

VYBRANÉ APLIKACE MATEMATICKÉ STATISTIKY.  
SYNOPSIS OF COURSE AT OU, WINTER SEMESTERS 2015/2016,  
2016/2017, SUMMER SEMESTER 2018/2019

PASHA ZUSMANOVICH

Each class lasted 1.5 hours. *Stuff typed in italic was not covered during summer semester 2018/2019.*

CLASS 1. (FEBRUARY 14, 2019)

Subject of statistics.

*Interesting applications of statistics: detection election frauds ([KSP], [KYHT]), patterns in citations [SR].*

Types of data: numerical and categorical. Statistical data always contains errors and is incomplete.

Bar charts and histograms.

Average, standard deviation, their meaning. Assuming  $\bar{x} = (x_1, \dots, x_n)$ ,

$$m(\bar{x}) = \frac{x_1 + \dots + x_n}{n}$$
$$\sigma(\bar{x}) = \sqrt{\frac{(x_1 - m(\bar{x}))^2 + \dots + (x_n - m(\bar{x}))^2}{n}}.$$

Median, quantiles (generalization of median).

A glimpse into R: installation, usage as calculator, assignments, 1-dimensional arrays, functions. `help()`, `example()`, `mean()`, `sd()`, `median()`, `quantile()`, `plot()`, `barplot()`. Drawing histograms in different ways.

A toy example: plot *and linear regression* of air pollution against temperature for a 24 hour period in Ostrava.

CLASS 2. (FEBRUARY 21, 2019)

Mode.

Discrete vs. continuous distributions.

Density function of the normal distribution:

$$f_{m,\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-m}{\sigma}\right)^2}.$$

Its importance, its properties: symmetry, maximum (by calculating derivative). Central Limit Theorem.

Density function of the logistic distribution:

$$f_{m,s}(x) = \frac{e^{-\frac{x-m}{s}}}{s(1 + e^{-\frac{x-m}{s}})^2}.$$

Its standard deviation:  $\sigma = \frac{\pi s}{\sqrt{3}}$ .

Density function of the Laplace distribution:

$$f_{m,b}(x) = \frac{1}{2b} e^{-\frac{|x-m|}{b}}.$$

Its standard deviation:  $\sigma = b\sqrt{2}$ .

---

*Date:* last modified May 9, 2019.

These distributions have, roughly, the same properties as a normal distribution: the same “bell-shaped” form, attain one maximum “in the middle” (average), and are symmetric, but logistic distribution is “heavier on tails”, and Laplace distribution has a sharp peak in the middle.

Fitting in R real data to normal, Laplace and logistic distributions.

### CLASS 3. (FEBRUARY 28, 2019)

Demonstration in R of the central limit theorem. (An alternative demonstration, using a different R code, can be found in [Cr1, §7.3.2]).

Skewness and kurtosis, their meaning (according to [Cr2, pp. 84–87]) and Wikipedia ([Sk<sub>w</sub>], [K<sub>w</sub>]).

$$skewness(\bar{x}) = \frac{\text{3rd moment}(\bar{x})}{\sigma(\bar{x})^3} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - m(\bar{x}))^3}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - m(\bar{x}))^2\right)^{\frac{3}{2}}};$$

$$kurtosis(\bar{x}) = \frac{\text{4th moment}(\bar{x})}{\sigma(\bar{x})^4} - 3 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - m(\bar{x}))^4}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - m(\bar{x}))^2\right)^2} - 3.$$

Kurtosis of a normal distribution is equal to 0.

“Paradoxes” in statistics (according to [T, §6.5]).

Weighted mean. `weighted.mean()` in R. The usual mean of a discrete statistical distribution  $(x_1, \dots, x_n)$  can be interpreted as a weighted mean, if we assume that all  $x_i$ ’s are pairwise distinct, and each appears with a frequency (probability)  $p_i$ . Then

$$m(\bar{x}) = p_1x_1 + \dots + p_nx_n = \frac{p_1x_1 + \dots + p_nx_n}{p_1 + \dots + p_n}$$

(as  $p_1 + \dots + p_n = 1$ ).

