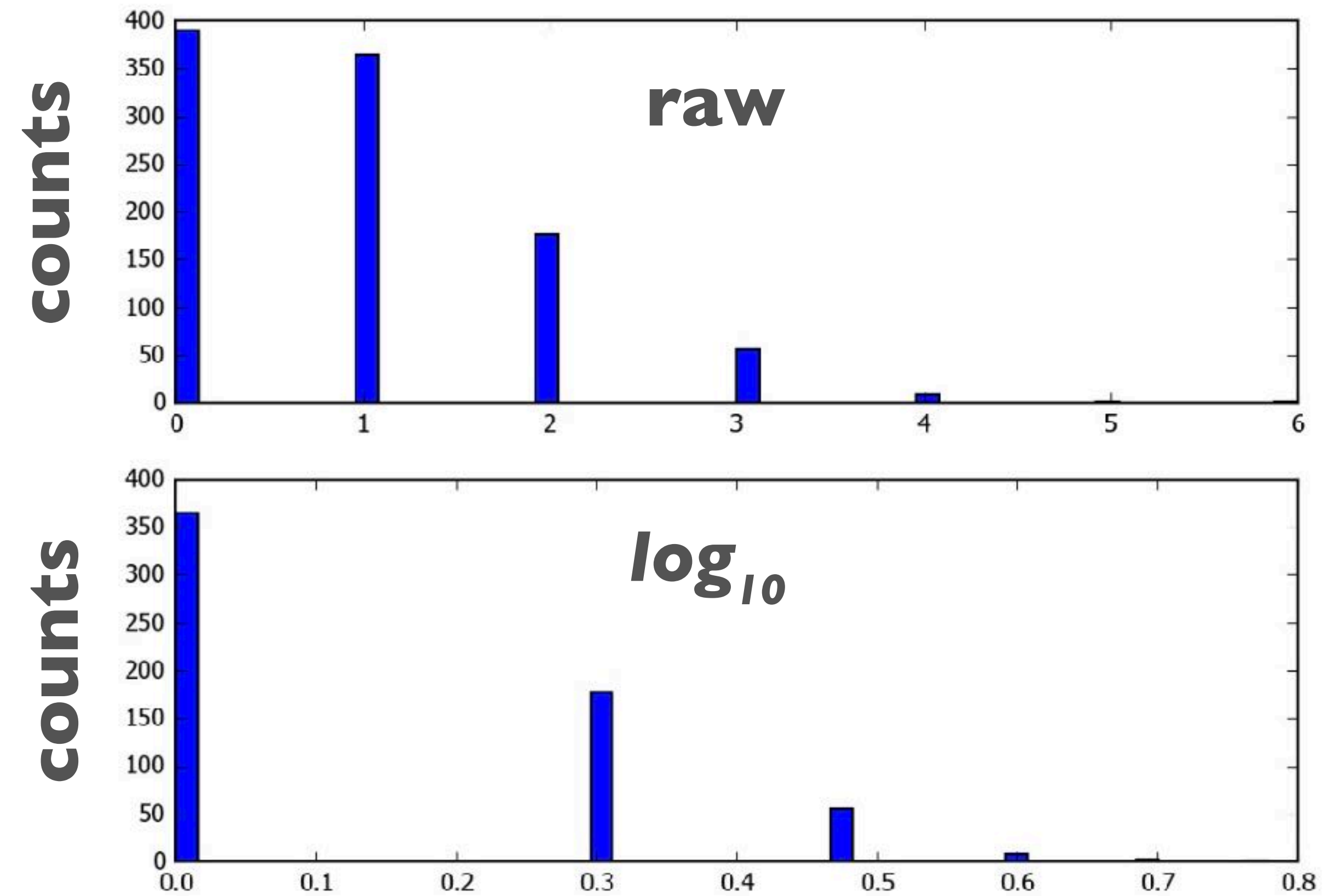
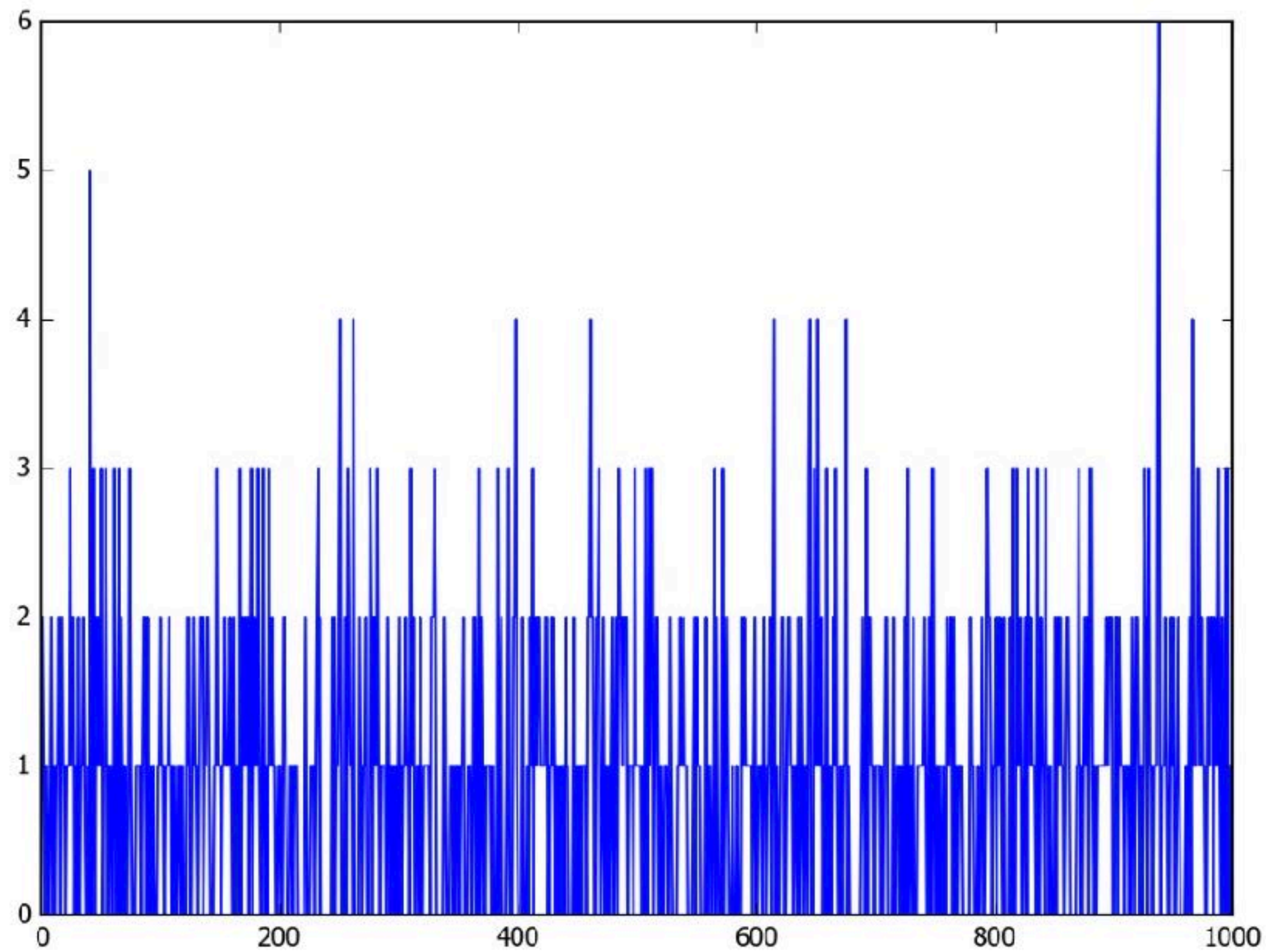
Lognormal



