

Přednáška I

AKM I

Lukáš Frýd



EVROPSKÁ UNIE
Evropské strukturální a investiční fondy
Operační program Výzkum, vývoj a vzdělávání

MŠMT
MINISTERSTVO ŠKOLSTVÍ,
MLÁDEŽE A TĚLOVÝCHOVY

Máme 3 základní pilíře

Odhad regresní rovnice
pomocí vhodné funkce

Ověření předpokladů

Statistická verifikace

Jaká je střední hodnota $E(X) = \mu$ výšky mužů v ČR?

Potřebujeme odhadnout střední hodnotu μ

Jak odhadneme μ ?

Aritmetický průměr

Medián

Jiný průměr atd.

Ovlivňuje důchod spotřebu?

$$C = \beta_0 + \beta_1 Y + \epsilon$$

Potřebujeme odhadnout podmíněnou střední hodnotu

$$E(C|Y) = \beta_0 + \beta_1 Y$$

Jak se v průměru mění C, když se mění Y

Jak odhadneme $E(C|Y)$?

Metoda nejmenších čtverců

Zobecněná metoda nejmenších čtverců

Dvoustupňová metoda nejmenších čtverců

Atd.

Který odhad použijeme a proč?

Požadujeme, aby měl nějaké vlastnosti?

Požadavky na odhad

\tilde{X} – odhad $E(X)$

b – odhad β

- **Nezkreslený (nestranný, nevychýlený) -**

$$E(\tilde{X}) = \mu$$

$$E(b) = \beta$$

- **Konzistentní**

$$\begin{array}{l} E(\tilde{X}) \xrightarrow{n \rightarrow \infty} \mu \\ \text{Var}(\tilde{X}) \xrightarrow{n \rightarrow \infty} 0 \end{array} \quad \text{pak } \tilde{X} \xrightarrow{P} \mu$$

$$\begin{array}{l} E(b) \xrightarrow{n \rightarrow \infty} \beta \\ \text{Var}(b) \xrightarrow{n \rightarrow \infty} 0 \end{array} \quad \text{pak } b \xrightarrow{P} \beta$$

