# Experimental Data Analysis
*in ©MATLAB*

## Lecture 4:
## Correlations, normality of data testing, parametric vs. non-parametric tests

Jan Rusz
Czech Technical University in Prague
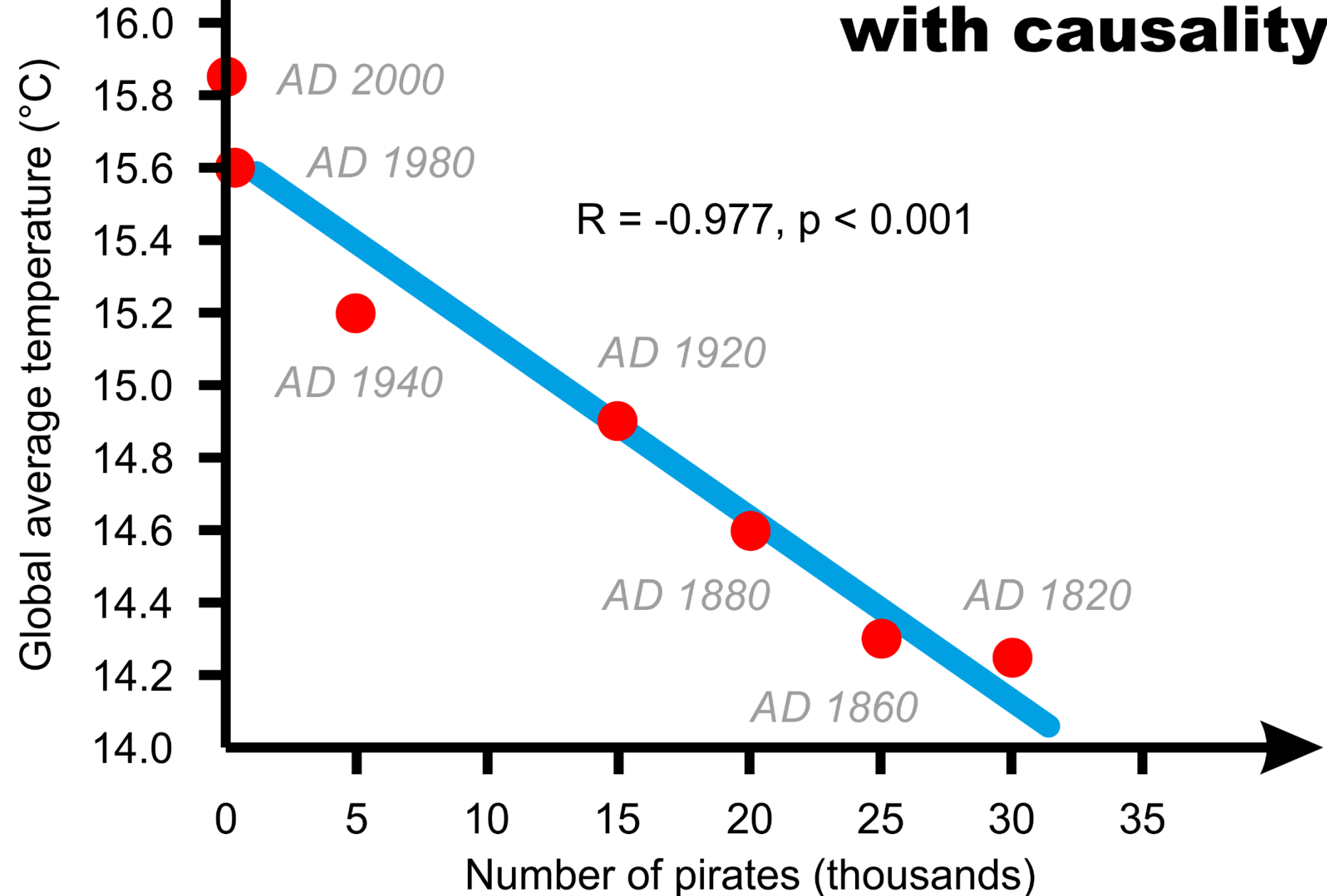
## Motivation