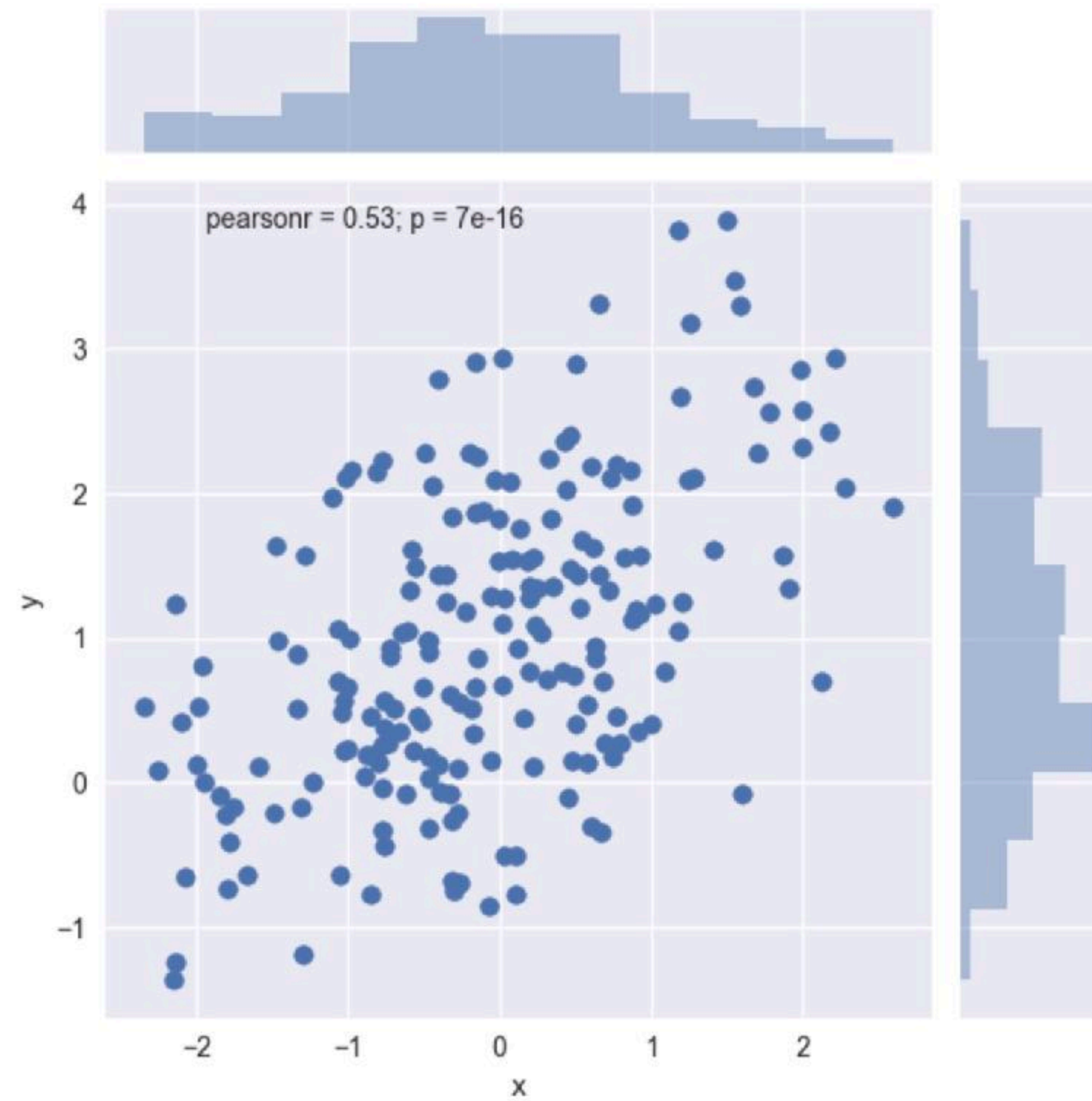
Exponential



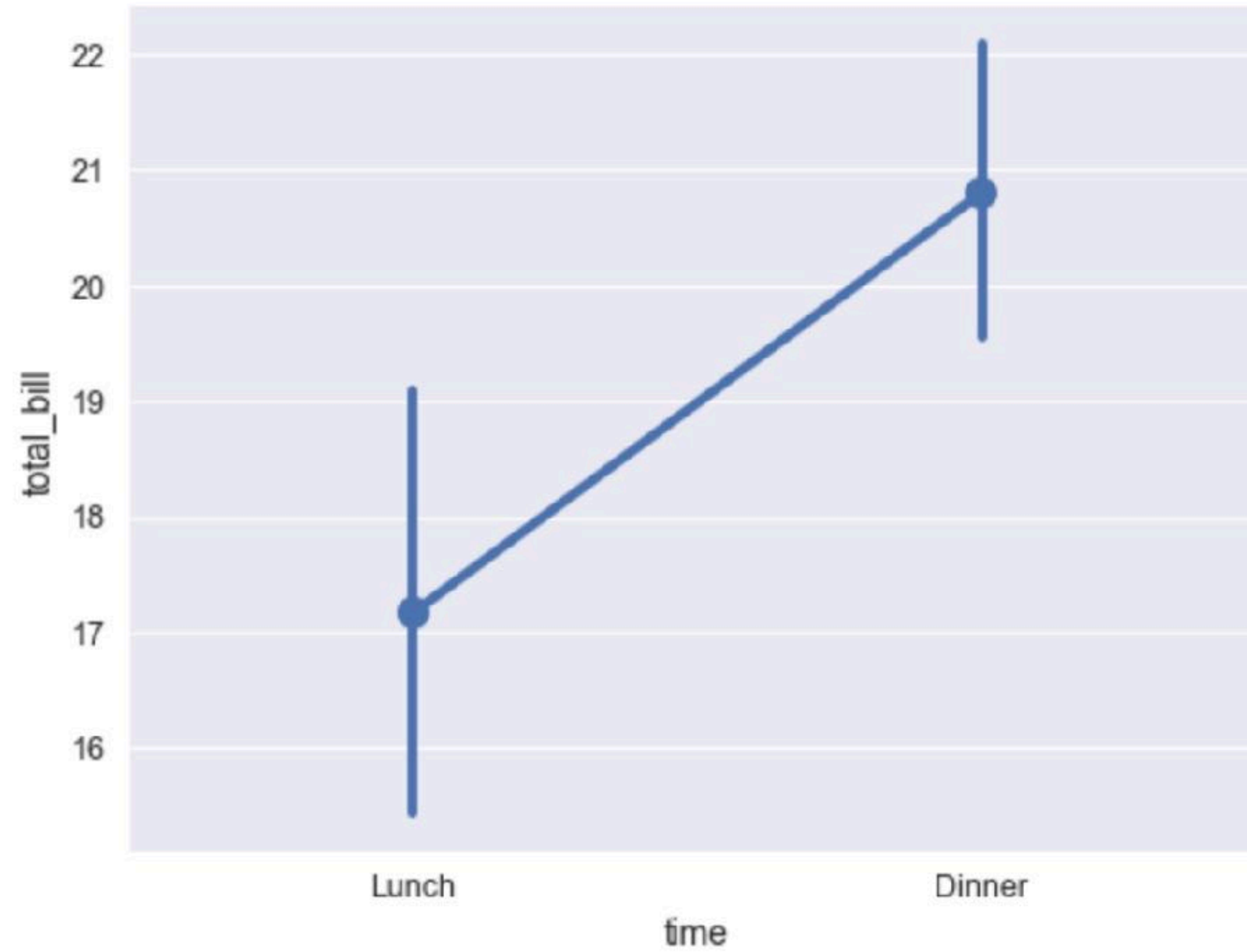
Poisson



Data distributions



Data distributions



Student's t-test

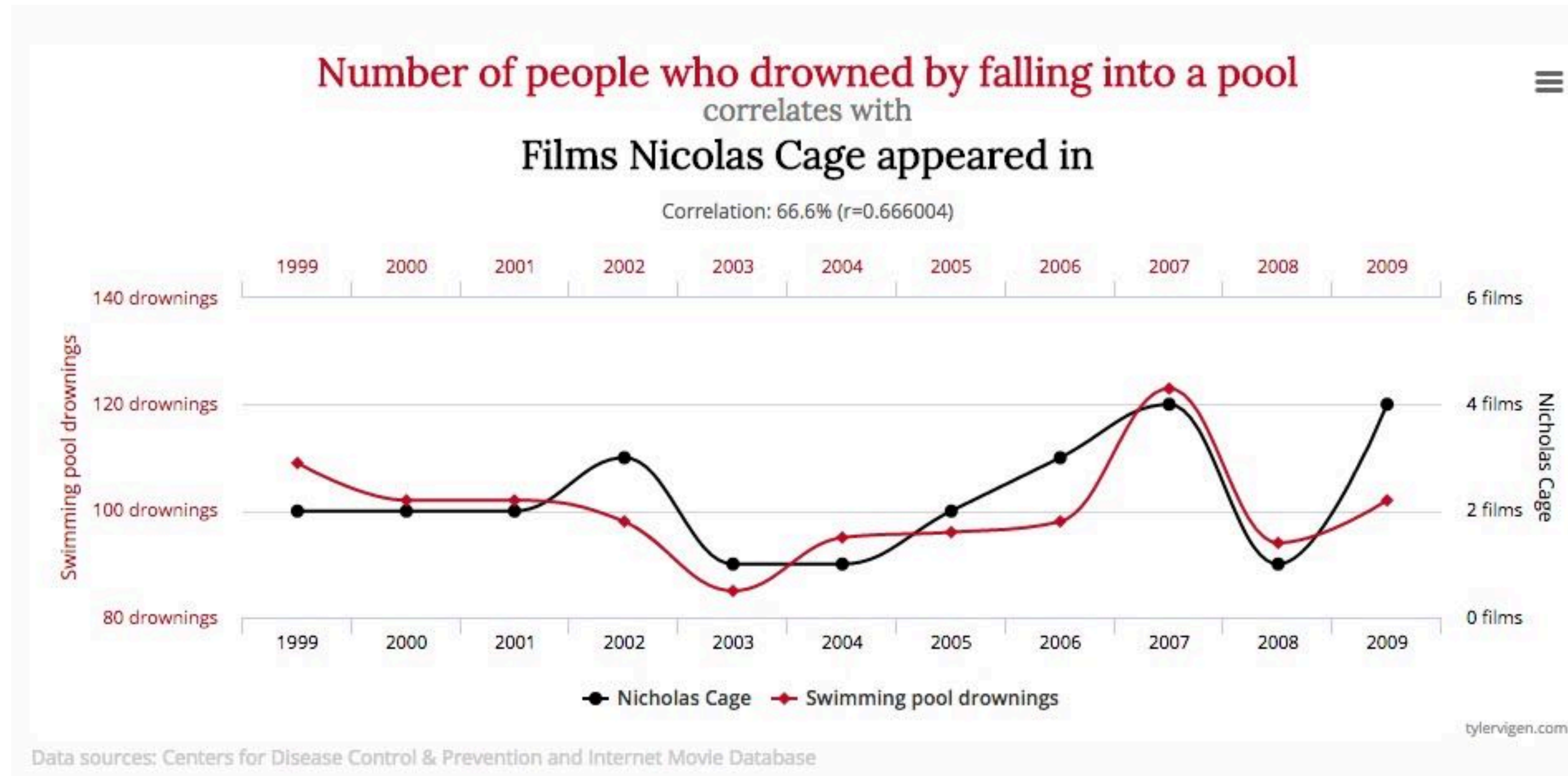


Student's t-test

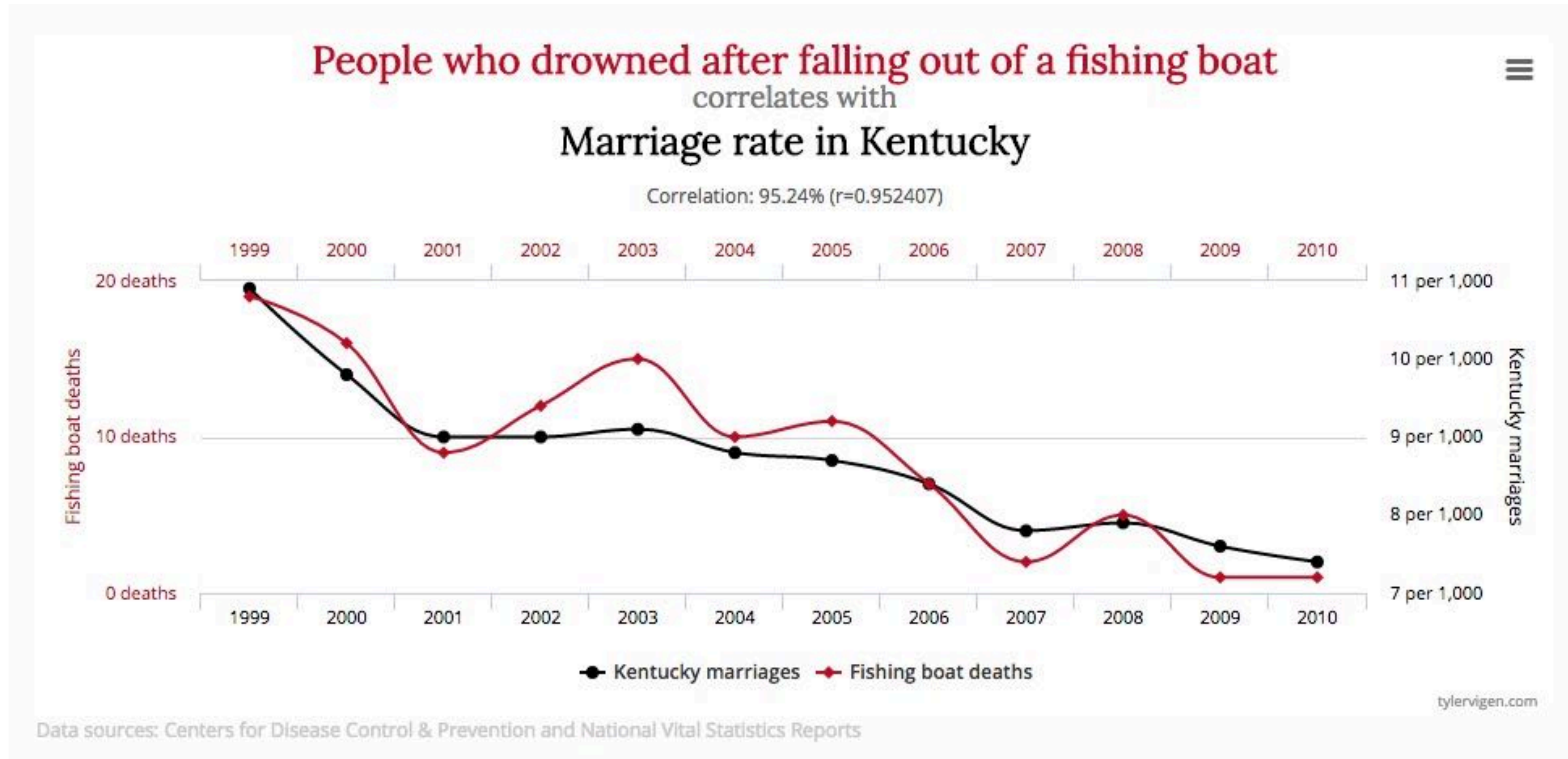
...brewing a batch of beer is a time consuming and expensive process, so in order to draw conclusions from experimental methods applied to just a few batches, Gosset had to figure out the role of chance in determining how a batch of beer had turned out. Guinness frowned upon academic publications, so Gosset had to publish his results under the modest pseudonym, “Student.”