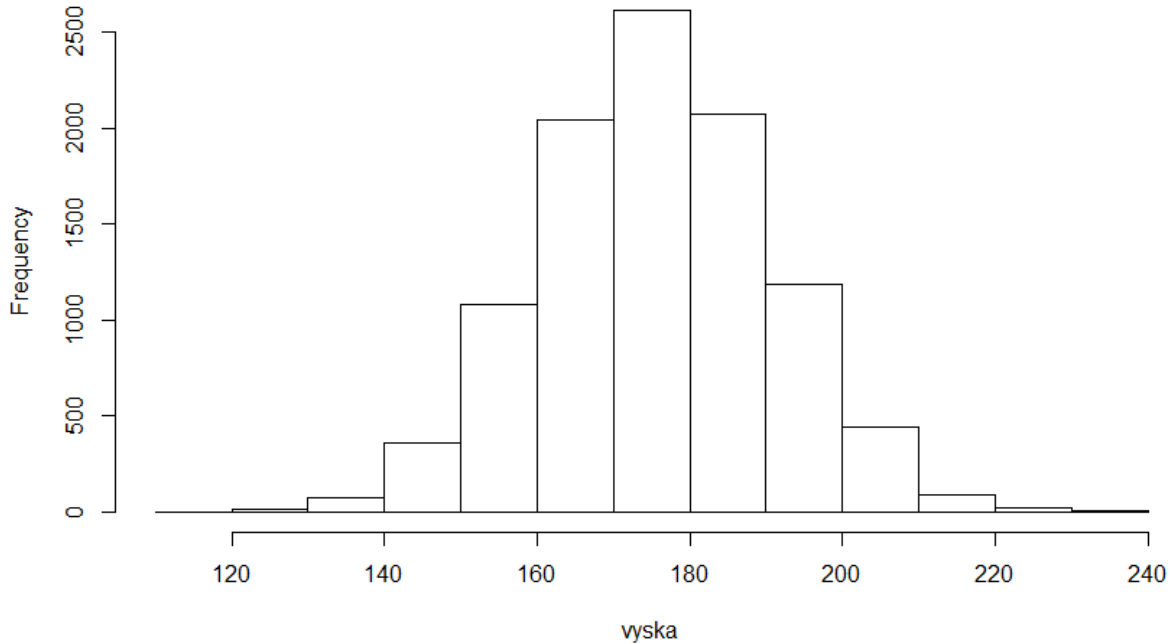
- **Vydatný**

$$\text{Var}(\tilde{X}) < \text{Var}(\bar{X})$$

$$\text{Var}(b) < \text{Var}(\bar{b})$$

Jak zjistím, že daná metoda odhadu nám bude poskytovat odhady s danými vlastnostmi?

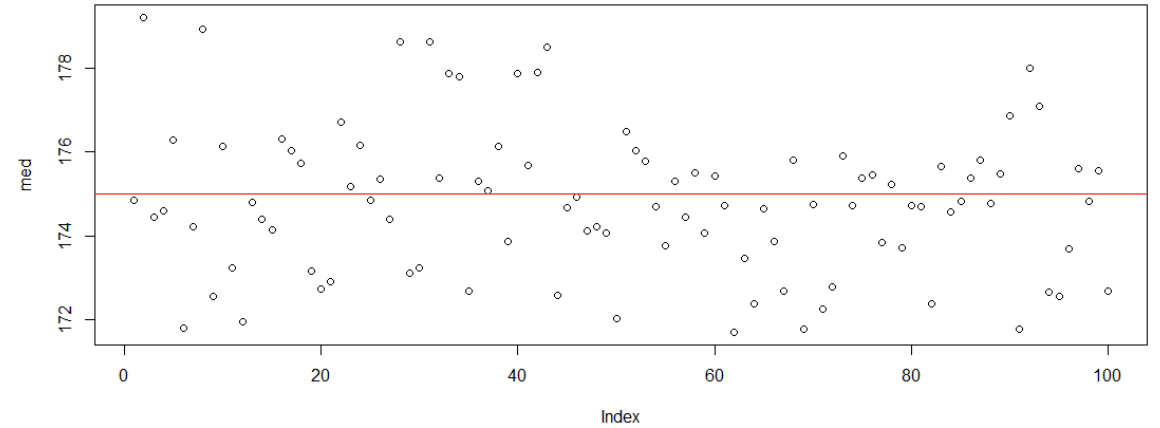
Histogram of vyska



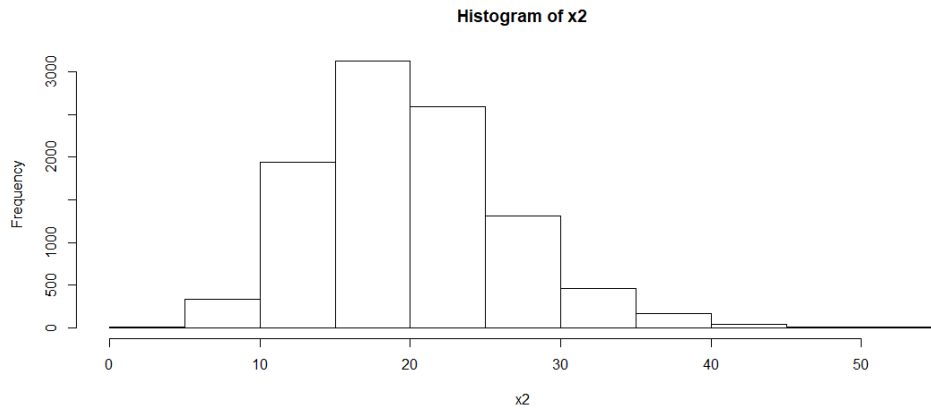
$$vyska \sim N(175, 225)$$

NEZNÁM – VYGENEROVÁVAL SEM DATA!!!!!!