### Association
Question: Can be increased blood pressure associated with stress?
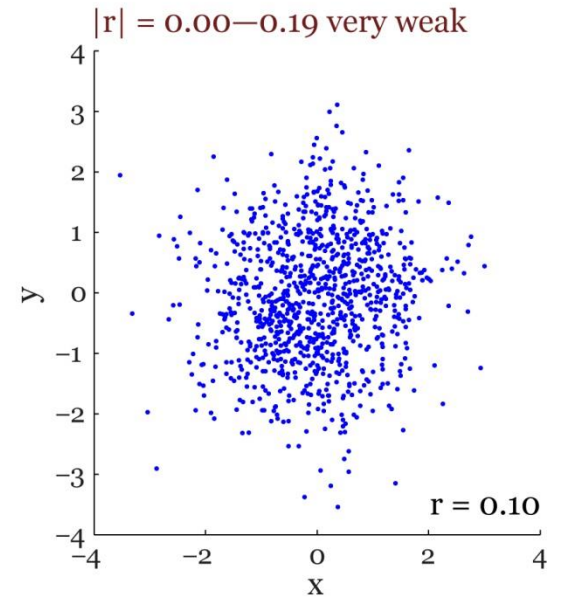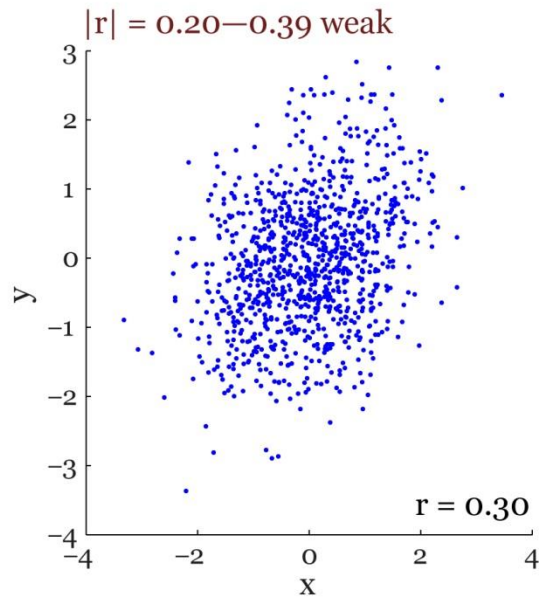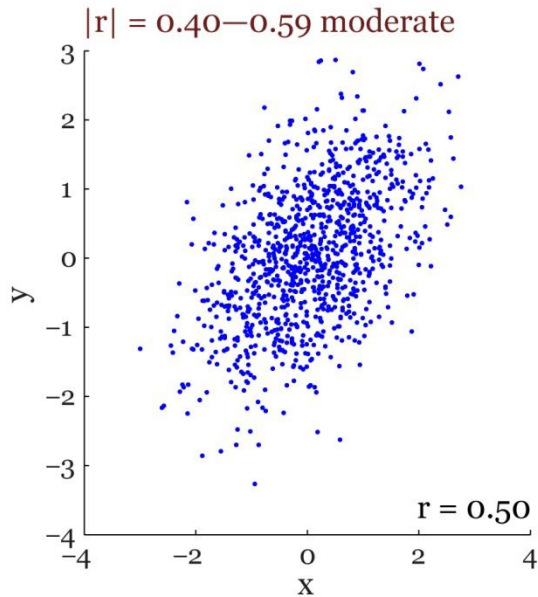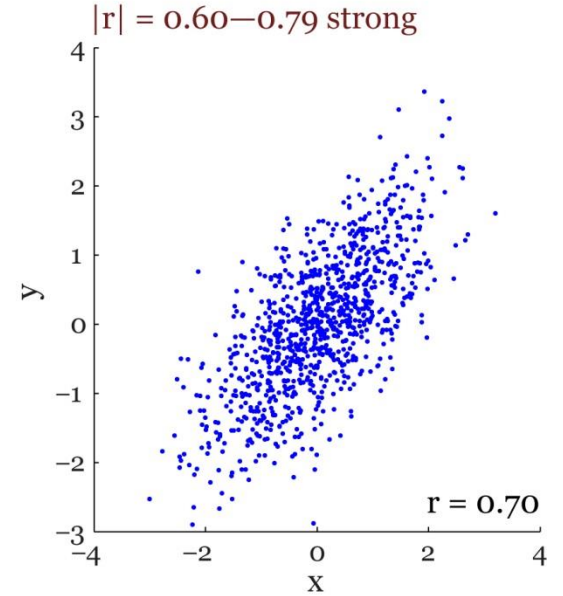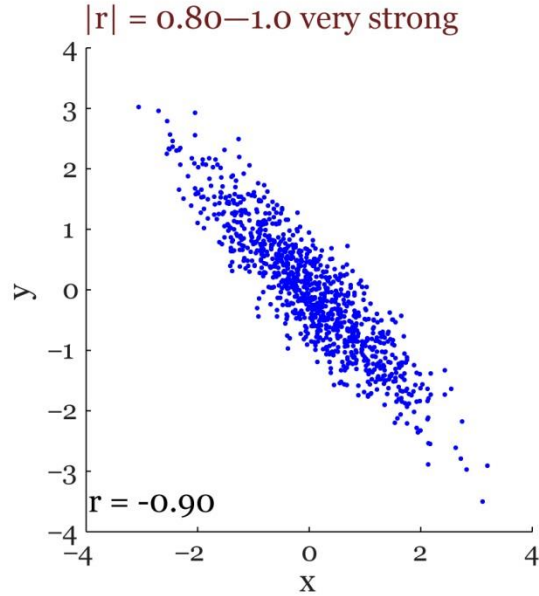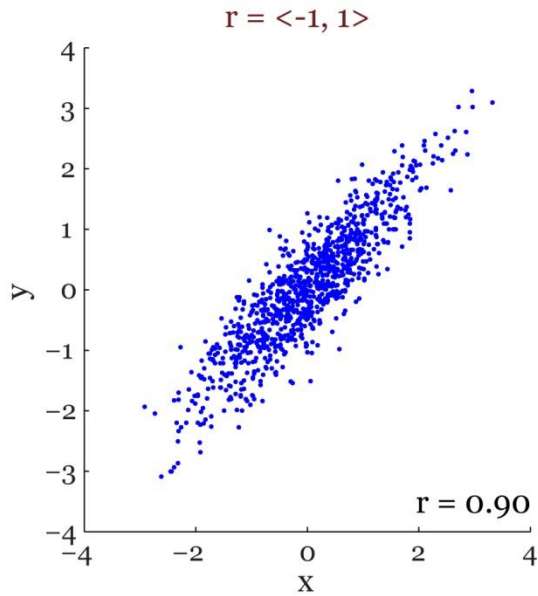Answer: Correlation analysis.

