

# Limits of statistical method

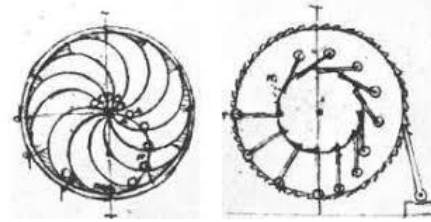
**Petr Paščenko**

**6. 1. 2022**

# Motivation: Limits of knowledge

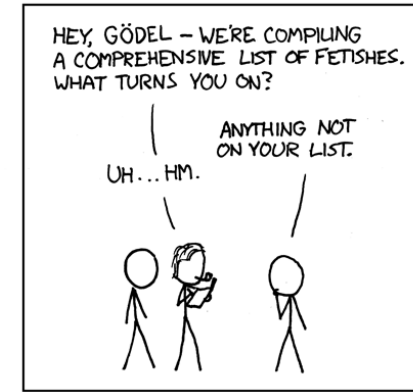
*It is good to know the limits...*

- › Second law of thermodynamics
  - *Heat does not spontaneously flow from a colder body to a hotter.*
- › Gödel's incompleteness
  - *No all truths are provable. Turing's halting problem*
- › Heisenberg uncertainty principle
  - *the position and the velocity of an object cannot both be measured exactly, at the same time, even in theory*
- › Many others
  - speed of light
  - P is not ? NP
  - ...



AUTHOR KATHARINE GATES RECENTLY ATTEMPTED TO MAKE A CHART OF ALL SEXUAL FETISHES.

LITTLE DID SHE KNOW THAT RUSSELL AND WHITEHEAD HAD ALREADY FAILED AT THIS SAME TASK.



# Osnova

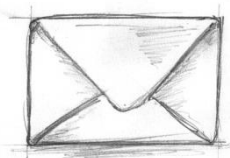
1. Mandelbrot 7 states of randomness and why it matters
2. Central limit theorem assumptions
3. Correlation and causation
4. Statistical paradoxes

# **Mandelbrot seven states of randomness and why it matters (a lot)**

# 10 most important days

PROFINIT

Imagine, you delete 10 most important days from your life...