V realitě existuje „nějaký neznámý“ proces, který určuje výšku



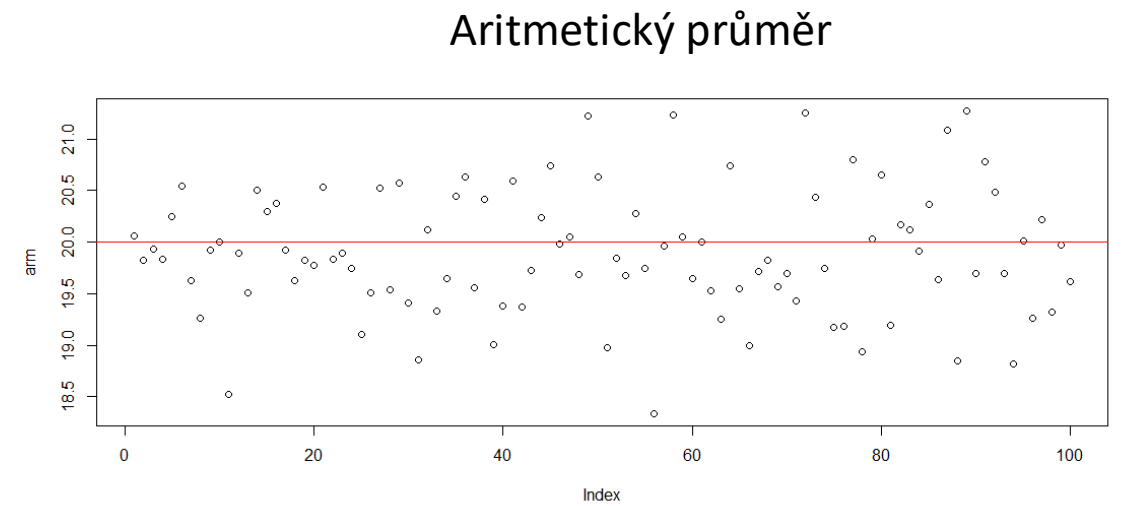
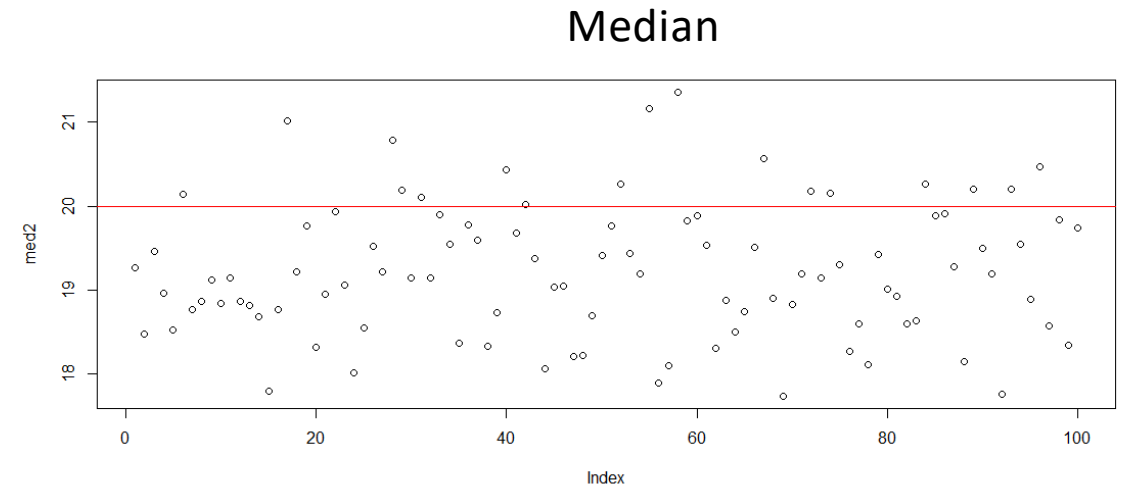
Pro odhad μ jsem použil medián



$$x \sim \chi^2(20)$$

Jaký předpoklad musí být splněn, aby nám
Medián dával nezkreslený odhad střední hodnoty?

Náhodná veličina musí mít symetrické rozdělení!