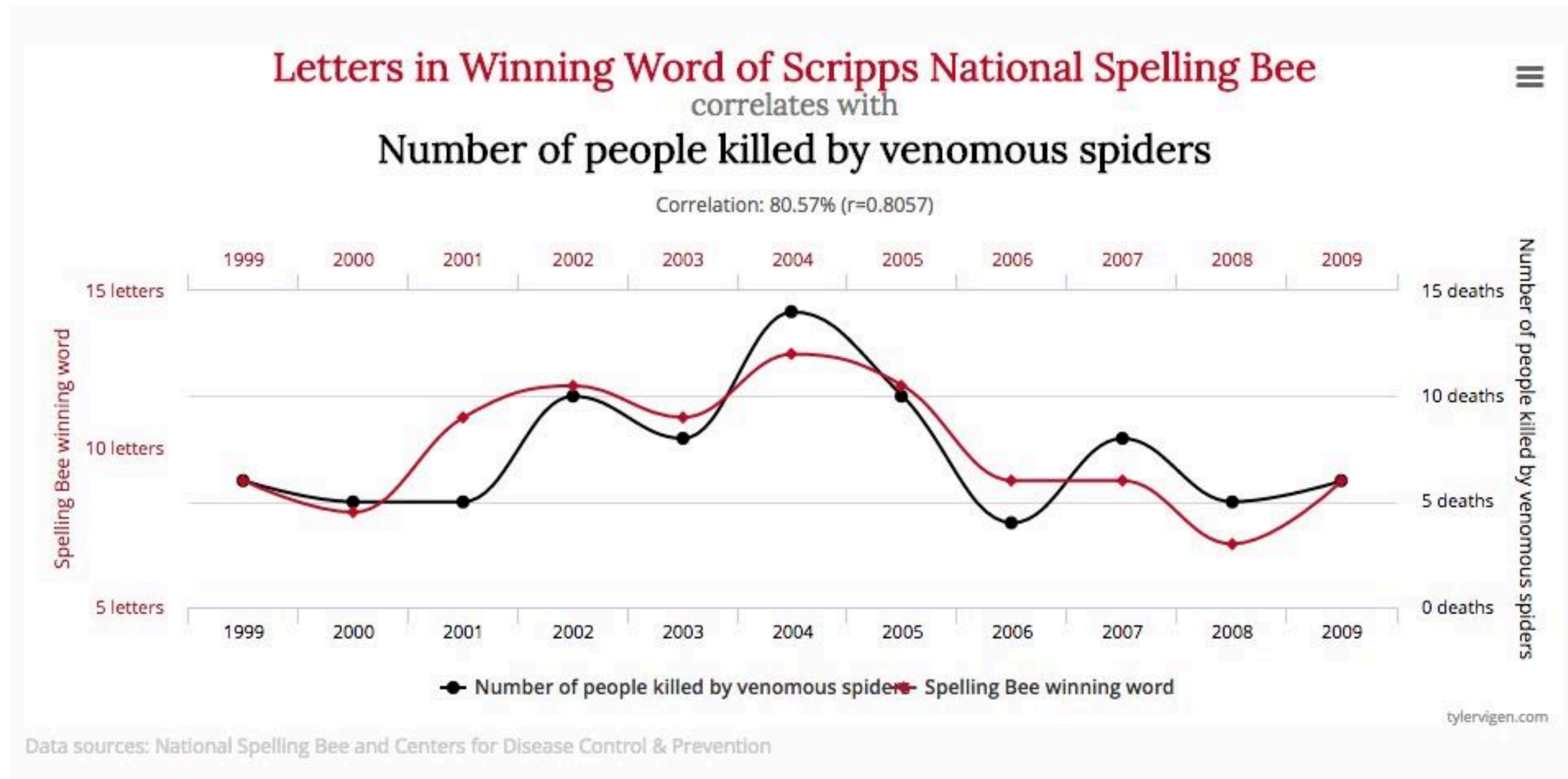
Don't trust your data (or narratives)



Don't trust your data (or narratives)



Don't trust your data (or narratives)



A line

$$y = mx + b$$

slope

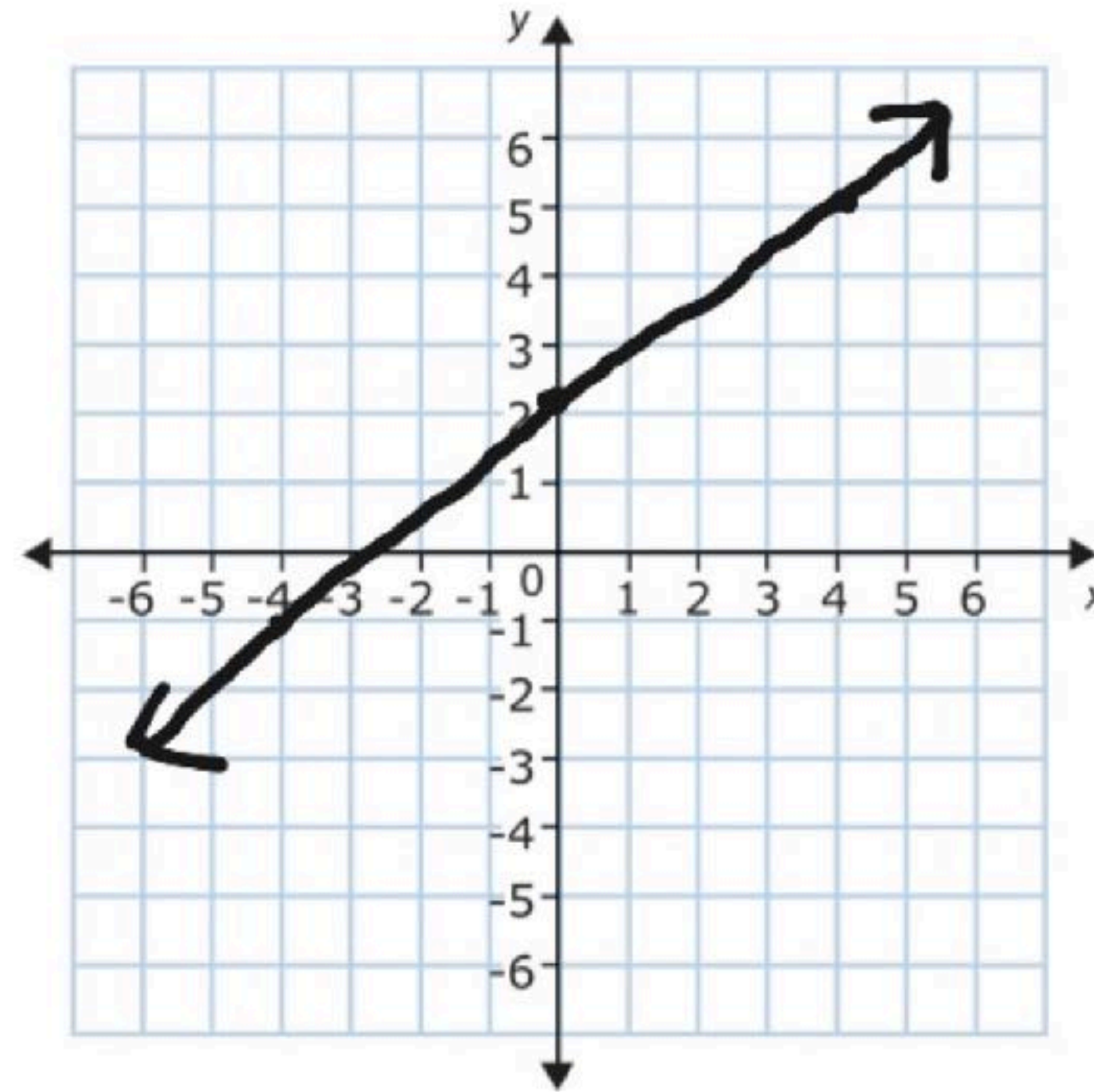
intercept

A line

$$y = mx + b$$

slope

intercept

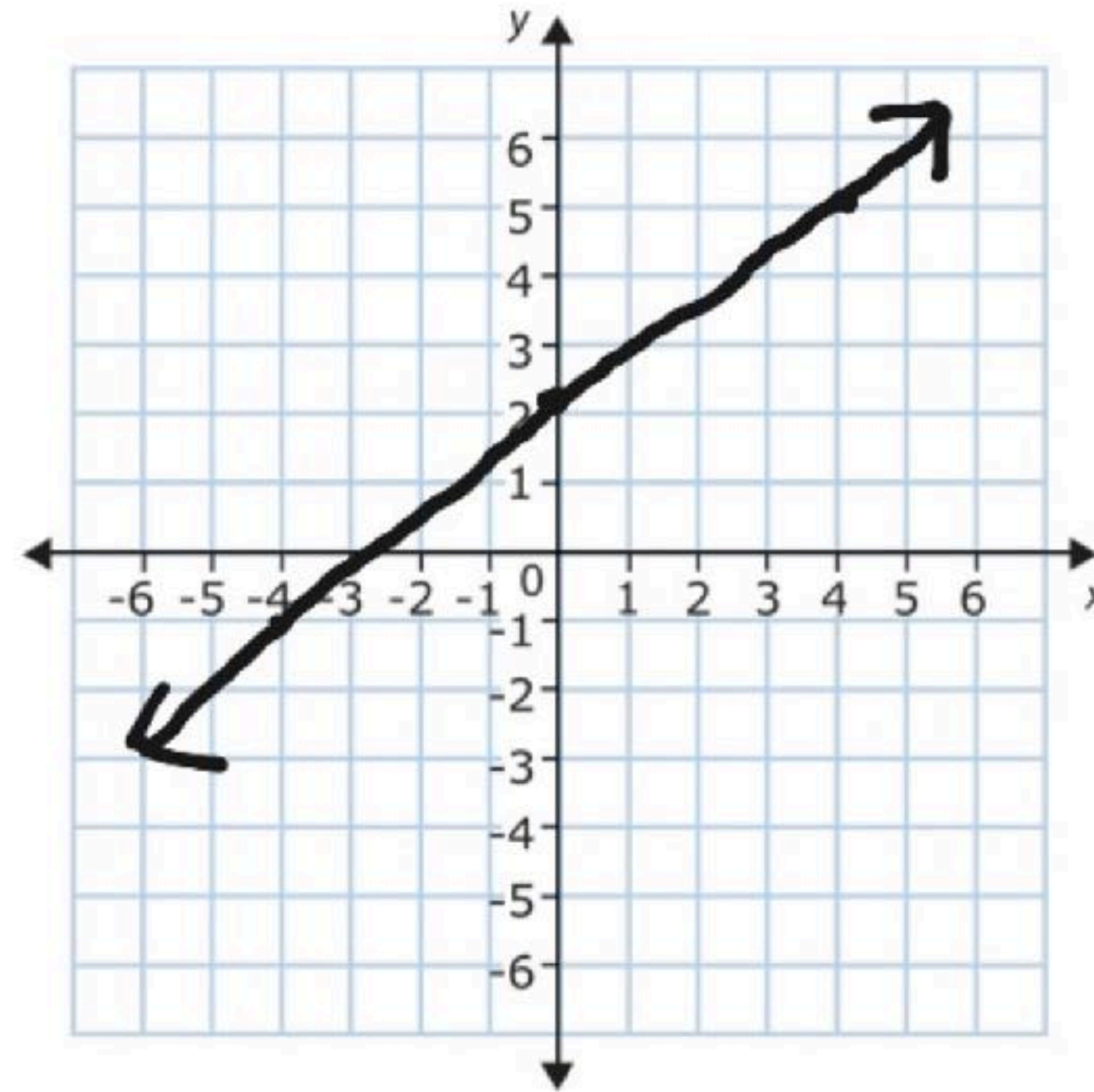


$$y = \frac{3}{4}x + 2$$

Ordinary Least Squares (OLS)