Computation of weighted population density: if the whole area is divided to  $n$  regions with population  $x_1, \dots, x_n$  and areas  $s_1, \dots, s_n$ , then the weighted population density is the weighted mean of densities per region, with weights given by population:

$$\frac{x_1 \frac{x_1}{s_1} + \dots + x_n \frac{x_n}{s_n}}{x_1 + \dots + x_n}.$$

Simpson’s paradox (according to [Si<sub>w</sub>]). UC Berkeley suitcase.

|              | men   | men admitted    | women | women admitted |
|--------------|-------|-----------------|-------|----------------|
| Department 1 | $m_1$ | $\lambda_1 m_1$ | $w_1$ | $\mu_1 w_1$    |
| Department 2 | $m_2$ | $\lambda_2 m_2$ | $w_2$ | $\mu_2 w_2$    |

It could be that  $\lambda_1 < \mu_1$  and  $\lambda_2 < \mu_2$ , but

$$\frac{\lambda_1 m_1 + \lambda_2 m_2}{m_1 + m_2} > \frac{\mu_1 w_1 + \mu_2 w_2}{w_1 + w_2}.$$

Another example of Simpson’s paradox often occurs in US election system, see, e.g. [W].

### CLASS 4. (MARCH 7, 2019)

Discussion of homeworks 1-2.

Confidence intervals for normal distribution. Standard error. (According to [D, pp. 63-64] and [C<sub>w</sub>]).

Using confidence intervals to determine sample size.

$$\text{Margin error} = \frac{\sigma}{\sqrt{n}} N_{1-\frac{1-\alpha}{2}},$$

where  $\alpha$  is confidence interval (say,  $\alpha = 0.95$ ), and  $N$  are quantiles for the normal distribution.  
 Q-Q plot of one data against another.  
 Tests for normality: normal scores, Q-Q plots (according to [CC, pp. 220–223]).

### CLASS 5. (MARCH 21, 2019)

Discussion of homework 3.

Correlation:

$$\begin{aligned} \text{cov}(\bar{x}, \bar{y}) &= \sum_{i=1}^n (x_i - m(\bar{x})) (y_i - m(\bar{y})) \\ \text{cor}(\bar{x}, \bar{y}) &= \frac{\text{cov}(\bar{x}, \bar{y})}{\sqrt{\text{cov}(\bar{x}, \bar{x}) \text{cov}(\bar{y}, \bar{y})}}. \end{aligned}$$

Properties of correlation:

$$\begin{aligned} \text{cor}(\bar{x}, \bar{x}) &= 1 \\ \text{cor}(\bar{x}, \bar{y}) &= \text{cor}(\bar{y}, \bar{x}) \\ -1 &\leq \text{cor}(\bar{x}, \bar{y}) \leq 1. \end{aligned}$$

The latter one follows from the Cauchy–Schwarz inequality:

$$\left( \sum_{i=1}^n a_i b_i \right)^2 \leq \left( \sum_{i=1}^n a_i^2 \right) \left( \sum_{i=1}^n b_i^2 \right).$$

Correlation between linearly dependent datasets is equal to 1 or  $-1$ .

Example of two datasets with correlation zero: let  $\bar{x}$  be any vector of even length whose alternating sum is zero, for example,  $x_1 = x_2, x_3 = x_4, \dots, x_{2n-1} = x_{2n}$ , and  $\bar{y}$  is oscillating, say,  $y_i = 1$  for  $i$  odd and  $y_i = 0$  for  $i$  even. Then  $m(\bar{y}) = \frac{1}{2}$ , and

$$\begin{aligned} \text{cov}(\bar{x}, \bar{y}) &= \sum_{i=1,3,\dots,2n-1} (x_i - m(\bar{x})) \left(1 - \frac{1}{2}\right) + \sum_{i=2,4,\dots,2n} (x_i - m(\bar{x})) \left(0 - \frac{1}{2}\right) \\ &= \frac{1}{2} (x_1 - x_2 + x_3 - x_4 + \dots + x_{2n-1} - x_{2n}) = 0, \end{aligned}$$

and hence  $\text{cor}(\bar{x}, \bar{y}) = 0$ .