Když se přeneseme zpět k naší regresi a odhadu podmíněné střední hodnoty.
Jaké podmínky musí být splněny, aby metoda nejmenších čtverců dávala odhady:

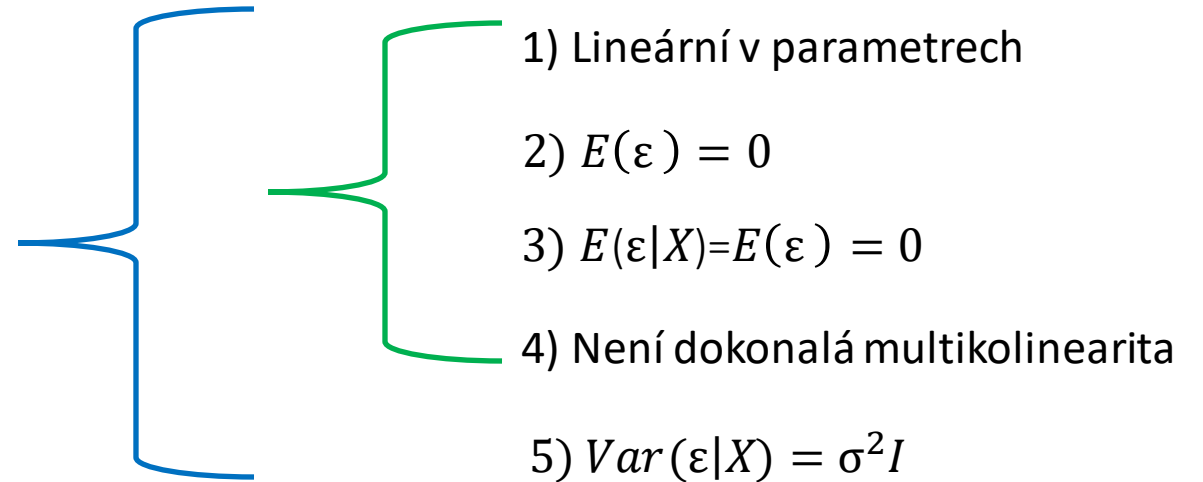
- **Nezkreslený (neustranný, nevychýlený)**
- **Konzistentní**
- **Vydatný**

Gauss-Markovovy předpoklady

Blue odhad

Nezkreslený odhad

Existuje více odhadů, které splňují 1,2,3,4
Hledáme ten nejlepší – s nejmenším rozptylem
Best linear unbiased estimators



```
set.seed(101)
sigma=500
N=10000
YD=rnorm(N,23000,1000)
e = rnorm(N, 0, sigma)
C=-1200+0.8*YD + e
```

```
# vygenerujeme náhodnou slozku se strední hodnotou 0 a rozptylem 1
# Populacní regresní funkce, proces, který generuje velikost wage
```

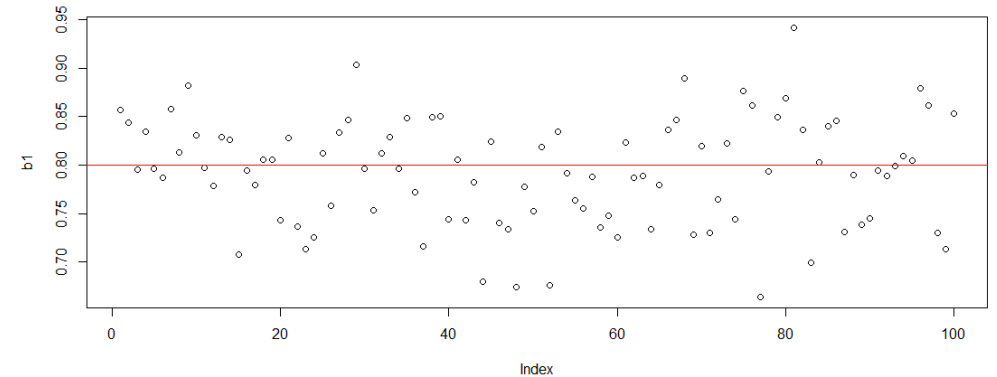
Musíme rozlišovat

$$C = \beta_0 + \beta_1 YD + \varepsilon$$

$$C = b_0 + b_1 \cdot YD + e$$

Populační regresní funkce

výběrová regresní funkce
(sample)



R-ko ukaž

Posledním krokem je statistická verifikace

(samozřejmě je nutná i ekonomická, ale k tomu Vám stačí selský rozum a znalost ekonomické teorie)

Když se vrátíme k odhadu průměrné výšky v populaci.