$$y = mx + b$$

slope *intercept*

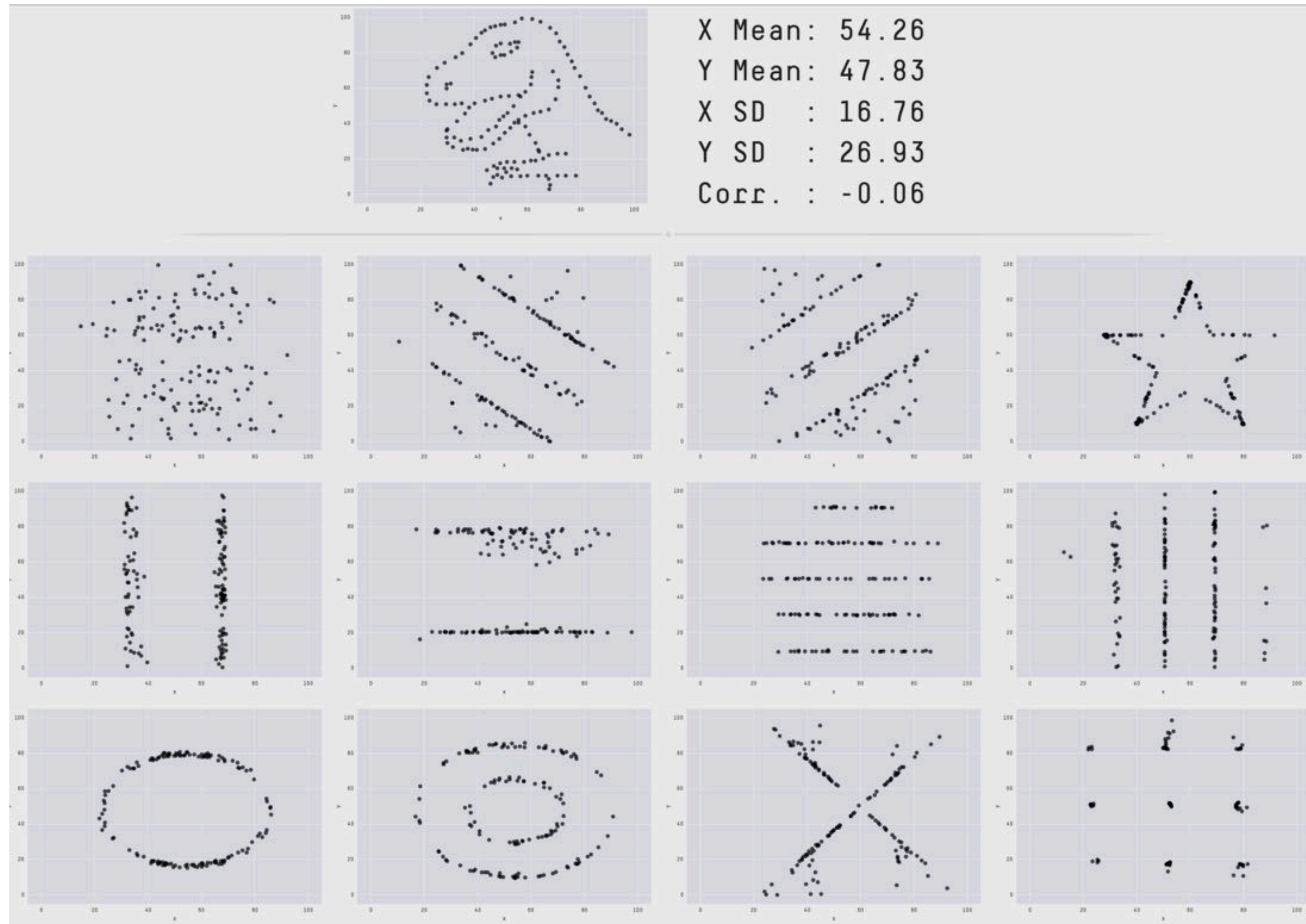


$$y = \frac{3}{4}x + 2$$

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

intercept *slope* *error*

Visualize your data!



Visualize your data!

Same Stats, Different Graphs:

Generating Datasets with Varied Appearance and Identical Statistics through Simulated Annealing



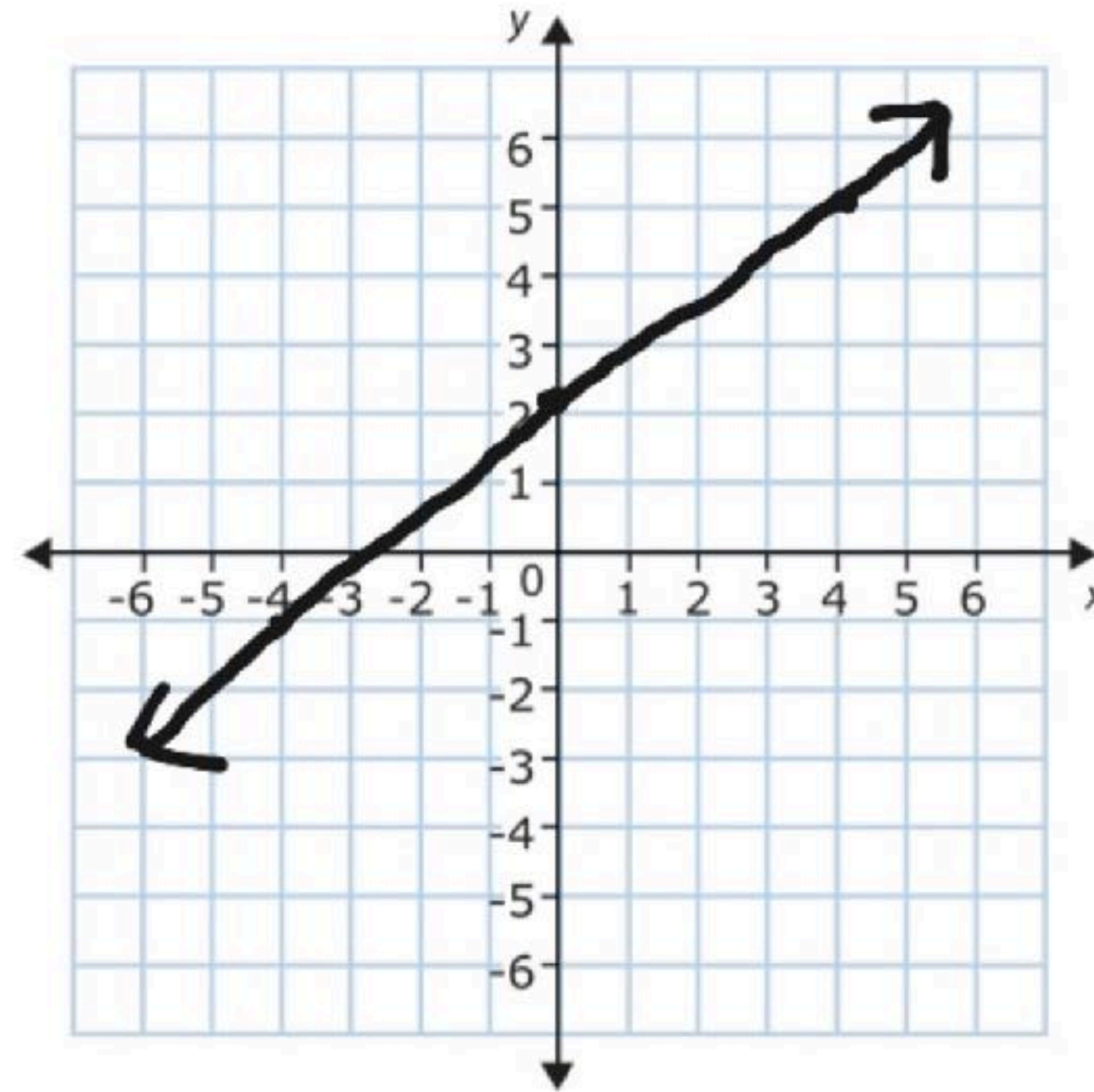
COGS 108

Data Science in Practice

Distributions and outliers

Ordinary Least Squares (OLS)

$$y = \overset{\text{slope}}{mx} + \underset{\text{intercept}}{b}$$



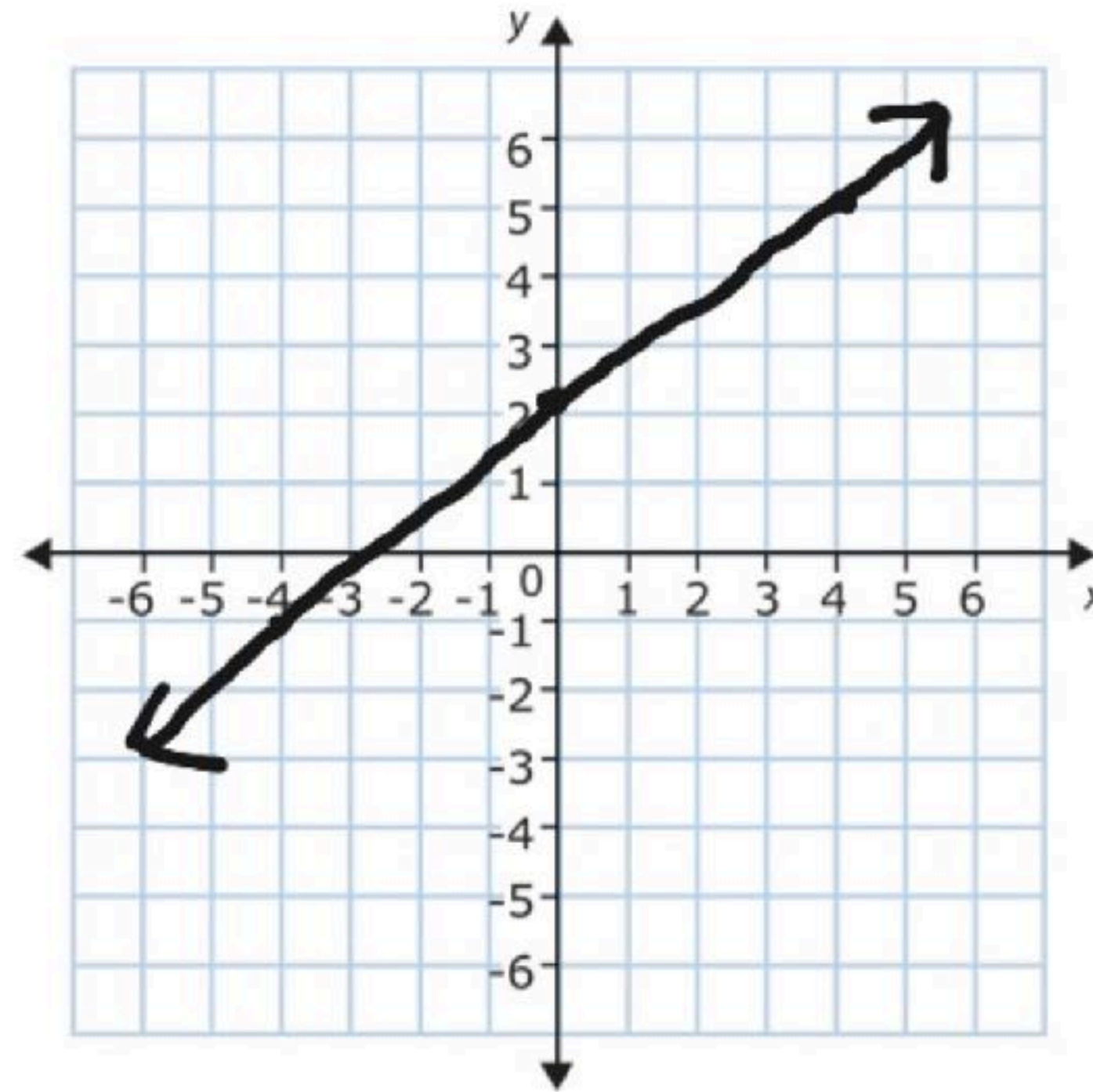
$$y = \frac{3}{4}x + 2$$

$$y_i = \underset{\text{slope}}{\beta_1} x_i + \underset{\text{intercept}}{\beta_0} + \overset{\text{error}}{\varepsilon_i}$$

Ordinary Least Squares (OLS)

$$y = mx + b$$

slope *intercept*



$$y = \frac{3}{4}x + 2$$

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

intercept *slope* *error*

Ordinary Least Squares (OLS)

$$y_i = \overset{\text{intercept}}{\beta_0} + \underset{\text{slope}}{\beta_1} x_i + \overset{\text{error}}{\varepsilon_i}$$

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

residual
error

slope
parameter

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

parameter
intercept

residual
error

slope
parameter

**Goal: find values for the parameters
to provide a "best" fit for the data**

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

parameter
intercept

residual
error

slope
parameter

**Goal: find values for the parameters
to provide a "best" fit for the data**

how?