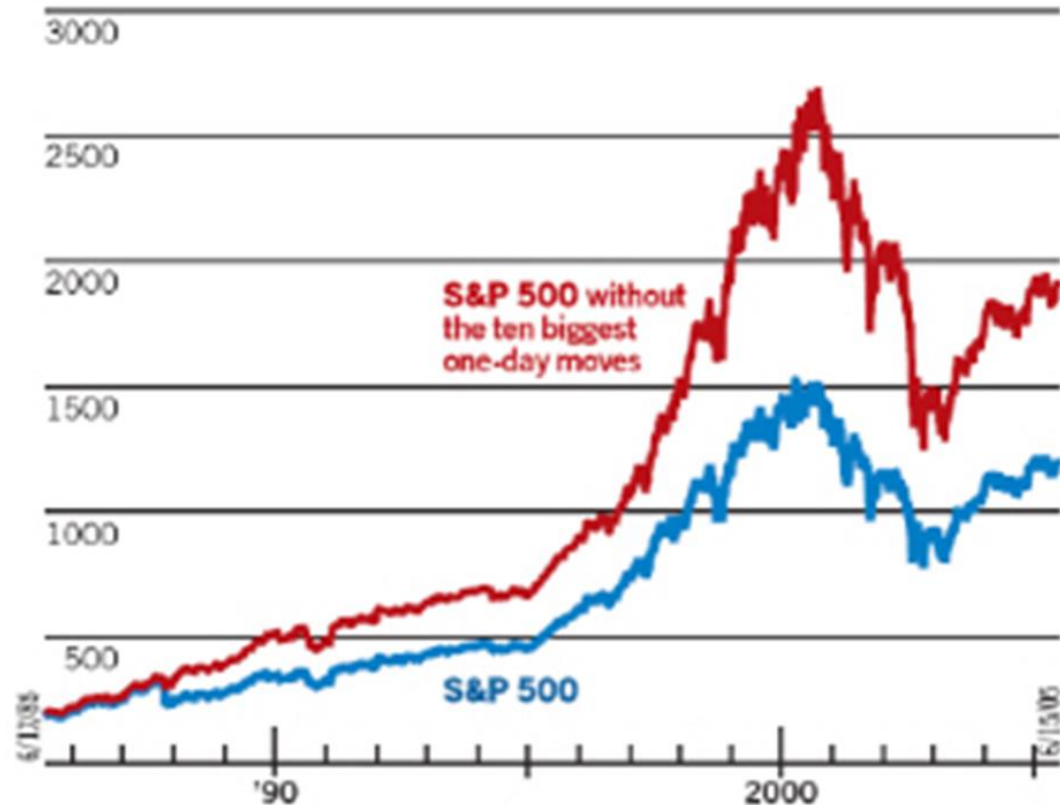


# 10 most important days

PROFITIT

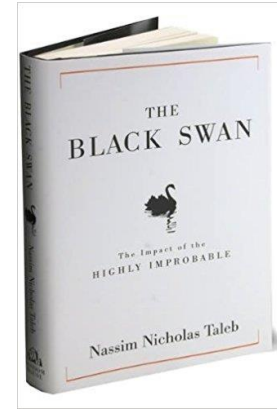
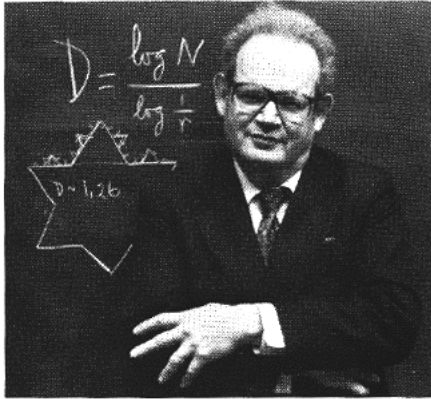
Imagine, you delete 10 most important days from the stock market.

- › out of 20 years
  - i.e. 4000 business days
  - i.e.  $< 0.25\%$

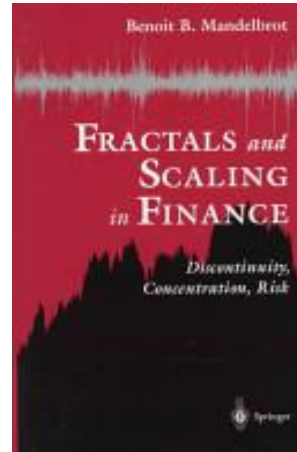


# Benoit Mandelbrot and Nassim Taleb

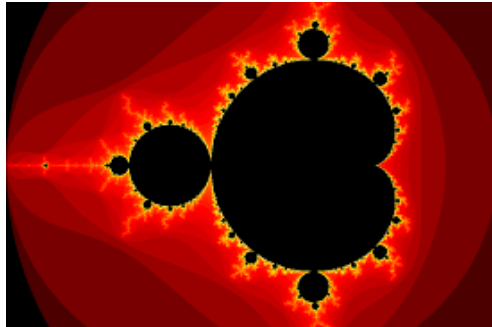
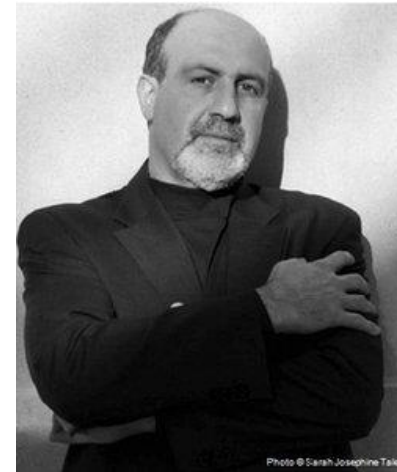
PROFITIT