Vzneseme hypotézu, zdali je průměrná výška mužů rovna 180 cm.

$$H_0: \mu = 180$$

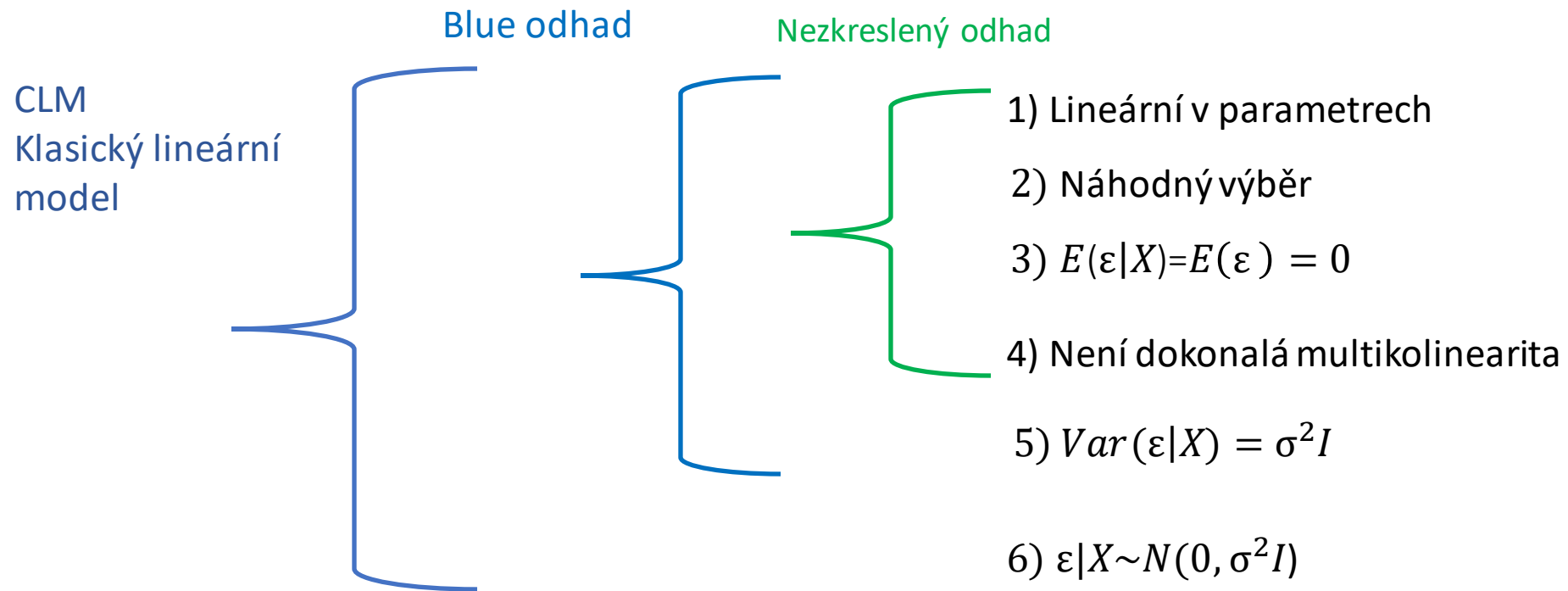
$$H_1: \mu \neq 180$$

U regrese budeme například testovat:

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

Gauss-markovovy předpoklady + chování náhodné složky



Proč GM?

Po splnění předpokladů GM – nemusíme hledat/použít jiných metod pro odhad

Nenajdeme lepší maximálně stejně „dobré“

ve skupině lineárních modelů tedy, ale to nám je zatím jedno

Body 1-5 jsou GM předpoklady

Předpoklad 6 již není GM!!!

1-6 klasický lineární model

Proč $\varepsilon|X \sim N(0, \sigma^2 I)$?