### Connection
- Correlation indicates a relationship not causality.
- We need to find a connection to say that relationship is causal (i.e. examine that hormonal response to stress can elevate blood pressure).
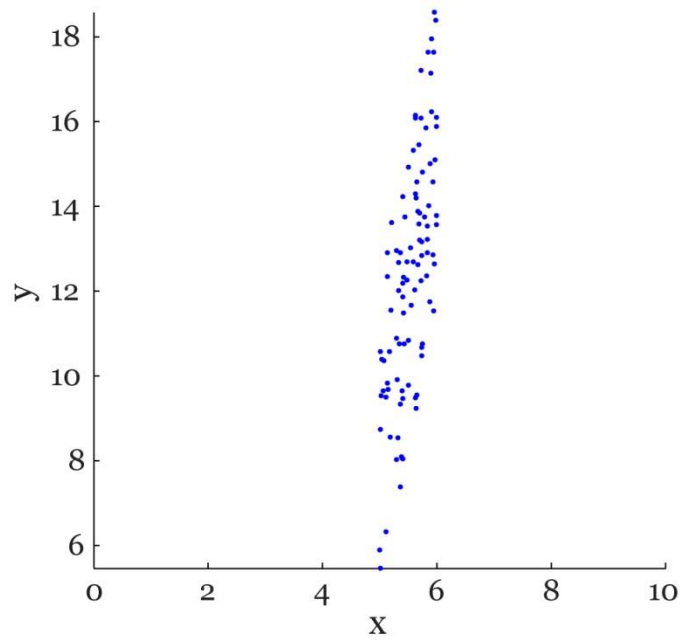
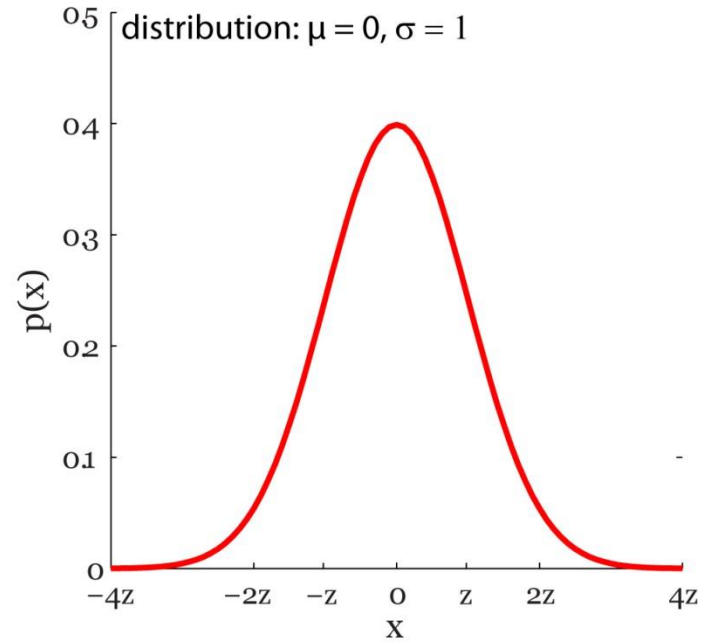**Do not confuse simultaneity with causality**