2010

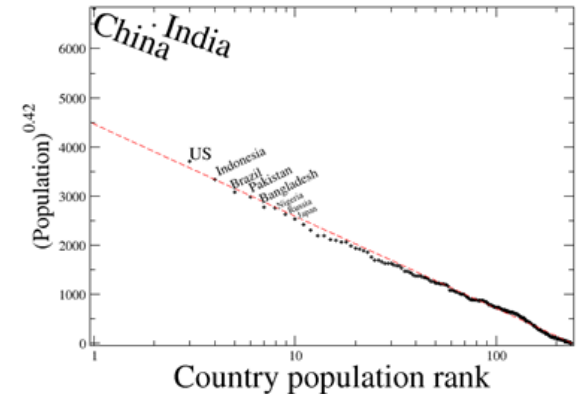
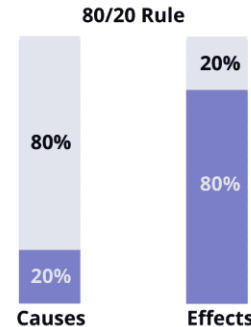
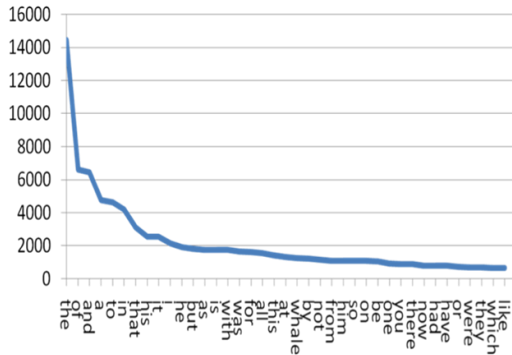
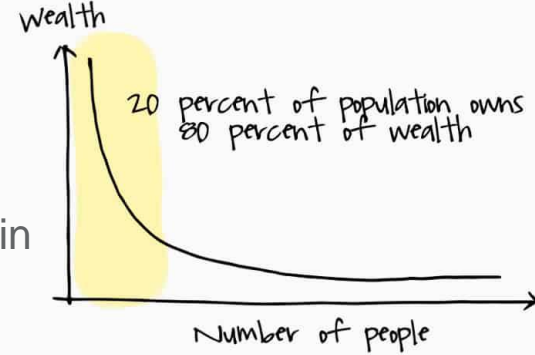


1997



# Pareto law and other empirical observations

- › **Pareto principle** (1890 – economy)
  - 80% of Italy's land was owned by 20% of the population
- › **Zipf's law** (1935 – mathematical linguistics)
  - the frequency of any word is inversely proportional to its rank in the frequency table
- › **Jackson's law**: size of human settlements





# Taleb: fable of two worlds

PROFINIT

## Mediocristan

- › Take 1000 random people on a stadium, sort them all by **weight**
- › Calculate average value in each group  
**75 kg**
- › Add a single most **heavy** / **wealthy** person on the planet  
**200 kg**
- › How does the average changed?

**75,1 kg**

## Extremistan

**wealth**

**\$ 120k**

**\$ 131T**

**\$ 131G**

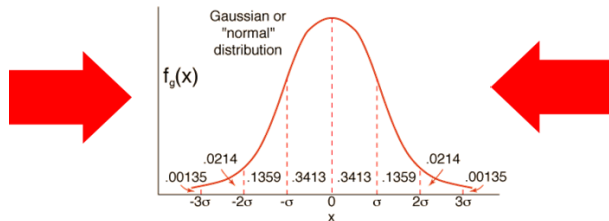


99,9999%

# The source proces on the background

## Mediocristan

- › Evolutionary search for optimum
  - size, weight, height, etc.
  - natural panalization of extremis
  - negative feedback loop
- › Aggregation
  - Central limit theorem
  - covergence to Normal distribution



## Extremistan

- › Winner takes all
  - join the winner
  - positive feedback loop
    - popularity, capital, gravity,...
- › Matthew effect
  - *'For unto every one that have shall be given, and he shall have abundance; but him that have not shall be taken, even that which he have.'* Matthew 25:29

# Key properties

## Mediocristan

Does not scale (dentist)

Physical limits

Physical measures (height)

Gaussian randomness

Typical is close to average

Winner takes a small piece

Common in history

Black swan robust

Extremes can be neglected

Easy to comprehend

Easy to predict

Slow gradual changes, continuity

## Extremistan

Does scale (google)

No limits

A number (wealth)

Power law (Pareto) randomness

No typical, no average

Winner takes (almost) all

Common in current era

Black swan vulnerable

Extremes is what matters

Tricky to comprehend

Impossible to predict

Phase changes, discontinuities

# Mandelbrot: Seven states of randomness

## Key concepts

- › **even portioning vs. concentration portioning**
  - Having N random addends from a distribution, are they of the same order of magnitude?
  - In other words, is maximum major portion of the sum?
- › **scale factor of order q**  $\alpha_q = E|(X)^q|$ 
  - root of degree q of a q-th moment
  - **finite** or **infinite** moment

# Mandelbrot: Seven states of randomness

PROFITIT

## › Mild randomness (long term even portioning, all moments finite)

### 1. Proper mild randomness (normal distribution)

- Even portioning for  $N = 2$ .