Use and misuse of correlation. “Correlation is not causation”.

Correlation matrix, its properties (symmetric, 1’s on the main diagonal).

### CLASS 6.

Discussion of homework 6 (R code demonstrating Central Limit Theorem for any distribution).

Iterative correlation matrices (according to [Ch]).

The case of  $2 \times 2$  matrices:

$$\begin{pmatrix} 1 & a \\ a & 1 \end{pmatrix} \mapsto \begin{pmatrix} 1 & t(a) \\ t(a) & 1 \end{pmatrix}$$

where

$$t(a) = \text{cor}((1, a), (a, 1)) = \frac{(1 - \frac{1+a}{2})(a - \frac{1+a}{2}) + (a - \frac{1+a}{2})(1 - \frac{1+a}{2})}{\sqrt{((1 - \frac{1+a}{2})^2 + (a - \frac{1+a}{2})^2)((a - \frac{1+a}{2})^2 + (1 - \frac{1+a}{2})^2)}} = -1$$

unless  $a \neq 1$ .

## CLASS 7.

Descriptive and inferential statistics.

Statistical models. Linear regression. Explanatory and response variables.

$$Y = \alpha + \beta X + \varepsilon.$$

$Y$  - response variable,  $X$  - predictor,  $\varepsilon$  - error term.

$$\varepsilon \sim N(0, \sigma^2).$$

Simple (one predictor) and multiple (several predictors) regressions (according to [AR, Chapter 7]).

Least squares. Derivation of coefficients for simple linear regression via least squares:

$$\beta = \text{cor}(\bar{x}, \bar{y}) \frac{\sigma(\bar{y})}{\sigma(\bar{x})}$$
$$\alpha = m(\bar{y}) - \beta m(\bar{x}).$$

Residuals. Standard deviation of residuals and test of residuals for normality as criteria for “goodness” of a linear model.

## CLASS 8.

Generalized additive models (according to [Cr2, pp. 146–148]).

## CLASS 9.

Clustering: Hierarchical, K-means, gravitational algorithms.

Examples: genetic analysis; transportation, traffic.

## CLASS 10.

Null and alternative hypotheses. Hypotheses testing,  $p$ -values (according to Pruim, pp. 71 onwards).

## REFERENCES

- [AR] J. Albert and M. Rizzo, *R by Example*, Springer, 2012.
- [Ch] C.-H. Chen, *Generalized association plots: information visualization via iterative generated correlation matrices*, *Statistica Sinica* 12 (2002), 7–29.
- [CC] Y. Cohen and J.Y. Cohen, *Statistics and Data with R*, Wiley, 2008.
- [Cr1] M.J. Crawley, *The R Book*, 2nd ed., Wiley, 2013.
- [Cr2] ———, *Statistics: An Introduction Using R*, 2nd ed., Wiley, 2014.
- [D] P. Dalgaard, *Introductory Statistics with R*, 2nd ed., Springer, 2008.
- [KYHT] P. Klimek, Y. Yegorov, R. Hanel, and S. Thurnera, *Statistical detection of systematic election irregularities*, *Proc. Nat. Acad. Sci. USA* **109** (2012), 16469–16473.
- [KSP] D. Kobak, S. Shpilkin, and M.S. Pshenichnikov, *Statistical anomalies in 2011-2012 Russian elections revealed by 2D correlation analysis*, arXiv:1205.0741v2.
- [SR] M.V. Simkin and V.P. Roychowdhury, *An introduction to the theory of citing*, *Significance* **3** (2006), *N*4, 179–181; arXiv:math/0701086.
- [T] T. Tao, *Compactness and Contradiction*, AMS, 2013.
- [W] S. Wang, *The great gerrymander of 2012*, *The New York Times*, February 2, 2013.

## WIKIPEDIA ARTICLES

[C<sub>w</sub>] *Confidence interval.*

[K<sub>w</sub>] *Kurtosis.*

[Si<sub>w</sub>] *Simpson's paradox.*

[Sk<sub>w</sub>] *Skewness.*

*Email address:* `pasha.zusmanovich@gmail.com`