Global average temperature (°C)

Number of pirates (thousands)

R = -0.977, p < 0.001

AD 2000
AD 1980
AD 1940
AD 1920
AD 1880
AD 1820
AD 1860

# Correlation values



| r = <-1, 1> | |r| = 0.80—1.0 very strong | |r| = 0.60—0.79 strong |
|---|---|---|
| r = 0.90 | r = -0.90 | r = 0.70 |

| |r| = 0.40—0.59 moderate | |r| = 0.20—0.39 weak | |r| = 0.00—0.19 very weak |
|---|---|---|
| r = 0.50 | r = 0.30 | r = 0.10 |

# Data (raw)



# Data (z-scored)

distribution: $\mu = 0$, $\sigma = 1$

# Data (raw)



# Data (z-scored)

r = 0.69

$$r = \frac{\sum_{i=1}^{n}\left(\dfrac{x_i - \mu_x}{\sigma_x}\right)\left(\dfrac{y_i - \mu_y}{\sigma_y}\right)}{n}$$

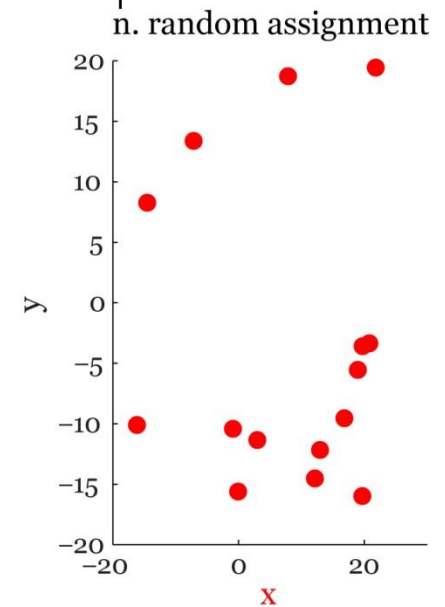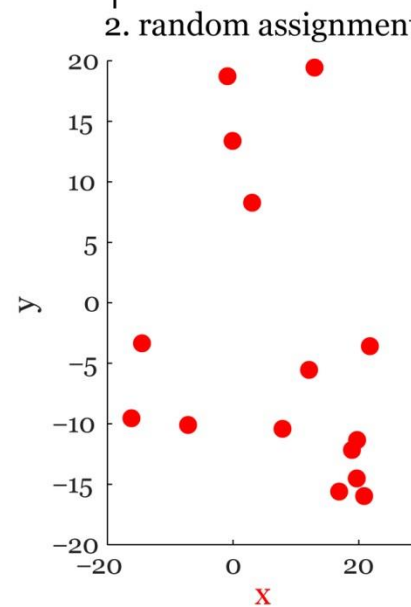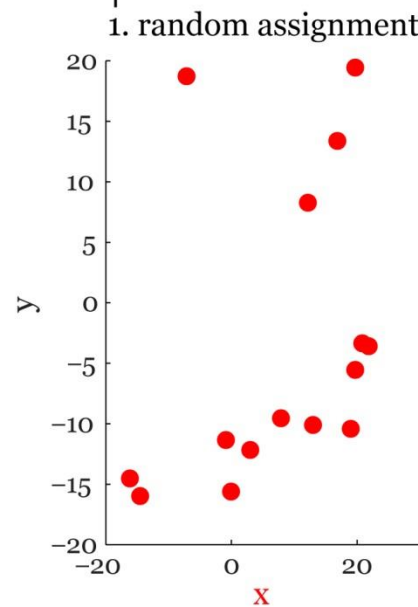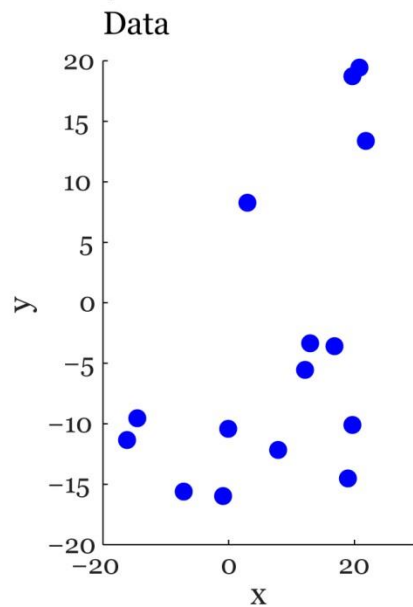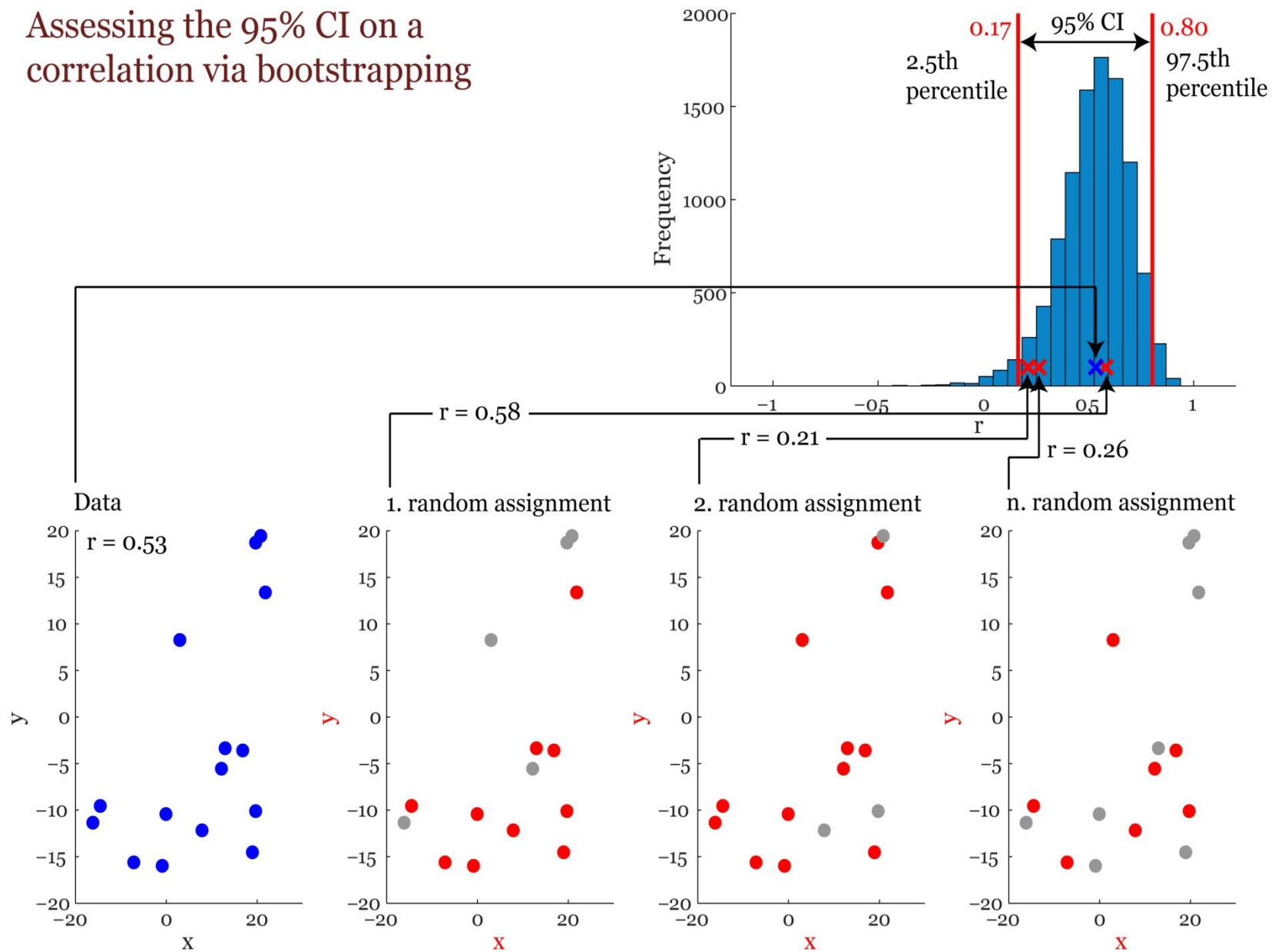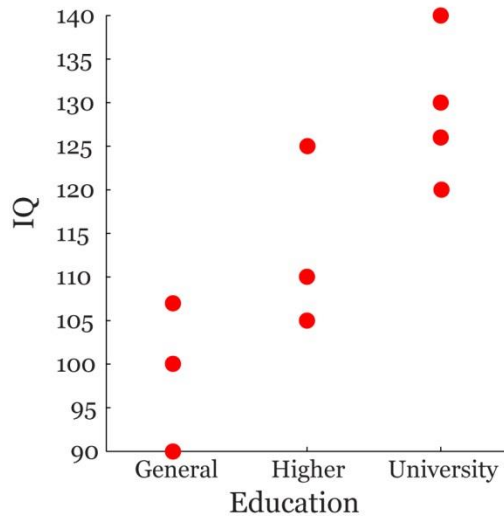# Assessing the statistical significance of correlation via randomization