### 2. Borderline mild randomness (exponential)

- Short term concentrated portioning, long term even portioning

## › Slow randomness (long term concentrated portioning, all moments finite)

### 3. Slow randomness with finite delocalized moments

### 4. Slow randomness with finite and localized moments (lognormal)

## › Wild randomness

### 5. Pre-wild randomness (pareto $\alpha > 2$ )

- infinite moments for  $q > 2$

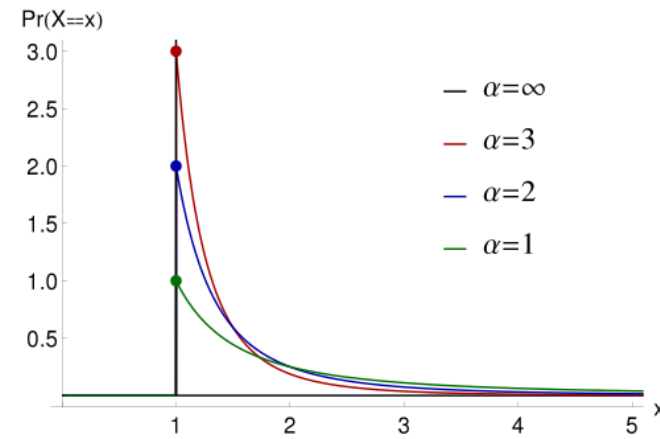
### 6. Wild randomness (pareto $\alpha \leq 2$ )

- infinite variance, i.e. non convergent sample variance

### 7. Extreme randomness (pareto $\alpha \leq 1$ )

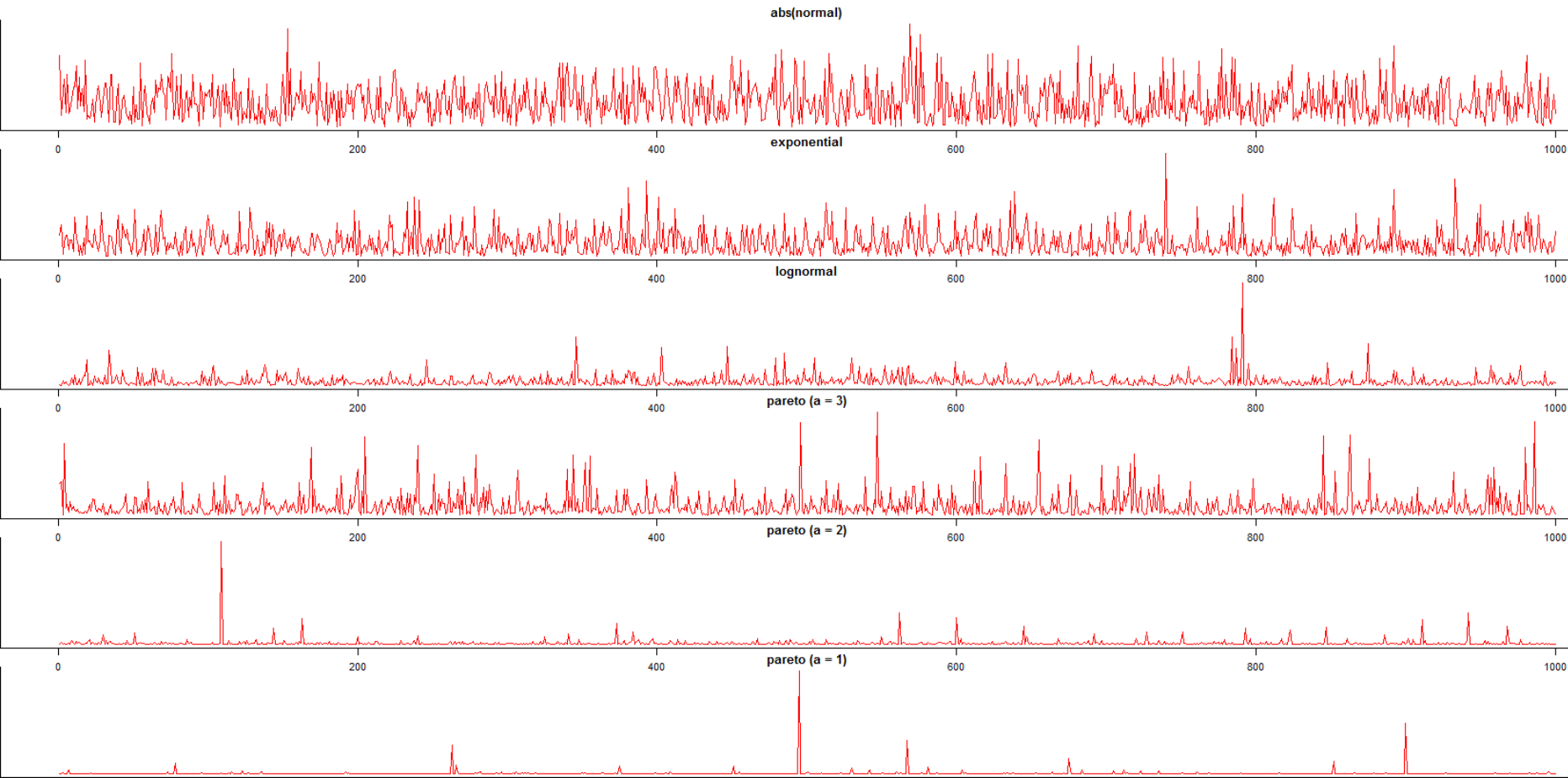
- infinite mean, i. e. non convergent sample mean