Jelikož pak bude platit $b \sim N(\beta, \Sigma)$

Nezkreslenost je vlastností „estimátoru“

Neříká, že konkrétní odhad je „blízko“ neznámému parametru

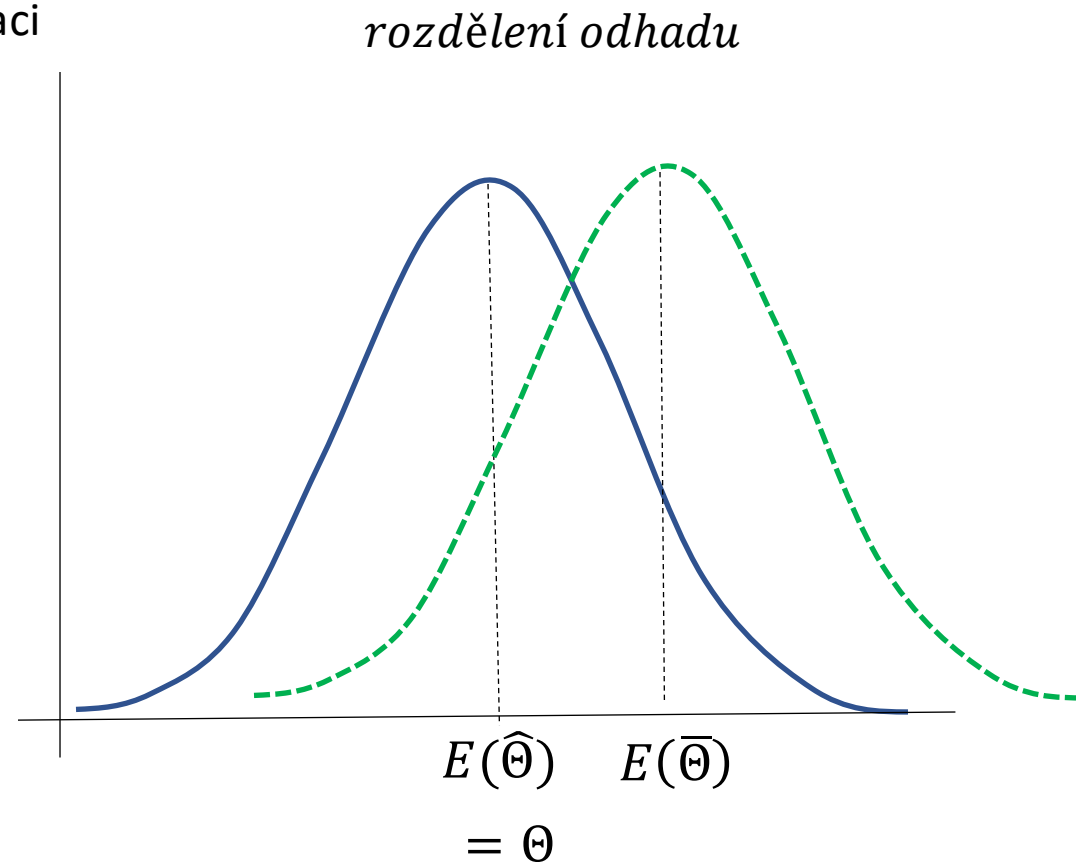
Musíme doufat, že výběrový vzorek je „co nejvíce podobný“ populaci

$\hat{\Theta}$ – odhad, výběrový rozptyl, aritmetický průměr, b – z OLS

Θ – populační rozptyl, střední hodnota, β

$E(\hat{\Theta}) > \Theta$ – odhady jsou nadhodnoceny

$E(\hat{\Theta}) < \Theta$ – odhady jsou podhodnoceny



metoda poskytuje nezkreslený odhad

V reálném životě samozřejmě neznáme hodnotu Θ . Ani si nemůžeme dělat n výběrů.

Tak jak si můžeme vědět, že nejsme úplně mimo?

Nemůžeme, ale víme, že když jsou SPLNĚNY DANÉ PŘEDPOKLADY tak se pohybujeme kolem hodnoty skutečného parametru

Proto tak záleží na rozptylu daného odhadu.