r = 0.53, p = 0.04

r = 0.34

r = -0.28

r = -0.07

Data

1. random assignment

2. random assignment

n. random assignment

# Assessing the 95% CI on a correlation via bootstrapping

# Spearman rank correlation



| Raw data | | Ranks | | Difference | |
|---|---|---|---|---|---|
| IQ | Education | $x_r$ | $y_r$ | $d = x_r-y_r$ | $d^2$ |
| 100 | General | 2 | 2 | 0 | 0 |
| 105 | Higher | 3 | 5 | -2 | 4 |
| 140 | University | 10 | 8.5 | 1.5 | 2.25 |
| 125 | Higher | 7 | 5 | 2 | 4 |
| 110 | Higher | 5 | 5 | 0 | 0 |
| 130 | University | 9 | 8.5 | 0.5 | 0.25 |
| 90 | General | 1 | 2 | -1 | 1 |
| 107 | General | 4 | 2 | 2 | 4 |
| 120 | University | 6 | 8.5 | -2.5 | 6.25 |
| 126 | University | 8 | 8.5 | -0.5 | 0.25 |

$$\sum d^2 = 22$$

$$r_s = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

$$r_s = 1 - \frac{6 \times 22}{10(100 - 1)}$$

$$r_s = 0.87$$

Pearson: r = 0.88
Spearman: r = 0.96

r = 0

## Correlations

Normal speaker is able to perform sustained vowel phonation for several seconds without voice breaks that represent impaired function of vocal folds. To verify if vocal fold function disability in patients with Huntington's disease (HD) corresponds to overall motor disability, researcher collected sustained phonations from 32 speakers with HD and performed analysis of maximum phonation time until voice breaks (mptvb). He also assessed every patient using clinical motor scale of Unified Huntington's Disease Raring Scale (UHDRS).
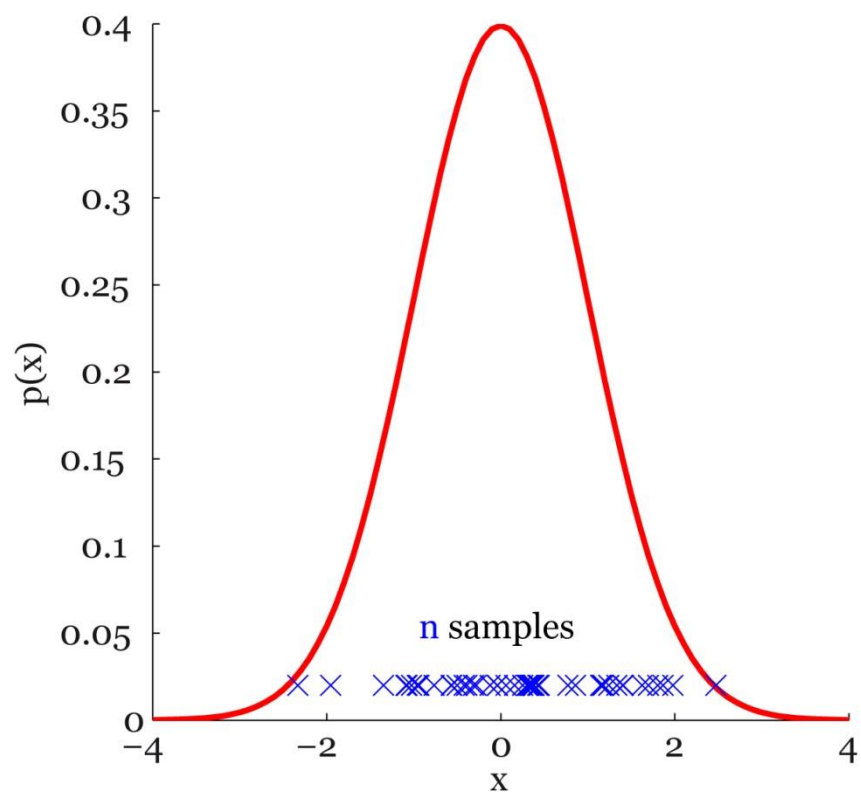