$$\frac{\alpha x_m^\alpha}{x^{\alpha+1}}$$



# Mandelbrot: Seven states of randomness

PROFINIT



# Pareto distribution, empirical parameter estimation

PROFINIT

Variable	Alpha
Word usage frequency	1,2
WWW visits per page (before FB)	1,4
Book title sell numbers	1,5
Earthquake magnitude	2,8
Moon crater size	2,14
Sun corona eruption sizes	0,8
War intensity	0,8
American citizen wealth	1,1
Surname frequency	1
Market movements	3 or less?
City sizes	1,3
Corporation sizes	1,5
Terrorist attack death counts	2

# Consequences of wild randomness

- › Statistical inference does not work
  - we can not infer the parameters of distributions from data
- › Central limit theorem does not work
  - we can not reduce the uncertainty by aggregation
- › Prediction does not work
  - our forecast is systematically underestimated
  - our confidence interval is underestimated as well
- › Black swan events
  - Unpredictable large scale events with usually negative consequences
  - Natural disasters, market crashes, political crises, epidemics, etc.
  - We can only prepare for foreseeable catastrophes



# Central limit theorem assumptions

# Central limit theorem and its assumptions

$$\lim_{n \rightarrow \infty} \frac{\sum_{i=1}^N X_i - n\mu}{\sqrt{\sigma^2 n}} \sim N(0,1)$$

## › Assumptions

### 1. X has finite mean and variance

### 2. X is iid

- **independent**

- random variables  $X_1 \dots X_n$  are independent on each other
- coins vs. sheeps

- **identically distributed**

- $X_1 \dots X_n$  are chosen from the same probabilistic distribution
- there is no phase change or any other discontinuity in the process
- almost never satisfied in practice
- stability of model testing etc.

# Example: local retail bank in a small town

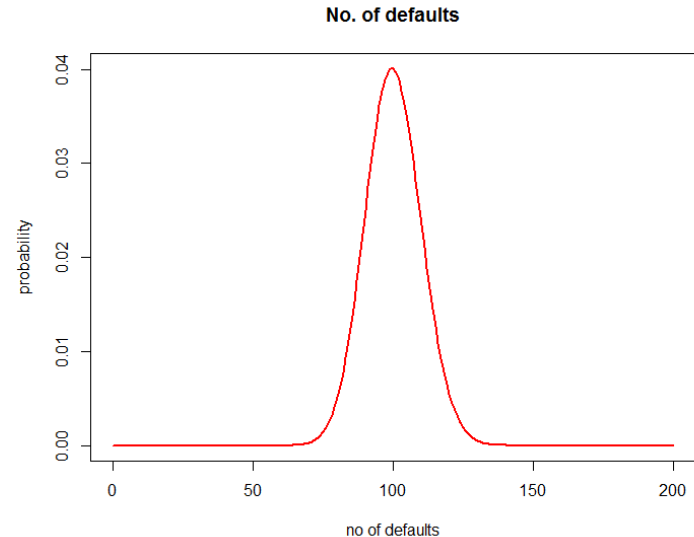
PROFITIT

- › A retail bank
  - 50k people, 10k mortgages, \$250k each
  - Priori probability of default is 1%
  - What is my expected worst case loss (68%, 98%, 99,9%)?

