# NMAI059 Probability and statistics 1
## Class 12

Robert Šámal

# Overview

# Sample mean & variance

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$$

$$\bar{S}_n^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{X}_n)^2$$

$$\widehat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X}_n)^2 \longrightarrow \text{unbiased}$$

lower MSE

# Maximal likelihood method, ML)

*(handwritten: p/f)*

*(handwritten: unknown parameter)* *(handwritten: known data)*

- **Maximal likelihood method:**
  choose $\vartheta$ that maximizes $L(x; \vartheta)$.  *(handwritten: $\to \hat{\vartheta}$ — usually done by differ.)*

- for convenience we put $\ell(x; \vartheta) = \log(L(x; \vartheta))$

- by independence of $X_1$, $X_2$, etc. we have

$$L(x; \vartheta) = L(x_1; \vartheta) \dots L(x_n; \vartheta)$$
$$\ell(x; \vartheta) = \ell(x_1; \vartheta) + \dots + \ell(x_n; \vartheta)$$

*(handwritten: do not know)*

*(handwritten: p = 9/20)*

| Bin(20,p) | 0.2 | 0.3 | 0.4 | 0.45 | 0.5 | 0.55 | 0.6 |
|---|---|---|---|---|---|---|---|
| 7 | 0.0545 | 0.1643 | 0.1659 | 0.1221 | 0.0739 | 0.0366 | 0.0146 |
| 8 | 0.0222 | 0.1144 | 0.1797 | 0.1623 | 0.1201 | 0.0727 | 0.0355 |
| 9 | 0.0074 | 0.0654 | 0.1597 | 0.1771 | 0.1602 | 0.1185 | 0.071 |
| 10 | 0.002 | 0.0308 | 0.1171 | 0.1593 | 0.1762 | 0.1593 | 0.1171 |
| 11 | 0.0005 | 0.012 | 0.071 | 0.1185 | 0.1602 | 0.1771 | 0.1597 |
| 12 | 0.0001 | 0.0039 | 0.0355 | 0.0727 | 0.1201 | 0.1623 | 0.1797 |
| 13 | 0 | 0.001 | 0.0146 | 0.0366 | 0.0739 | 0.1221 | 0.1659 |
| 14 | 0 | 0.0002 | 0.0049 | 0.015 | 0.037 | 0.0746 | 0.1244 |

*(handwritten: know → 9)*

*(handwritten: 20 measurements, indep. 0/1, p — prob of 1)*

*(handwritten: 9 said yes)*

*(handwritten: fixed max.)*

*(handwritten: depends on p)*

*(handwritten: pmf $P(X = k) = \binom{20}{k} p^k (1-p)^{20-k}$ )*

# ML – further illustration

$N(\mu, \sigma^2)$    $\vartheta = (\mu, \sigma)$

$x = (x_1, \ldots, x_n)$ numbers — realizations of $X_1, \ldots, X_n \sim N(\mu_c, \sigma^2)$

$$f_{X_i}(x_i) = \frac{1}{\sqrt{2\pi}\,\sigma}\, e^{-\left(\frac{x_i - \mu}{\sigma}\right)^2/2} \quad \ldots \text{ formula for pdf of } \uparrow$$

$$L(x_i; \vartheta) \implies \ell(x_i; \vartheta) = -\left(\frac{x_i - \mu}{\sigma}\right)^2/2 - \log \sigma - \log \sqrt{2\pi}$$

$$\ell(x; \vartheta) = -\sum_{i=1}^{n} \left(\frac{x_i - \mu}{\sigma}\right)^2/2 - n\log\sigma - n\log\sqrt{2\pi}$$

FIND $\vartheta_{ML}(\omega)$ max maximizing this

TO DO THIS DIFFERENT!

$$\frac{\partial \ell}{\partial \mu} = +\sum_{i=1}^{n} 2\left(\frac{x_i - \mu}{\sigma}\right)\frac{1}{2}\cdot\frac{-1}{\sigma} = \frac{1}{\sigma^2}\sum_{i=1}^{n}(x_i - \mu)$$

$\boxed{EQ1} = 0 \implies \mu = \frac{x_1 + \ldots + x_n}{n} = \bar{x}_n$

$\boxed{\hat{\mu} := \bar{x}_n}$

$$\frac{\partial \ell}{\partial \sigma} = +\sum_{i=1}^{n} \frac{(x_i - \mu)^2}{\sigma^3}\frac{(-1)}{2} - \frac{n}{\sigma} = 0$$