Pearson:
$r = -0.47, p = 0.01$

Spearman:
$r = -0.61, p < 0.001$

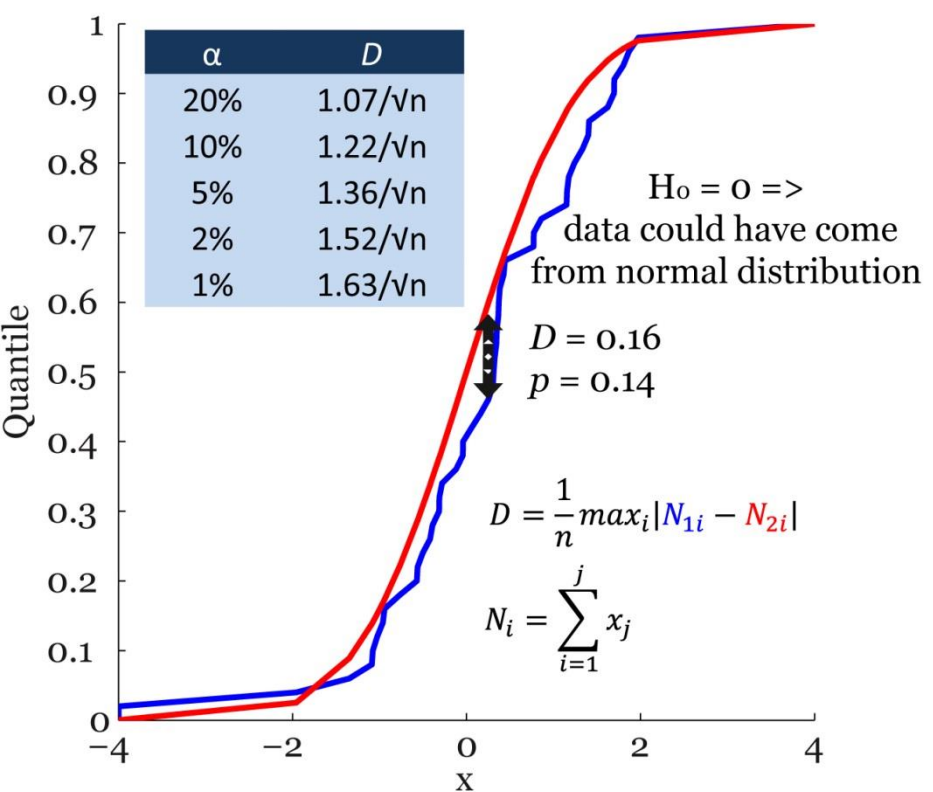Spearman correlation is more powerful due to non-normally distributed data

# Testing normality of the data using goodness-of-fit test: Kolmogorov-Smirnov

## Probability density function (PDF)
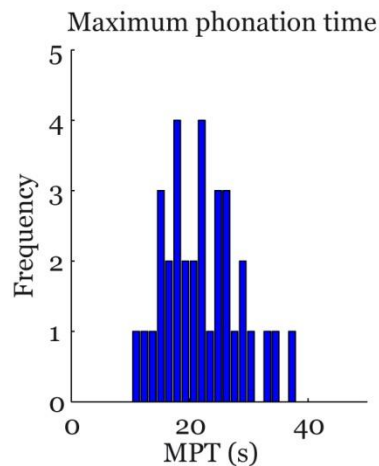
## Cumulative distribution function (CDF)



| α | D |
|-----|--------|
| 20% | $1.07/\sqrt{n}$ |
| 10% | $1.22/\sqrt{n}$ |
| 5% | $1.36/\sqrt{n}$ |
| 2% | $1.52/\sqrt{n}$ |
| 1% | $1.63/\sqrt{n}$ |

$H_0 = 0 =>$
data could have come
from normal distribution

$D = 0.16$
$p = 0.14$

$$D = \frac{1}{n} max_i |N_{1i} - N_{2i}|$$
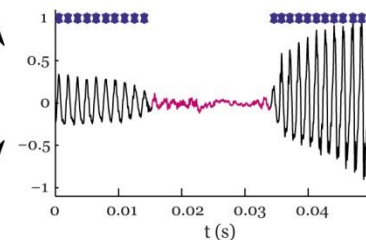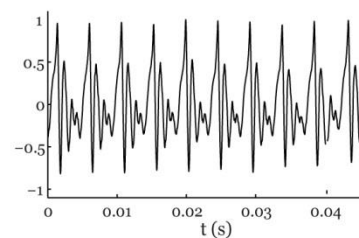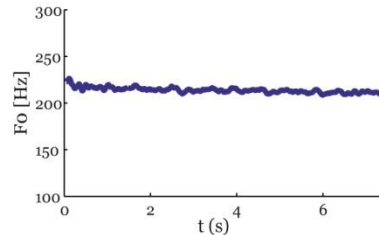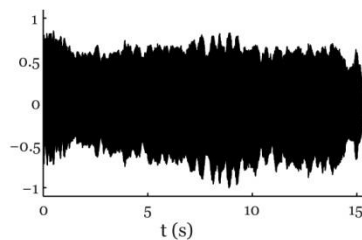
$$N_i = \sum_{i=1}^{j} x_j$$