## › Binomial distribution

- $p_0 = 0.01$
- $EX = p_0 \cdot N = 0.01 \cdot 10\,000 = 100$
- $sd(X) = \sqrt{N \cdot p_0 \cdot (1 - p_0)} \cong 10$
- Number of defaults in worst case:

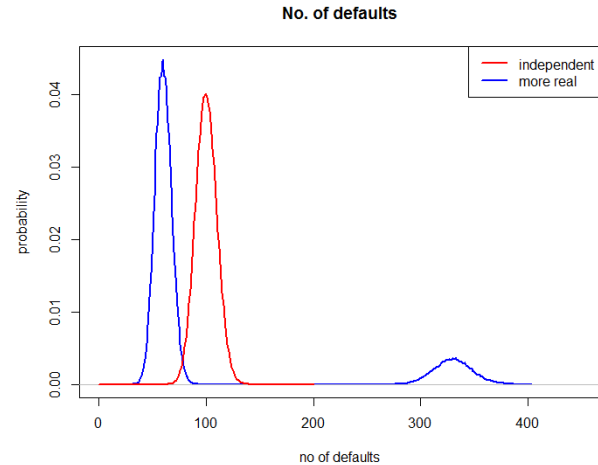
Probability	68%	98%	99,9%
Worst Case	110	120	130



# Example: local retail bank in a small town

- › Small city
  - half of the people work for 1 factory
  - probability of bankruptcy 15%

Situation	Number	Probability	Default
banc., fac.	5000	15%	5%
banc., non fac.	5000	15%	1,6%
non banc., fac.	5000	85%	0,4%
non banc., non fac.	5000	85%	0.8%

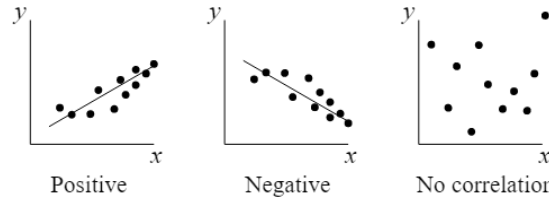


–  $EX = 5000 \cdot (0.15 \cdot (0.05 + 0.016)) + 0.15 \cdot (0.04 + 0.08) = 100$

Probability	68%	98%	99,9%
Worst Case independent	110	120	130
Worst Case real	66	357	375

# Statistical paradoxes

# Correlation and Causation

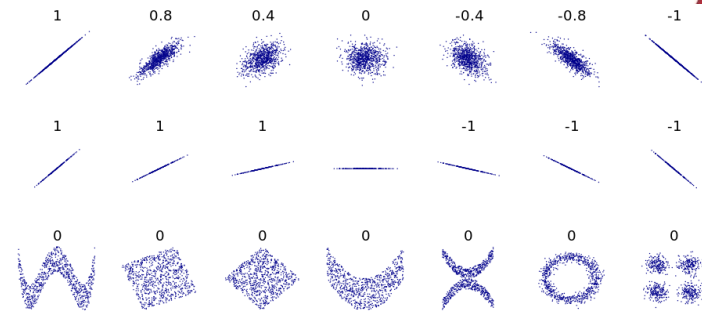


- > Correlation
  - linear dependency of variables: A and B
- > Causation = any form of dependence
  - visiting lectures implies passing the exam
- > Correlation does not imply causation

$$\rho_{X,Y} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E((X - \mu_X)(Y - \mu_Y))}{\sigma_X \sigma_Y},$$



- > Correlation without causation
- > Causation without correlation



# Simpson paradox

- › What is the vaccine effectiveness?

$$e = 1 - \frac{P(S|V)}{P(S|\neg V)}$$

$$e = 1 - \frac{\frac{5,3}{100\,000}}{\frac{16,4}{100\,000}} = 67,5\%$$

Severe cases		Efficacy
Not Vax per 100k	Fully Vax per 100k	vs. severe disease
214 16.4	301 5.3	<b>67.5%</b>