$\boxed{EQ2} \implies \sigma^2 = \frac{1}{n}\sum(x_i - \mu)^2$

$\boxed{\hat{\sigma}^2 = \frac{1}{n}\sum(x_i - \bar{x}_n)^2}$

# Overview

# Interval estimation

*unknown parameter*

- ▶ Instead of estimating by one number we compute from our data an interval $[\hat{\Theta}^-, \hat{\Theta}^+]$

## Definition

*Let $\hat{\Theta}^-$, $\hat{\Theta}^+$ be random variables that depend on the random sample $X = (X_1, \ldots, X_n)$ from distribution $F_\vartheta$. These random variables describe a $1 - \alpha$ confidence interval, if*

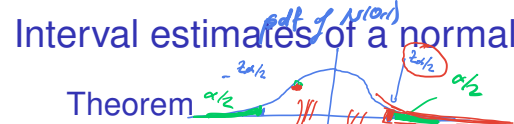$$P(\hat{\Theta}^- \leq \vartheta \leq \hat{\Theta}^+) \geq 1 - \alpha.$$

$A$      $B$

*NOT A PROB. STATEMENT ABOUT $\vartheta$ !*

- ▶ these are two-sided estimates
- ▶ one-sided: $[\hat{\Theta}^-, \infty)$ or $(-\infty, \hat{\Theta}^-]$

*$\vartheta$ IS A FIXED PARAMETER THAT WE DON'T KNOW*

# Interval estimates of a normal variable

*edf of N(0,1)*

*we measure temper.*
*constr. of thermometer → $\sigma$*
*read temp. $\vartheta$*

$-z_{\alpha/2}$   $z_{\alpha/2}$   $\alpha/2$   TRUTH
$\alpha/2$

**Theorem**
$X_1, \ldots, X_n$ *random sample from* $N(\vartheta, \sigma^2)$.
$\sigma$ **is known**, *we need to estimate* $\vartheta$, *we choose* $\alpha \in (0,1)$.
*Let* $\Phi(z_{\alpha/2}) = 1 - \alpha/2$. *We put* $\hat{\Theta}_n := \bar{X}_n$ *and*

$z_{\alpha/2} := \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)$

$\bar{X}_n \sim N\left(\vartheta, \frac{\sigma^2}{n}\right)$

$z = \frac{\bar{X}_n - \vartheta}{\sigma/\sqrt{n}} \sim N(0,1)$

$$C_n := [\hat{\Theta}_n - z_{\alpha/2}\frac{\sigma}{\sqrt{n}}, \quad \hat{\Theta}_n + z_{\alpha/2}\frac{\sigma}{\sqrt{n}}]$$

*Then* $P(C_n \ni \vartheta) = 1 - \alpha$.

$\bar{X}_u - \bigcirc$   $\bar{X}_u$   $\bar{X}_u + \bigcirc$

**Důkaz.**

$P(C_n \ni \vartheta) = P\left(\left|\bar{X}_n - \vartheta\right| \le z_{\alpha/2}\frac{\sigma}{\sqrt{n}}\right)$

$= P\left(\left|\frac{\bar{X}_n - \vartheta}{\sigma/\sqrt{n}}\right| \le z_{\alpha/2}\right) = P\left(-z_{\alpha/2} \le z \le z_{\alpha/2}\right) = \Phi(z_{\alpha/2}) - \Phi(-z_{\alpha/2})$

$z := ?$

$\cdot \left(1 - \frac{\alpha}{2}\right) - \frac{\alpha}{2} = 1 - \alpha$

# Interval estimates using CLT

## Theorem

$X_1, \ldots, X_n$ *random sample from a distribution with mean* $\vartheta$ *and variance* $\sigma^2$.

*not necessary* $N(\vartheta, \sigma^2)$

$\boxed{\sigma \textbf{ is known}}$ *we need to estimate* $\vartheta$, *we choose* $\alpha \in (0,1)$.
*Let* $\Phi(z_{\alpha/2}) = 1 - \alpha/2$. *We put* $\hat{\Theta}_n := \bar{X}_n$ *and*

$$C_n := [\hat{\Theta}_n - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \quad \hat{\Theta}_n + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$$

*not necess normal*

$$Z_n := \frac{\bar{X}_n - \vartheta}{\sigma/\sqrt{n}}$$

*Then* $\lim_{n \to \infty} P(C_n \ni \vartheta) = 1 - \alpha$.