Ordinary Least Squares (OLS)

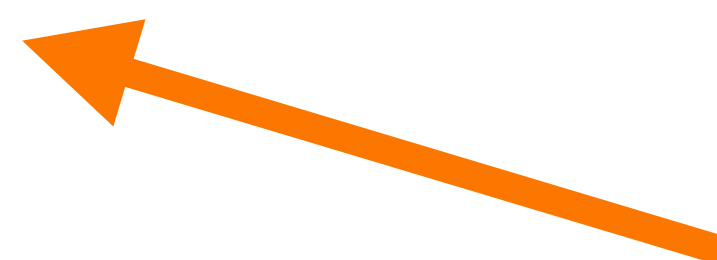
parameter
intercept

residual
error

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

slope
parameter

minimize this!



Ordinary Least Squares (OLS)

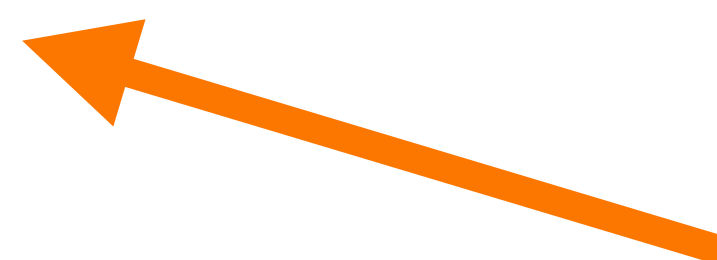
parameter
intercept

residual
error

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

slope
parameter

minimize this!



how?

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

slope
parameter

residual
error

$$\min \sum_{i=1}^n \varepsilon_i^2$$

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

slope
parameter

residual
error

$$\min \sum_{i=1}^n \varepsilon_i^2$$

why squared?

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

slope
parameter

residual
error

$$\min \sum_{i=1}^n \varepsilon_i^2$$

$$y_i - \beta_0 - \beta_1 x_i = \varepsilon_i$$

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

slope
parameter

residual
error

$$\min \sum_{i=1}^n \varepsilon_i^2$$

$$y_i - \beta_0 - \beta_1 x_i = \varepsilon_i$$

$$\min \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

residual
error

slope
parameter

$$\beta_0 = \mu_y - \beta_1 \mu_x$$

$$\beta_1 = r_{xy} \frac{\sigma_y}{\sigma_x}$$

Ordinary Least Squares (OLS)

parameter
intercept

residual
error

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

slope
parameter

**WHAT DOES THIS
MEAN???**

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

residual
error

slope
parameter

$$\beta_0 = \mu_y - \beta_1 \mu_x$$

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

slope
parameter

residual
error

$$\beta_0 = \mu_y - \beta_1 \mu_x$$

**when we see the average value of x ,
we predict the average value of y**

Ordinary Least Squares (OLS)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

parameter
intercept

slope
parameter

residual
error

$$\beta_0 = \mu_y - \beta_1 \mu_x$$

**when we see the average value of x ,
we predict the average value of y**

$$\beta_1 = r_{xy} \frac{\sigma_y}{\sigma_x}$$

Ordinary Least Squares (OLS)

parameter
intercept

residual
error

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

slope
parameter

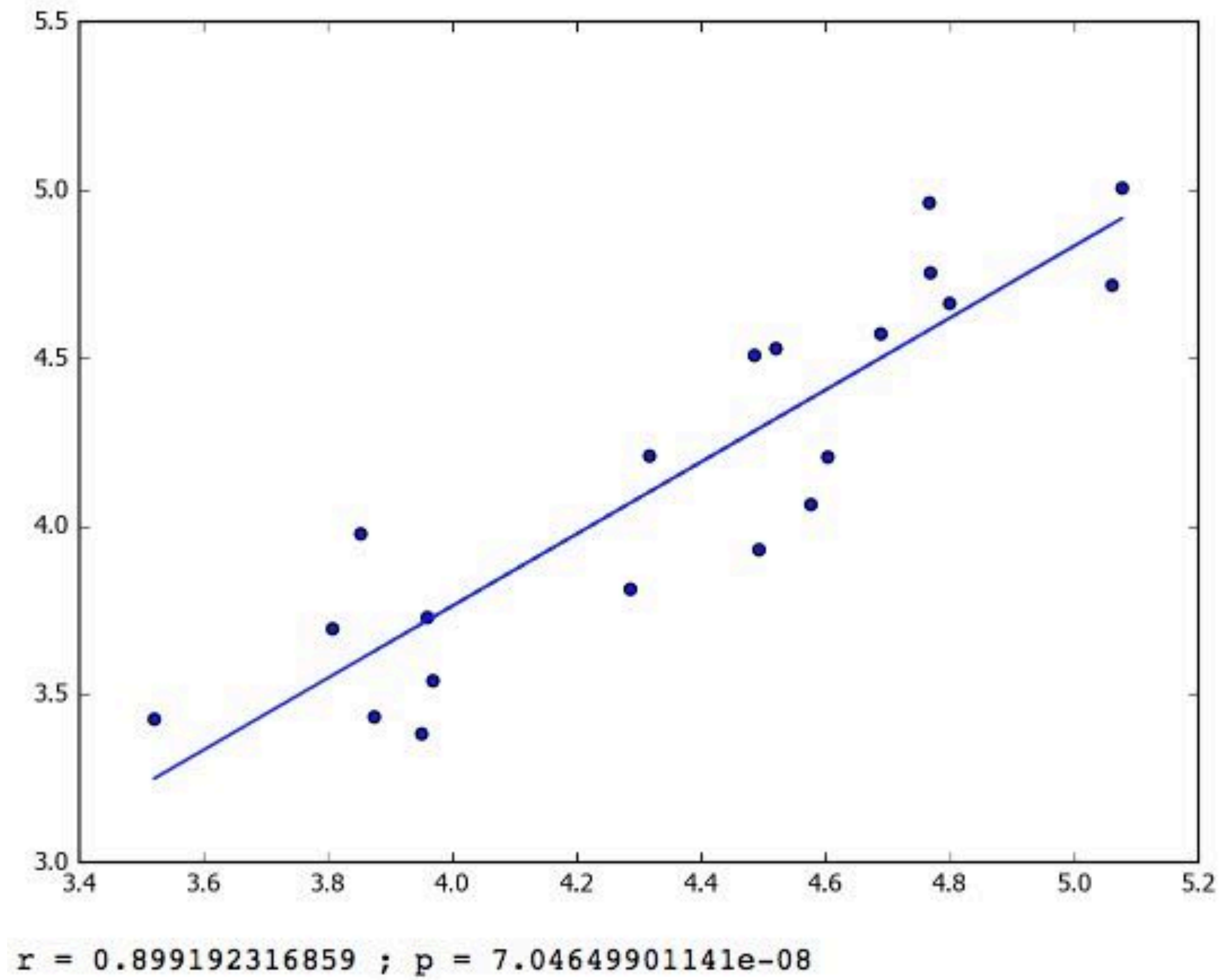
$$\beta_0 = \mu_y - \beta_1 \mu_x$$

**when we see the average value of x ,
we predict the average value of y**

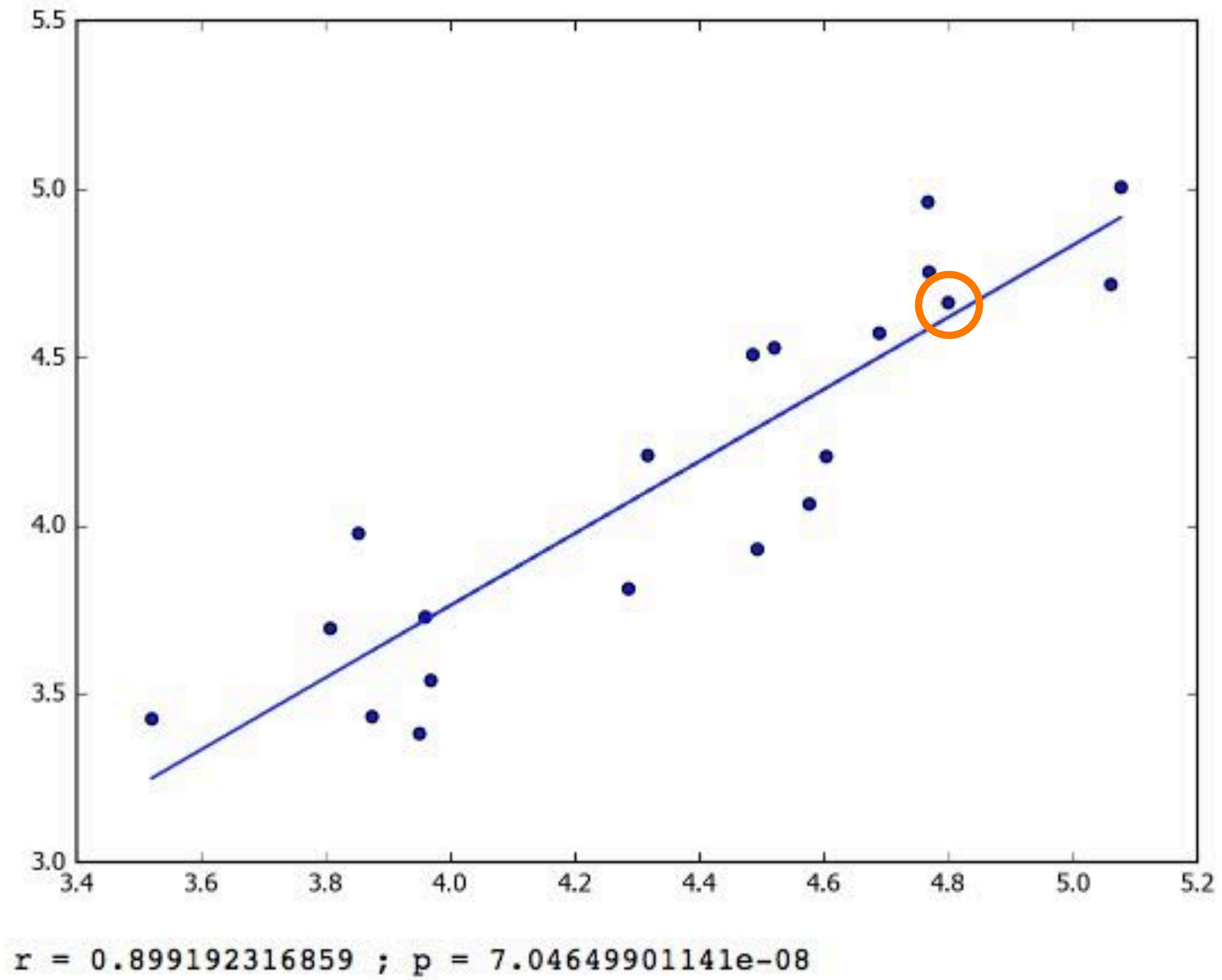
$$\beta_1 = r_{xy} \frac{\sigma_y}{\sigma_x}$$