n samples

CDF is integral of PDF evaluated from $-\infty$ to x

# Testing the data normality using Kolmogorov-Smirnov test
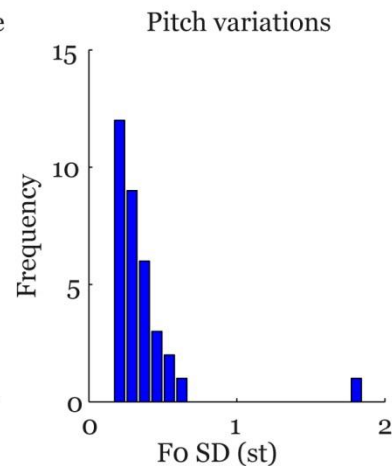
Researcher wants to verify the normality of data based upon 4 measures extracted from sustained phonation of 34 healthy speakers.



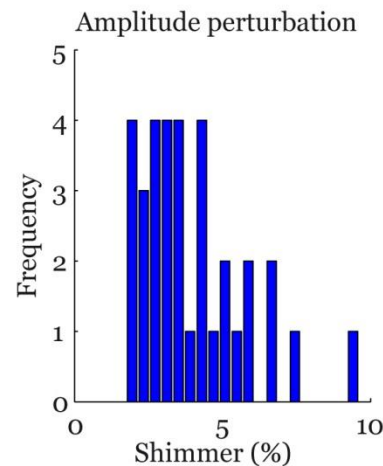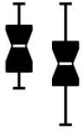|  |  |  |  |
|---|---|---|---|
| Maximum phonation time | Pitch variations | Amplitude perturbation | Degree of unvoiced segments |
| Normal distribution | Normal distribution with outliers | Log-normal distribution | Non-normal distribution |
| $D = 0.11, p = 0.76$ Do not reject $H_0$! | $D = 0.25, p = 0.02$ Reject $H_0$! | $D = 0.14, p = 0.52$ Do not reject $H_0$! | $D = 0.39, p < 0.001$ Reject $H_0$! |

# 1) Parametric vs. non-parametric data

| | **Parametric** | | **Non-parametric** | |
|---|---|---|---|---|
| **Expected distribution:** | Normal | | No limits | |
| **Expected variance:** | Homogeneous | | No limits | |
| **Typical scale:** | Interval | $y = x + d$ <br> dates (years), temperature (°C), IQ scale | Nominal | Non-comparable <br> YES/NO, colours, gender, phone numbers |
| | Ratio | $y = k \cdot x + d$ <br> velocities, lengths, age... | Ordinal | Can be sorted in terms of „*greater*" or „*less*" <br> education |
| **Central measures:** | Mean (SD) | | Median (IQR) | |
| **Advantages:** | More powerful | | Less vulnerable to outliers | |
| **Disadvantages:** | More samples needed | | Less powerful | |

# 2) Tests for normality

**Visual inspection**
- Relatively robust

**Chi-Square test** (chi2gof)
- Ties may be problematic
- Unsuitable for small samples

**Kolmogorov-Smirnov test** (kstest)
- Work for small samples (optimal > 50)
- Ties are not problem
- Lower power

**Shapiro-Wilk test** (not available in Matlab)
- Highest power among all tests of normality
- Most commonly used
- Samples < 50

**Liliefors test** (lillietest)
- Higher power than KS test
- Not effective for small samples (n >80)

**Anderson-Darling test** (adtest)
- Higher power than KS test

**Cramér-von-Mises test** (not available in Matlab)
- Higher power than KS test

# 3) Parametric vs. non-parametric tests
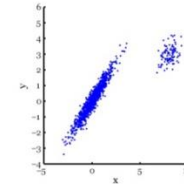
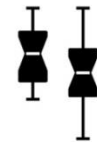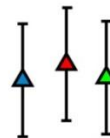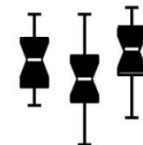|  | Parametric | Non-parametric |
|---|---|---|
| Correlations: | Pearson *corr* | Spearman *corr* |
| Independent 2 groups: | Independent t-test *ttest2* | Mann-Whitney U test *ranksum* |
| Paired 2 „levels": | Paired t-test *ttest* | Wilcoxon signed-rank test *signrank* |
| Independent > 2 groups: | Analysis of variance *anova1* | Kruskal-Wallis test *kruskalwallis* |
| Paired > 2 „levels": | Repeated measures analysis of variance *ranova* | Friedman test *friedman* |