B&T CLT:
$$Z_n \xrightarrow{d} N(0,1)$$

Proof $P(C_n \ni \vartheta) = P\left( |\bar{X}_n - \vartheta| \leq z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right)$

$= P\left( |Z_n| \leq z_{\alpha/2} \right) = F_{Z_n}(z_{\alpha/2}) - F_{Z_n}(-z_{\alpha/2}) = ?$

$\xrightarrow{n \to \infty} \phi(z_{\alpha/2}) - \phi(-z_{\alpha/2}) = 1 - \alpha$

# Student $t$-distribution

- $\bar{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i$ ... sample mean
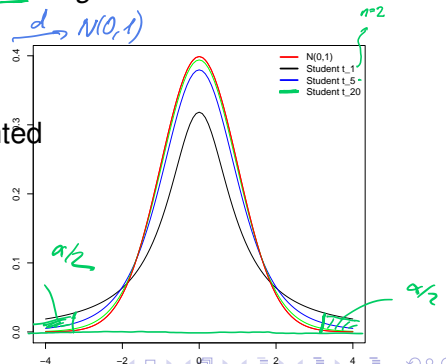- $\widehat{S}_n^2 = \frac{1}{n-1}\sum_{i=1}^{n}(X_i - \bar{X}_n)^2$ ... sample variance

*sample std. dev.*

$\widehat{S}_n = \sqrt{\widehat{S}_n^2} = \sqrt{\frac{1}{n-1}\sum (X_i - \bar{X}_n)^2}$

*does not depend on $\mu$ & $\sigma$*

- Let $X_1, \ldots, X_n \sim N(\mu, \sigma^2)$
- Then we know that $\frac{\bar{X}_n - \mu}{\sigma / \sqrt{n}} \sim N(0,1)$
- *Student $t$-distribution with $n-1$ degrees of freedom* is the distribution of r.v. $\dfrac{X_n - \mu}{\widehat{S}_n / \sqrt{n}}$

  $\xrightarrow{d} N(0,1)$

- Its cdf will be denoted $\Psi_{n-1}$
  It is tabulated, and implemented
  by computer sofware,
  in R: **pt**(x,n−1)



$n=2$

| | |
|---|---|
| — | N(0,1) |
| — | Student t_1 |
| — | Student t_5 |
| — | Student t_20 |

$\alpha/2$          $\alpha/2$

# Int. estimates of normal variable using Student $t$

**Theorem**

$X_1, \ldots, X_n$ *random sample from* $N(\vartheta, \sigma^2)$.

$\sigma$ **is not known**, *we need to estimate* $\vartheta$, *we choose* $\alpha \in (0,1)$.
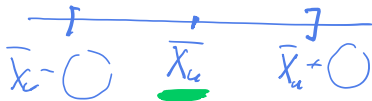
*Let* $\Psi_{n-1}(z_{\alpha/2}) = 1 - \alpha/2$. *We put* $\hat{\Theta}_n = \bar{X}_n$,

$\hat{S}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$ *and*

$$C_n := [\hat{\Theta}_n - z_{\alpha/2} \frac{\hat{S}_n}{\sqrt{n}}, \quad \hat{\Theta}_n + z_{\alpha/2} \frac{\hat{S}_n}{\sqrt{n}}]$$

*Then* $P(C_n \ni \vartheta) = 1 - \alpha$.

$\hat{S}_n$ instead of $\sigma$

$z_{\alpha/2}$ determined

using $\Psi_{n-1}$

inst. of $\Phi$

$\bar{X}_n - \bigcirc \qquad \bar{X}_n \qquad \bar{X}_n + \bigcirc$

$Z \sim \Psi_{n-1}$

$\underline{\text{Proof}}$

$P(|\bar{X}_n - \vartheta| \leq z_{\alpha/2} \frac{\hat{S}_n}{\sqrt{n}})$

$P(|\frac{\bar{X}_n - \vartheta}{\hat{S}_n/\sqrt{n}}| \leq z_{\alpha/2}) = \Psi_{n-1}(z_{\alpha/2}) - \Psi_{n-1}(-z_{\alpha/2}) = \ldots = 1 - \alpha$