**when the input increases by σ_x , the
prediction increases by the x,y correlation**

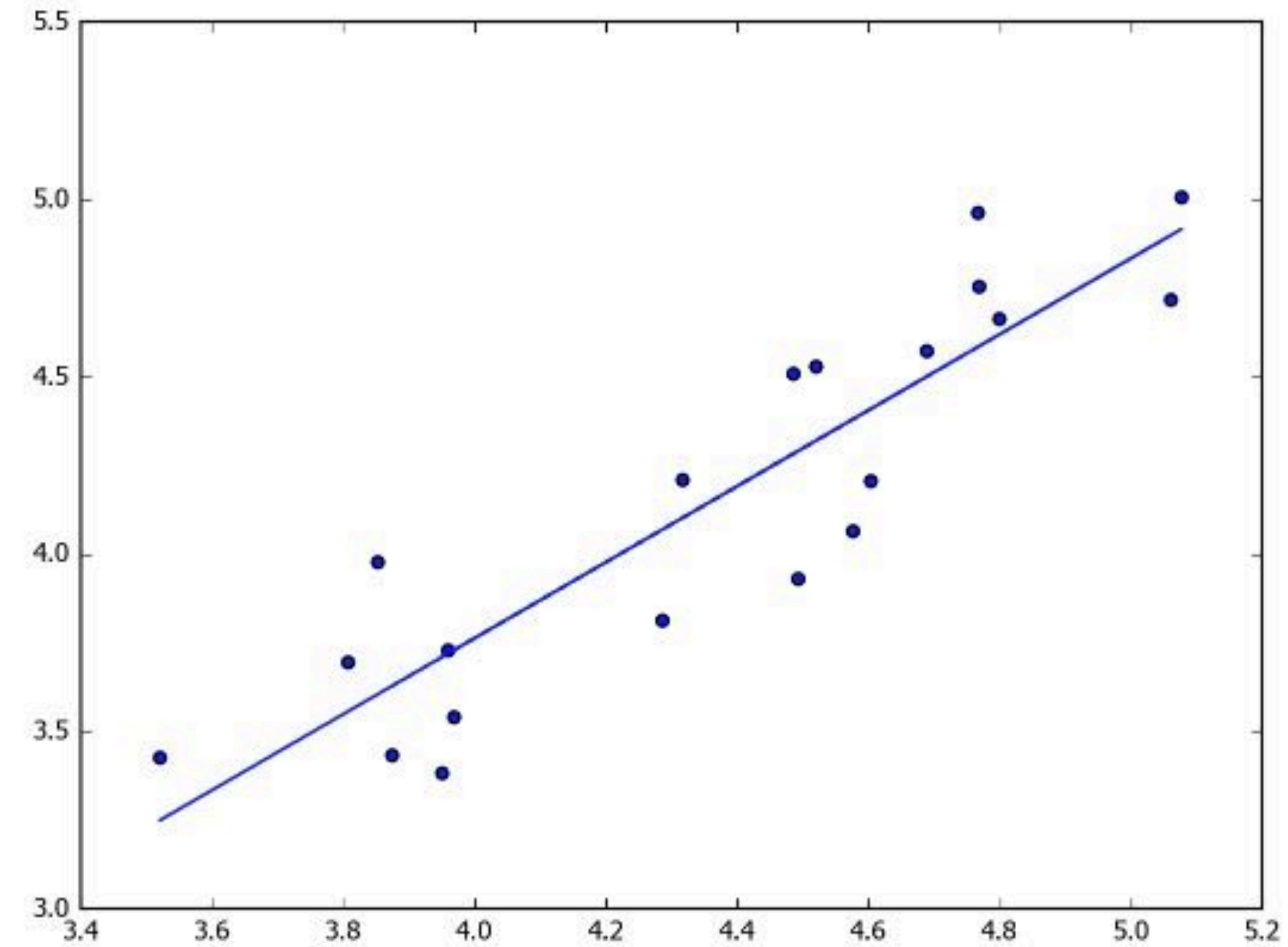
Outliers!



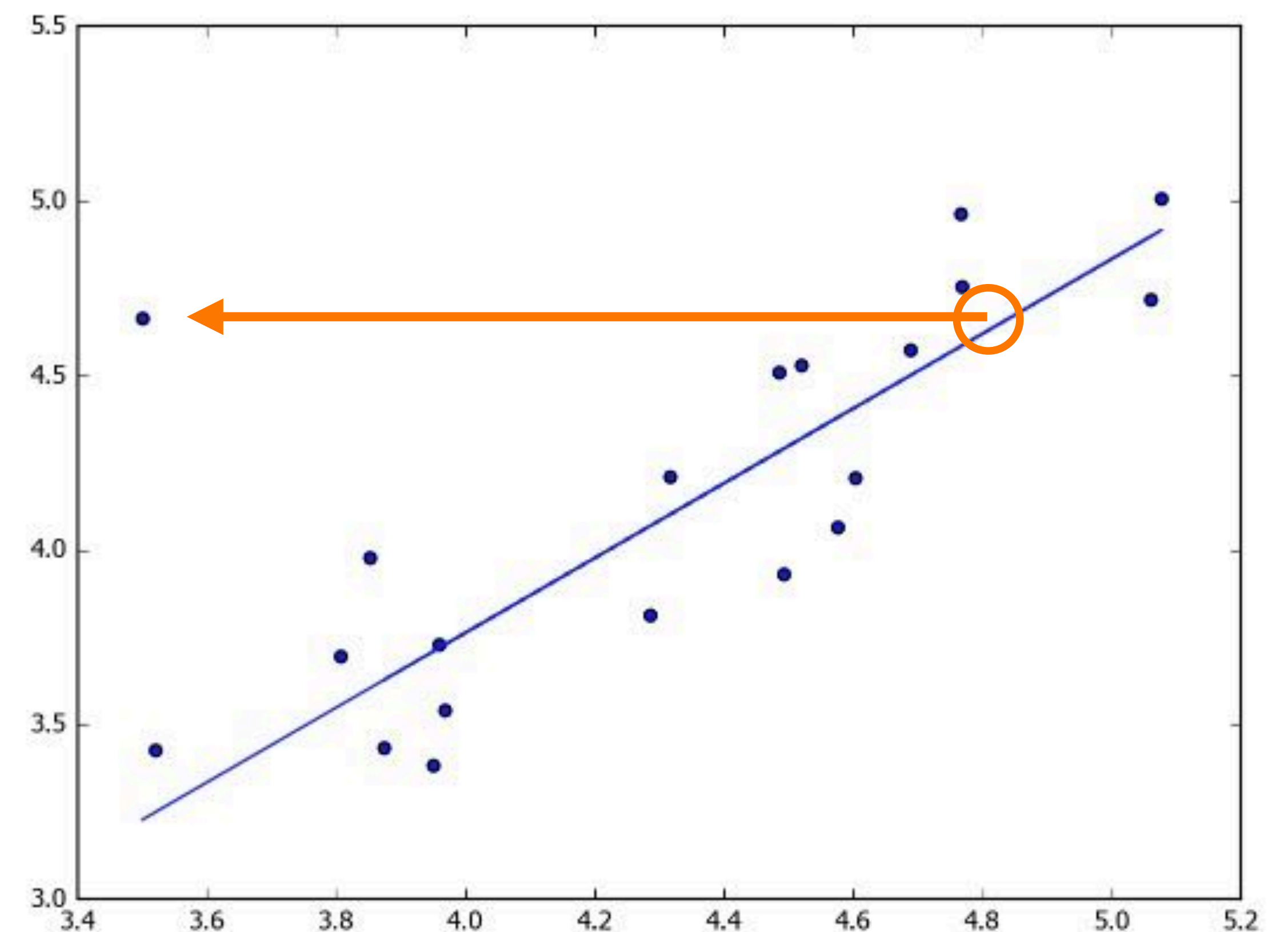
Outliers!



Outliers!