# Simpson paradox

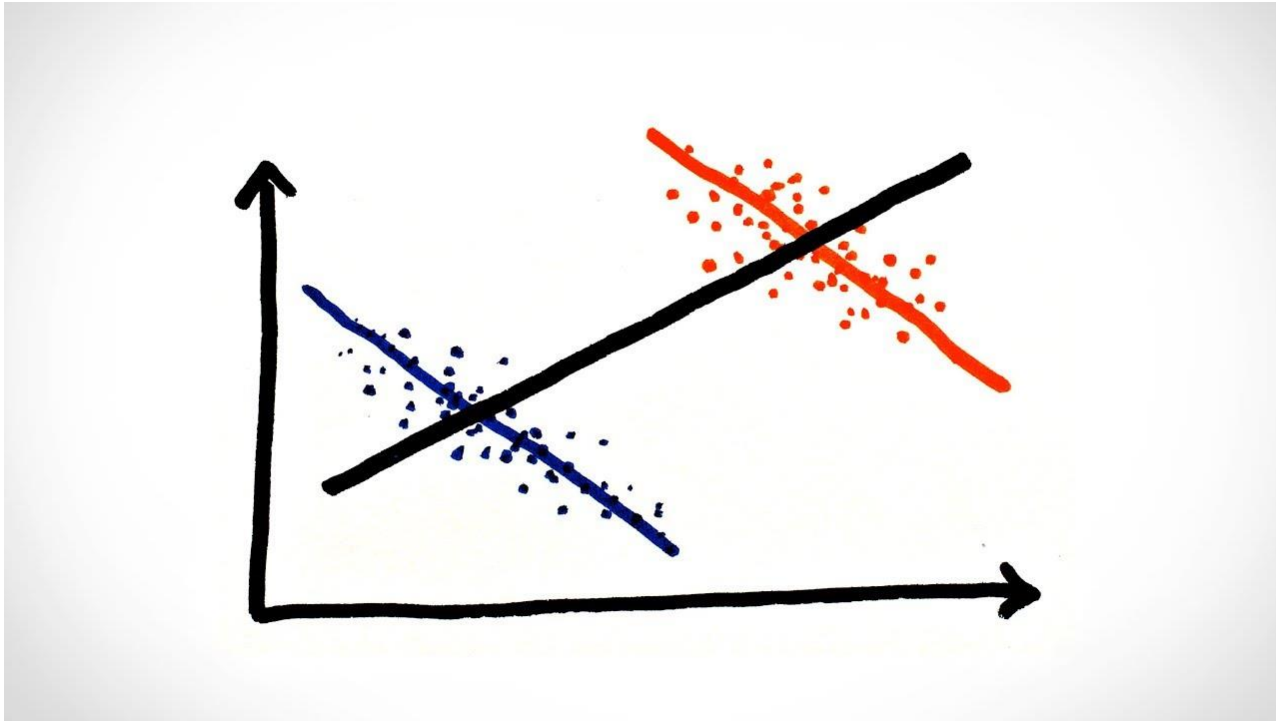
PROFINIT

Age	Population (%)		Severe cases		Efficacy vs. severe disease
	Not Vax %	Fully Vax %	Not Vax per 100k	Fully Vax per 100k	
All ages	1,302,912 <b>18.2%</b>	5,634,634 <b>78.7%</b>	214 <b>16.4</b>	301 <b>5.3</b>	<b>67.5%</b>
<50	1,116,834 <b>23.3%</b>	3,501,118 <b>73.0%</b>	43 <b>3.9</b>	11 <b>0.3</b>	<b>91.8%</b>
>50	186,078 <b>7.9%</b>	2,133,516 <b>90.4%</b>	171 <b>91.9</b>	290 <b>13.6</b>	<b>85.2%</b>

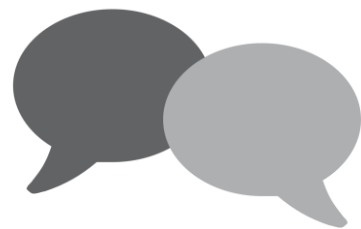
- › The classes are imbalanced
  - in both severe cases and vaccination



# Simpson paradox



- › smoking versus life expectancy for male and female



## Diskuze