$\frac{\bar{X}_n - \vartheta}{\hat{S}_n/\sqrt{n}} =: ?$

# Overview

# Intro to Hypothesis testing

- ▶ Is our coin fair?   $H_0$ : yes
- ▶ Is our die fair?   $H_0$ : yes
- ▶ Is the modified code faster then original?   $H_0$ : no
- ▶ Is the medical treatment X good? (Better than placebo, better than Y, ... )   $H_0$ : no
- ▶ Are left-handed people better at boxing?   $H_0$ : no


- ▶ two hypothesis: $H_0$, $H_1$
- ▶ $H_0$ – *null hypothesis* – default, conservative model, "unsurprising"
- ▶ $H_1$ – *alternative hypothesis* – alternative model "remarkable fact", if true
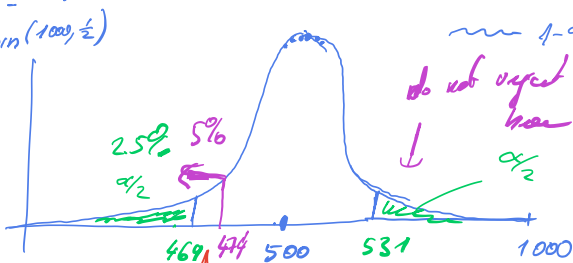
# Hypothesis testing – illustration

$n \approx 1000$

- We want to test, if a coin is fair.
- We toss it $n$-times, we get head $S$-times.
- If $\underline{|S - n/2|}$ is too large, we declare the coin not to be fair.

$P(\underline{S} = 0) = 2^{-n}$  very small

pmf $Bin(1000, \frac{1}{2})$

$\leadsto 1-\alpha$  interval est.

do not reject

$h_0$

$\alpha = 5\%$

$25\%$   $5\%$

$\frac{\alpha}{2}$

$\frac{\alpha}{2}$

469  474   500      531      1000

472

Remark Suppose, we prefer heads.

say  unfair

if  $S < \frac{n}{2} - c$

# Hypothesis testing – illustration

▶ We want to test, if a coin is fair.

▶ $H_0$: it is fair

▶ $H_1$: not fair ("Scientists discovered, that casino XY uses loaded coin.")

▶ Results: Reject $H_0$/don't reject $H_0$

▶ Type I error: false rejection. We reject $H_0$, even if it is true. Embarassing.

▶ Type II error: false non-rejection. We don't reject $H_0$, even if it is false. Unused opportunity.

▶ Need to find $k$ such that we will reject $H_0$ if $|S - n/2| > k$.

# Hypothesis testing – general approach

- ▶ We choose an appropriate statistical model.
- ▶ We choose *significance level* $\alpha$: prob. of false rejection of $H_0$. Typically $\alpha = 0.05$ (medicine/psychology – much less in high-energy physics).
- ▶ We determine *test statistics* $S = h(X_1, \ldots, X_n)$, that we will determine from the measured data.
- ▶ We determine *rejection region* – set $W$.
- ▶ We measure $x_1, \ldots, x_n$ – so-called realizations of $X_1, \ldots, X_n$.
- ▶ Decision rule: we reject $H_0$ iff $h(x_1, \ldots, x_n) \in W$.
- ▶ $\alpha = P(h(X) \in W; H_0)$
- ▶ $\beta = P(h(X) \notin W; H_1) \ldots$ *strength of the test*

- ▶ often we do not choose $\alpha$ in advance but compute so-called *p-value:* minimal $\alpha$, for which we would reject $H_0$.

# Hypothesis testing – an example

- $X_1, \ldots, X_n$ random sample from $N(\vartheta, \sigma^2)$
- $\sigma^2$ known
- $H_0 : \vartheta = 0 \qquad H_1 : \vartheta \neq 0$

# Hypothesis testing – an example

- $X_1, \ldots, X_{n_1}$ random sample from $Ber(\vartheta_X)$
- $Y_1, \ldots, Y_{n_2}$ random sample from $Ber(\vartheta_Y)$
- $H_0 : \vartheta_X = \vartheta_Y \qquad H_1 : \vartheta_X \neq \vartheta_Y$

# $p$-hacking

- ▶ we first gain data, then look for interesting stuff
- ▶ – given enough data, there will be random coincidences
- ▶ even worse, we may test, until we get the desired outcome
- ▶ *reproducibility* – after exploratory analysis of the data we make an independent measurement and a confirmatory analysis.
- ▶ or we split the data in advance to a part for hypothesis formation and part for verification . . . simple example of cross validation



**Letters in winning word of Scripps National Spelling Bee**
correlates with
**Number of people killed by venomous spiders**