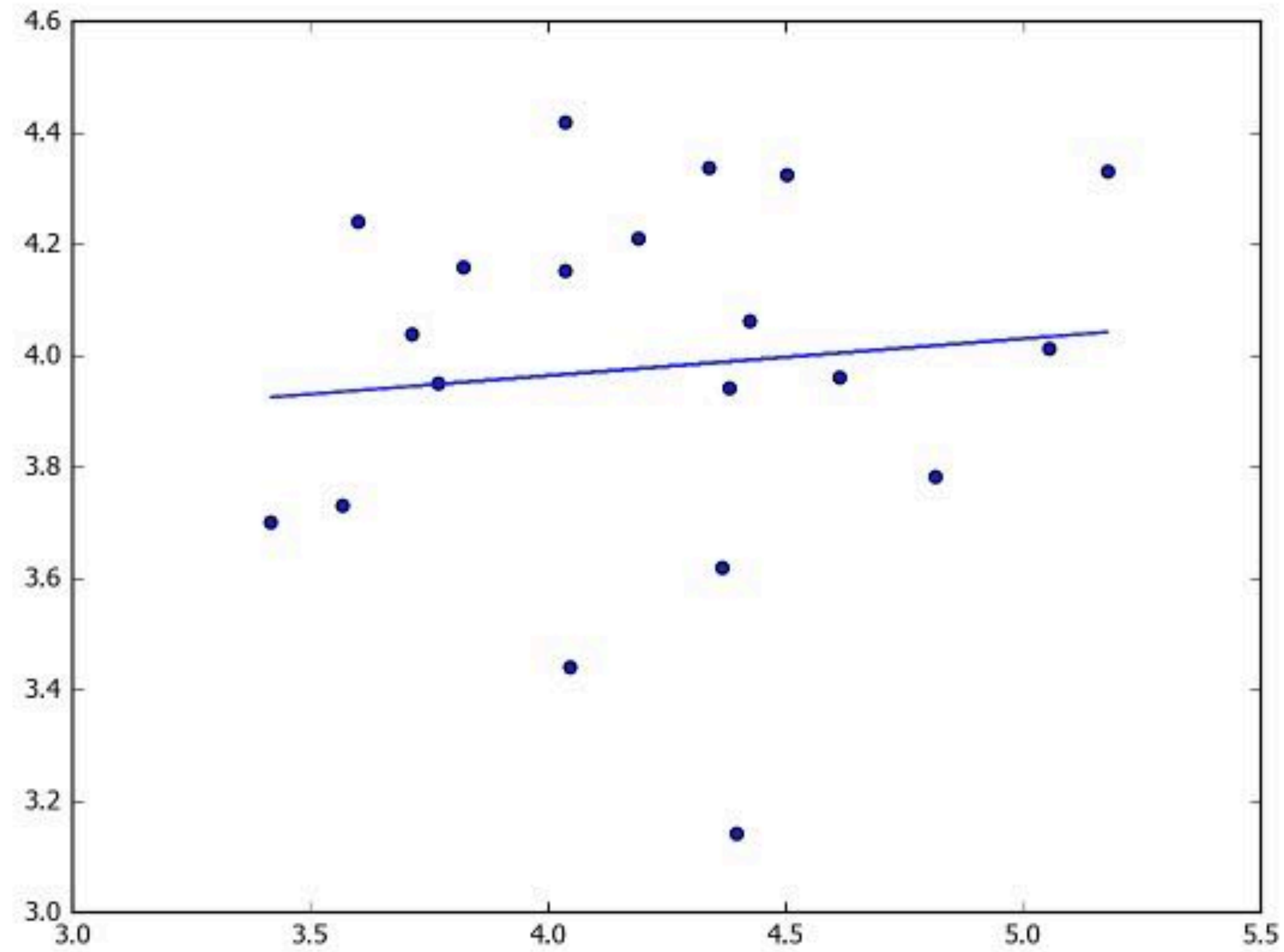


$r = 0.899192316859$; $p = 7.04649901141e-08$



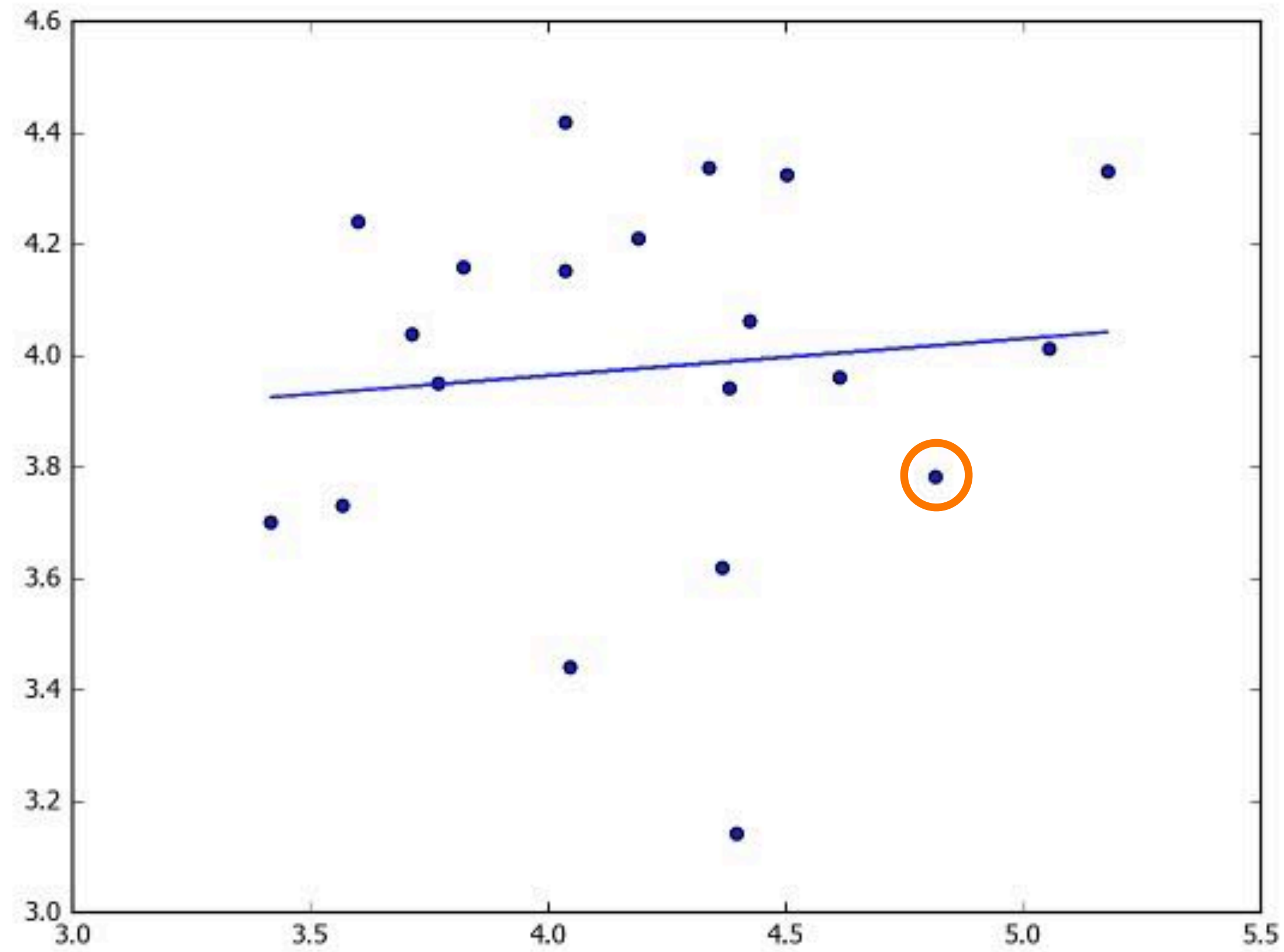
0.709892176011 0.000454396182138

Outliers!



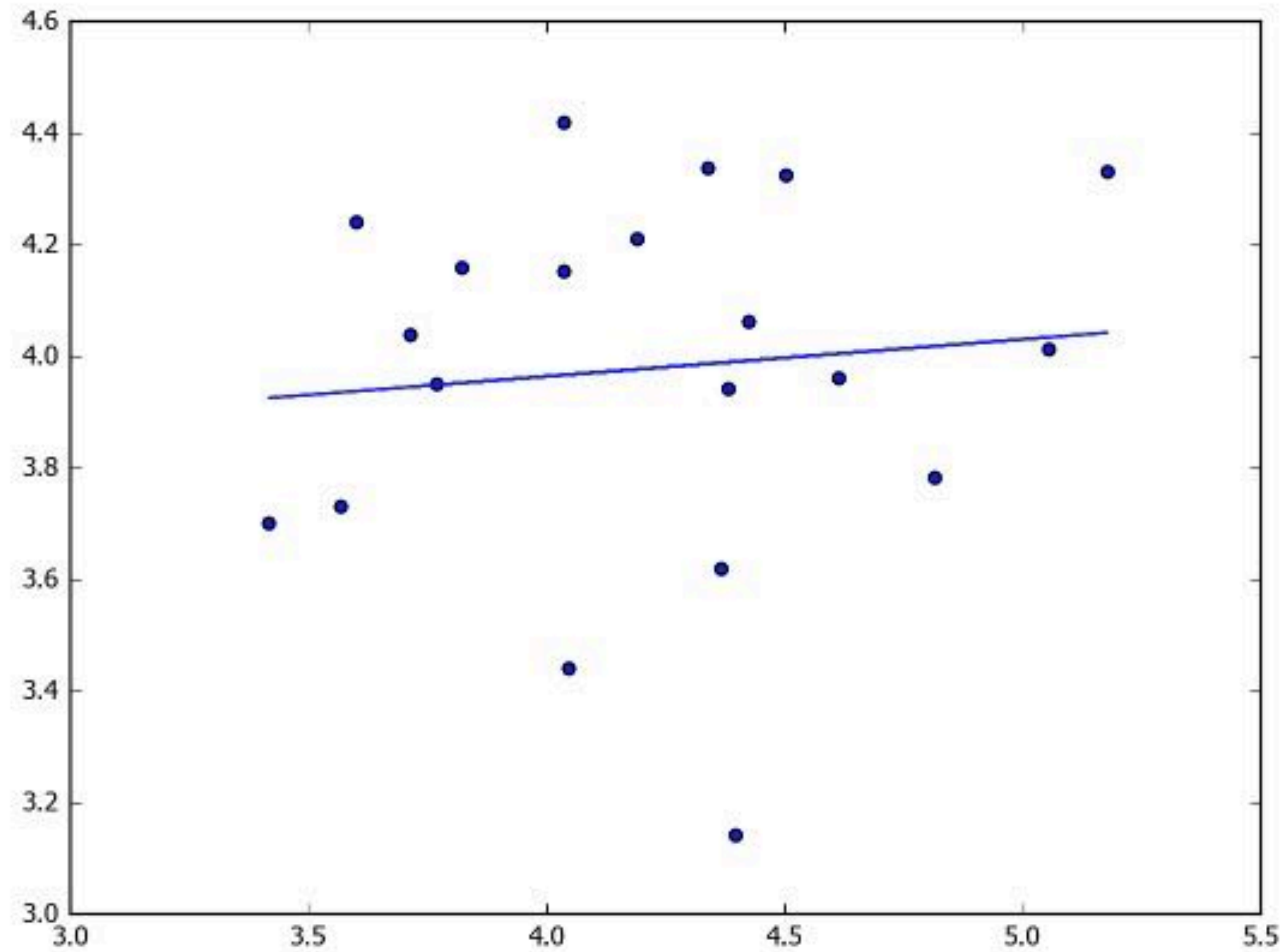
$r = 0.0982029696446$; $p = 0.680415082529$

Outliers!

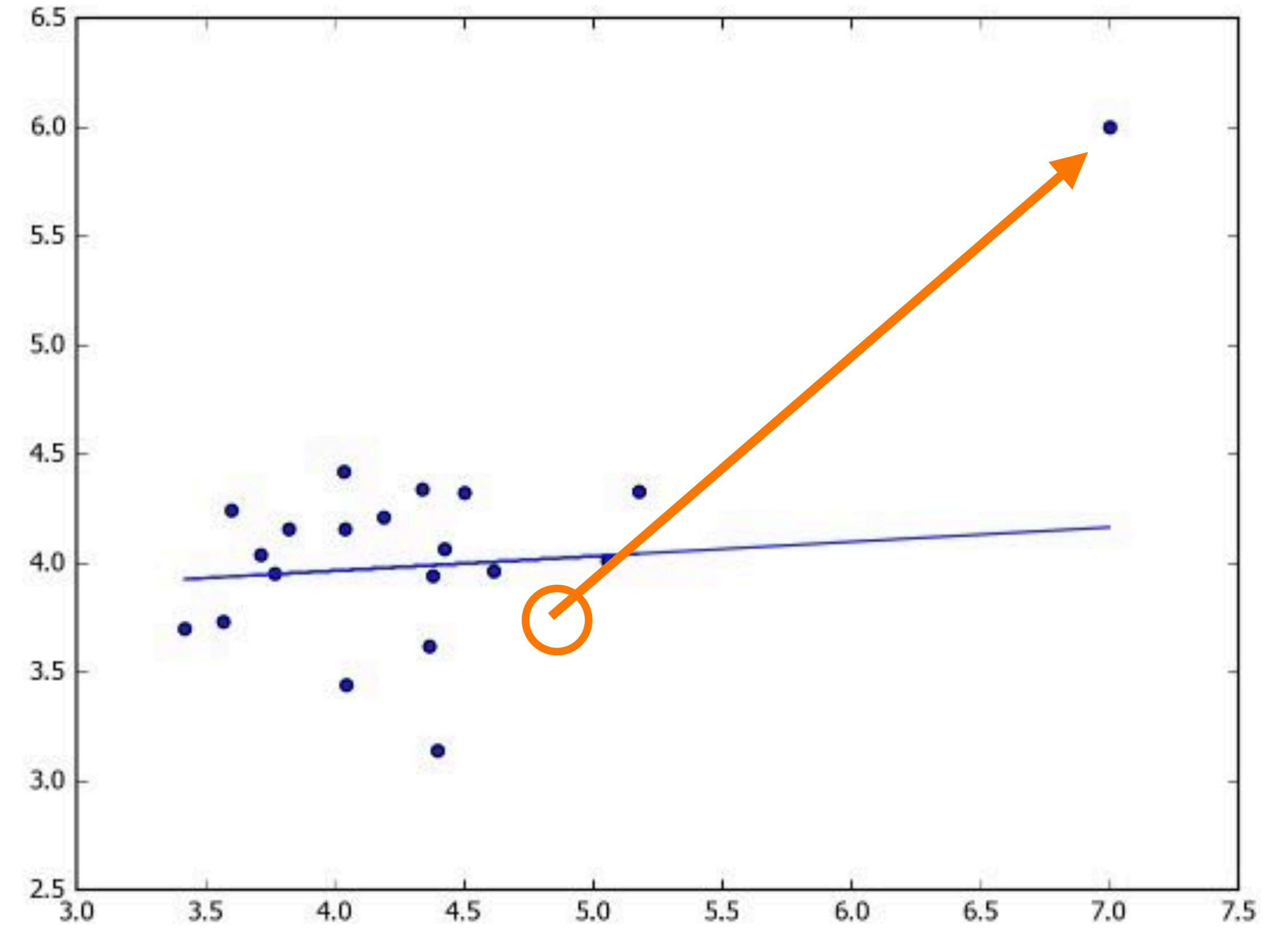


$r = 0.0982029696446$; $p = 0.680415082529$

Outliers!

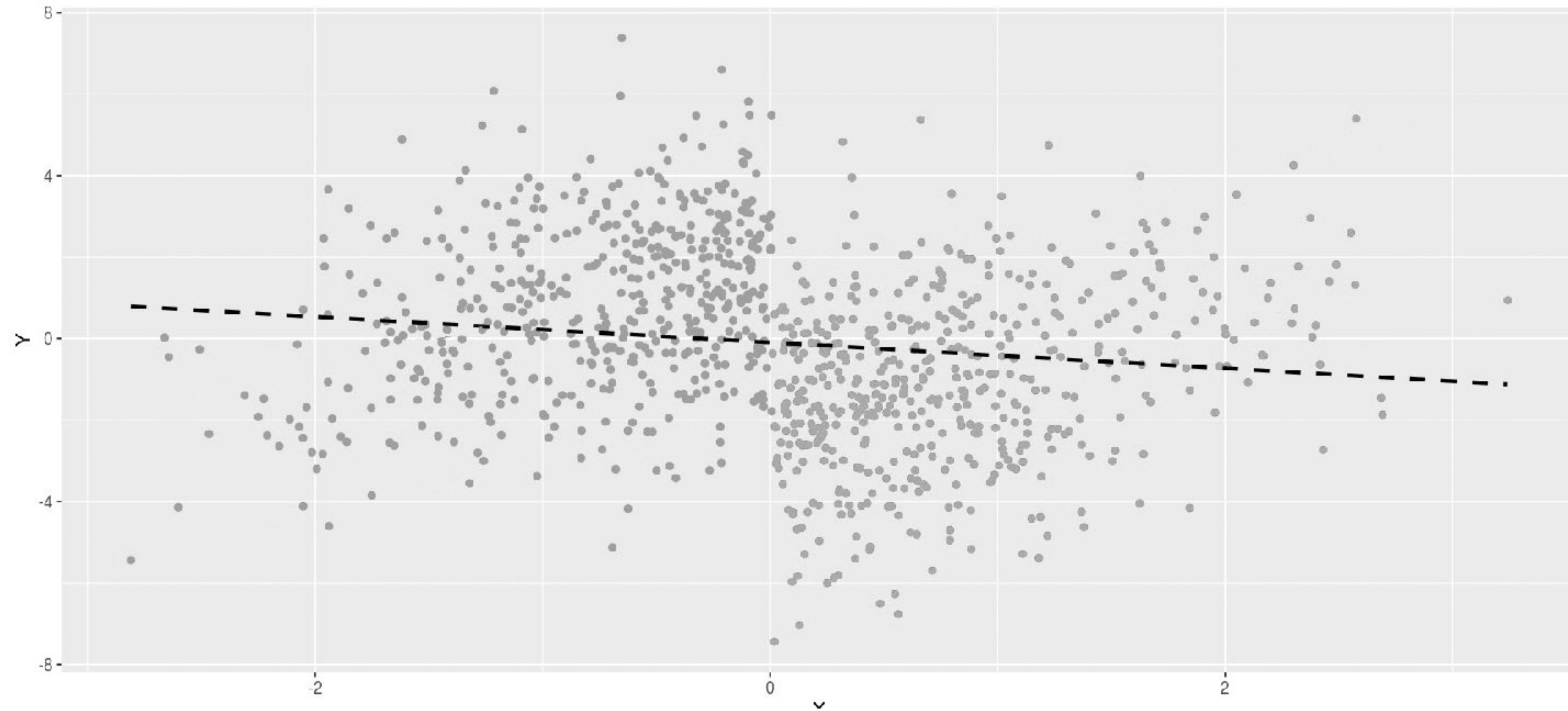


$r = 0.0982029696446$; $p = 0.680415082529$

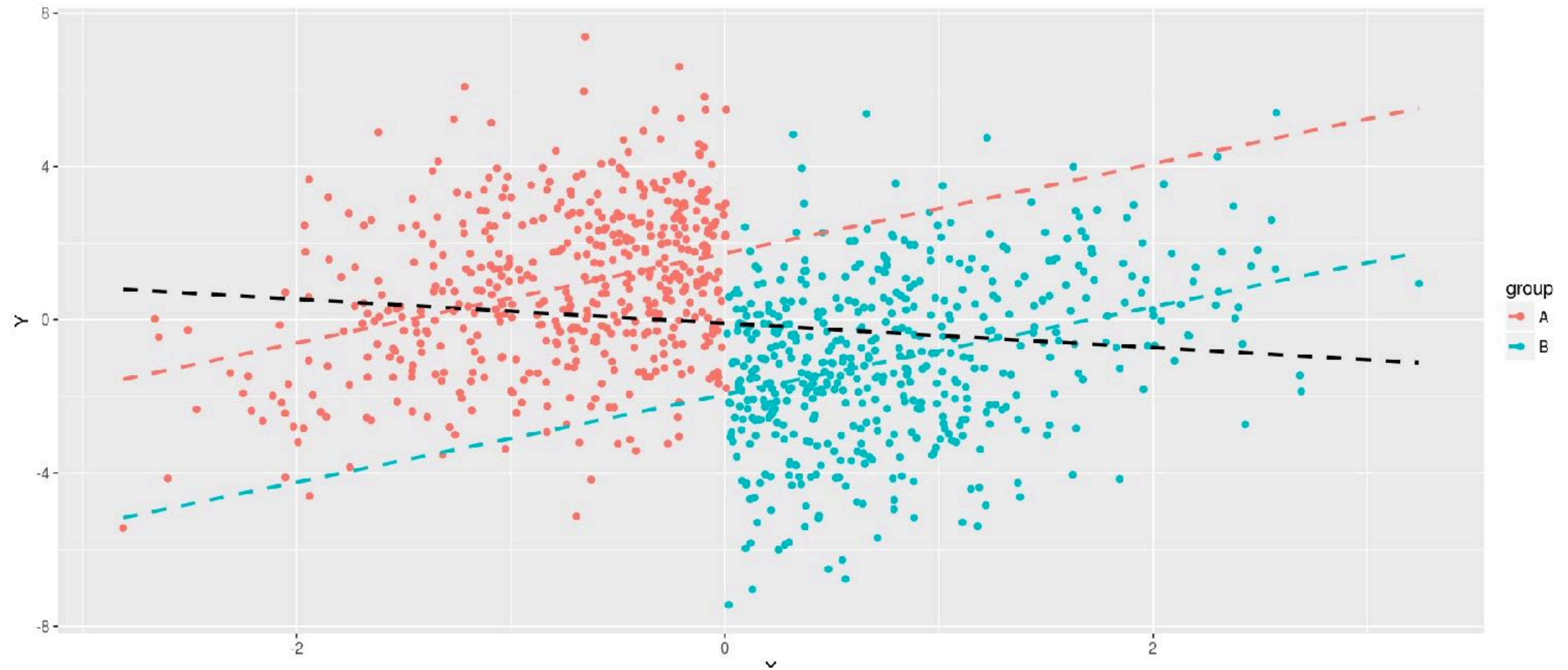


$r = 0.702033103659$; $p = 0.000559678419262$

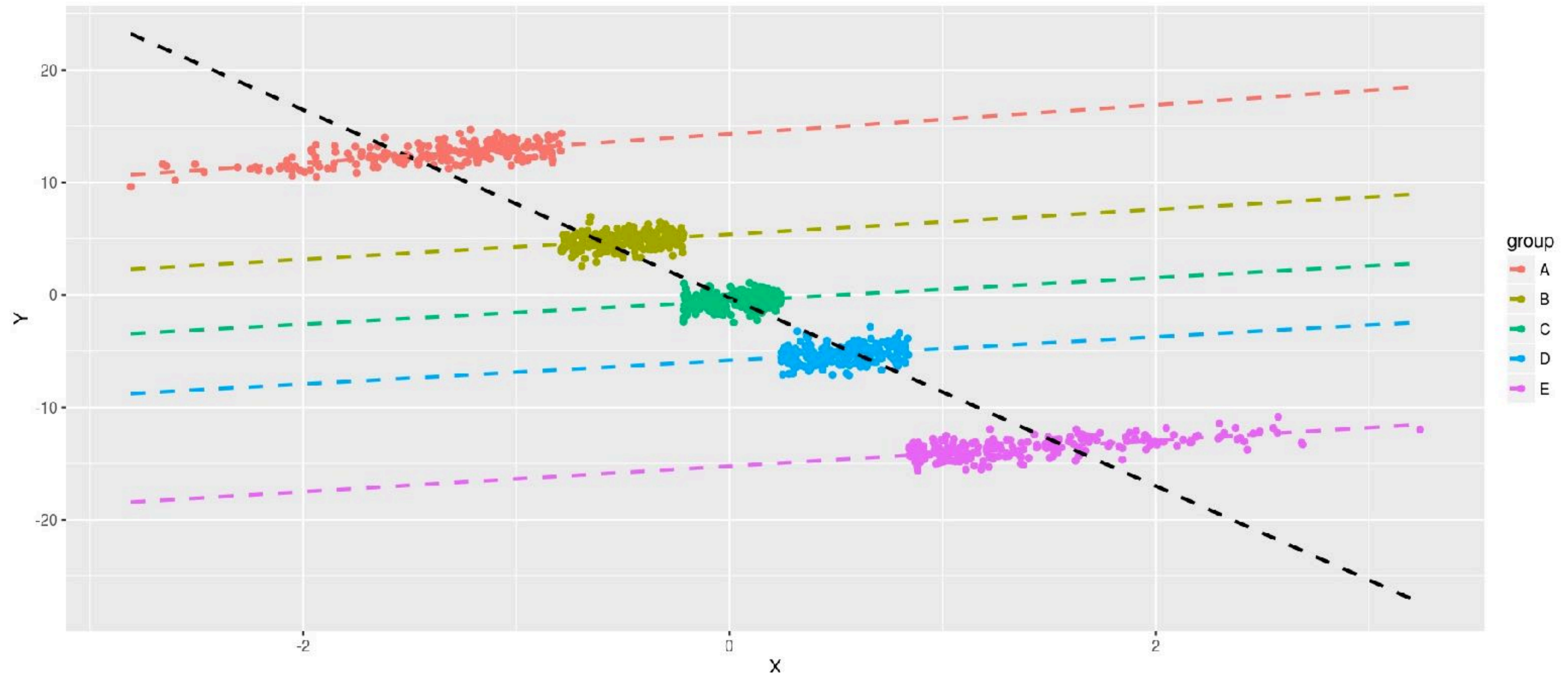
Simpson's Paradox



Simpson's Paradox



Simpson's Paradox



Simpson's Paradox - Real world example

Example 2—Airlines on-time data: Here are numbers of flights on time and delayed for two airlines at five airports in June 1991. The table shows that Alaska Airlines outperforms America West at all 5 cities. If you collapse the table over city, it appears that America West outperforms Alaska.

	Alaska Airlines			America West Airlines		
	On time	Delayed	Delay%	On time	Delayed	Delay%
LA	497	62	11.1%	694	117	14.4%
Phoenix	221	12	5.4%	4840	415	7.9%
San Diego	212	20	8.6%	383	65	14.5%
San Fran.	503	102	16.9%	320	129	28.7%
Seattle	1841	305	14.2%	201	61	23.3%
Total	3274	501	13.3%	6438	787	10.9%

Simpson's Paradox - Intuition

- **Totally contrived example:**
 - You and a friend are both avid gamers

Simpson's Paradox - Intuition

- **Totally contrived example:**
 - You and a friend are both avid gamers
 - You are also both kind of competitive

Simpson's Paradox - Intuition

- **Totally contrived example:**
 - You and a friend are both avid gamers
 - You are also both kind of competitive
 - You make a friendly wager: you'll both play a question-answering game twice, and whomever has the best record, wins

Simpson's Paradox - Intuition

- **Totally contrived example:**
 - You and a friend are both avid gamers
 - You are also both kind of competitive
 - You make a friendly wager: you'll both play a question-answering game twice, and whomever has the best record, wins
- **Day 1:**
 - you: 98/99 (98.99%)
- **Day 1:**
 - them: 1/1 (**100%**)

Simpson's Paradox - Intuition

- **Totally contrived example:**
 - You and a friend are both avid gamers
 - You are also both kind of competitive
 - You make a friendly wager: you'll both play a question-answering game twice, and whomever has the best record, wins
- **Day 1:**
 - you: 98/99 (98.99%)
- **Day 2:**
 - you: 0/1 (0%)
- **Day 1:**
 - them: 1/1 (**100%**)
- **Day 2:**
 - them: 1/99 (**1.01%**)

Simpson's Paradox - Intuition

- **Totally contrived example:**
 - You and a friend are both avid gamers
 - You are also both kind of competitive
 - You make a friendly wager: you'll both play a question-answering game twice, and whomever has the best record, wins
- **Day 1:**
 - you: 98/99 (98.99%)
- **Day 2:**
 - you: 0/1 (0%)
- **Total:**
 - you: 98/100 (**98%**)
- **Day 1:**
 - them: 1/1 (**100%**)
- **Day 2:**
 - them: 1/99 (**1.01%**)
- **Total:**
 - them: 2/100 (2%)

Bradley Voytek, Ph.D.
UC San Diego

Department of Cognitive Science
Neurosciences Graduate Program
Halicioglu Data Science Institute

bvoytek@ucsd.edu
@bradleyvoytek