Jako estimátor střední hodnoty populačního rozdělení $E(X)$
Použijeme aritmetický průměr výběrového souboru

$$E(\bar{X}) = \mu?$$

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

$$E(\bar{X}) = E\left(\frac{1}{n} \sum_{i=1}^n X_i\right)$$

$$E(\bar{X}) = \frac{1}{n} E\left(\sum_{i=1}^n X_i\right)$$

$$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n E(X_i)$$

$$E(X_i) = \mu$$

$$E(\bar{X}) = \frac{1}{n} \sum_{i=1}^n \mu$$

$$E(\bar{X}) = \frac{1}{n} n \cdot \mu = \mu$$

NEzkreslený odhad

Předpokládáme výběr z výška $X \sim N(\mu, \sigma^2)$

Každý 1 výběr ze základního souboru, tedy prvek výběrového souboru
můžeme chápat jako 1 z možných hodnot náhodné veličiny $X (X_i)$

Každá z těchto n náhodných veličin má stejné rozdělení $X_i \sim N(\mu, \sigma^2)$

$$\hat{X} = \frac{1}{n-1} \sum_{i=1}^n X_i \quad E(\hat{X}) = \mu?$$

$$E(\hat{X}) = E\left(\frac{1}{n-1} \sum_{i=1}^n X_i\right)$$

$$E(\hat{X}) = \frac{1}{n-1} E\left(\sum_{i=1}^n X_i\right)$$

$$E(\hat{X}) = \frac{1}{n-1} \sum_{i=1}^n E(X_i)$$

$$E(\hat{X}) = \frac{1}{n-1} \sum_{i=1}^n \mu$$

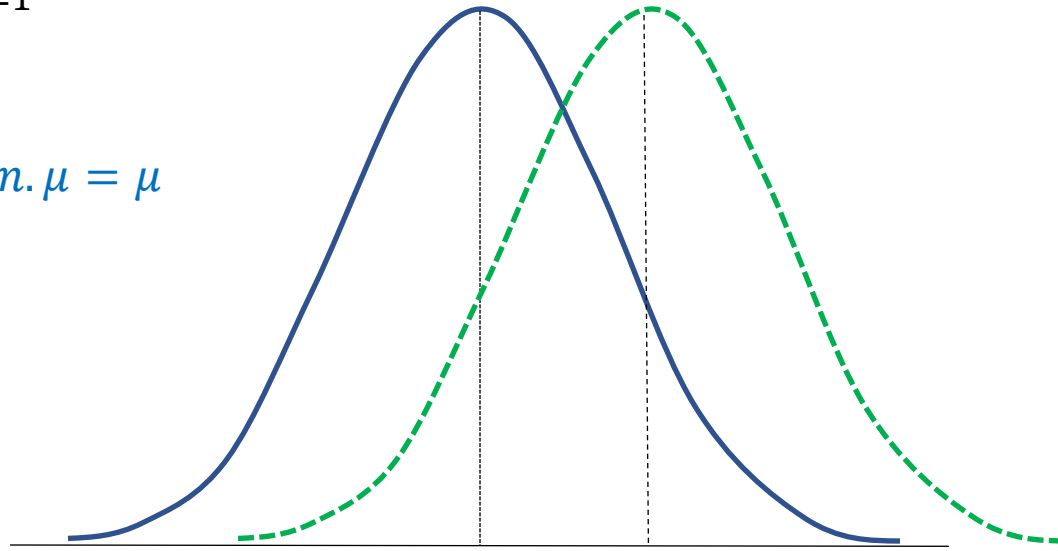
$$E(\hat{X}) \neq \mu$$

$$E(\hat{X}) = \frac{1}{n-1} n \cdot \mu = \frac{n}{n-1} \mu$$

Zkreslený odhad

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

$$E(\bar{X}) = \frac{1}{n} n \cdot \mu = \mu$$



$$E(\bar{X}) = \mu \quad E(\hat{X})$$

$$E(\hat{X}) \neq \mu$$

zkreslený odhad

$$\hat{X} = \frac{1}{n-1} \sum_{i=1}^n X_i$$

$$E(\hat{X}) = \frac{1}{n-1} n \cdot \mu = \frac{n}{n-1} \mu$$

$$E(\bar{X}) < E(\hat{X})$$

Konzistentní odhad

Nezkreslenost pro konečný výběr

Vybíráme n -vzorků o konkrétní velikosti

Konzistentní odhad „large sample“ „asymptotic“

Daný vzorek rozšiřujeme do nekonečna

Zpřesní se odhad pokud budeme zvyšovat velikost vzorku?

Bude menší rozptyl?

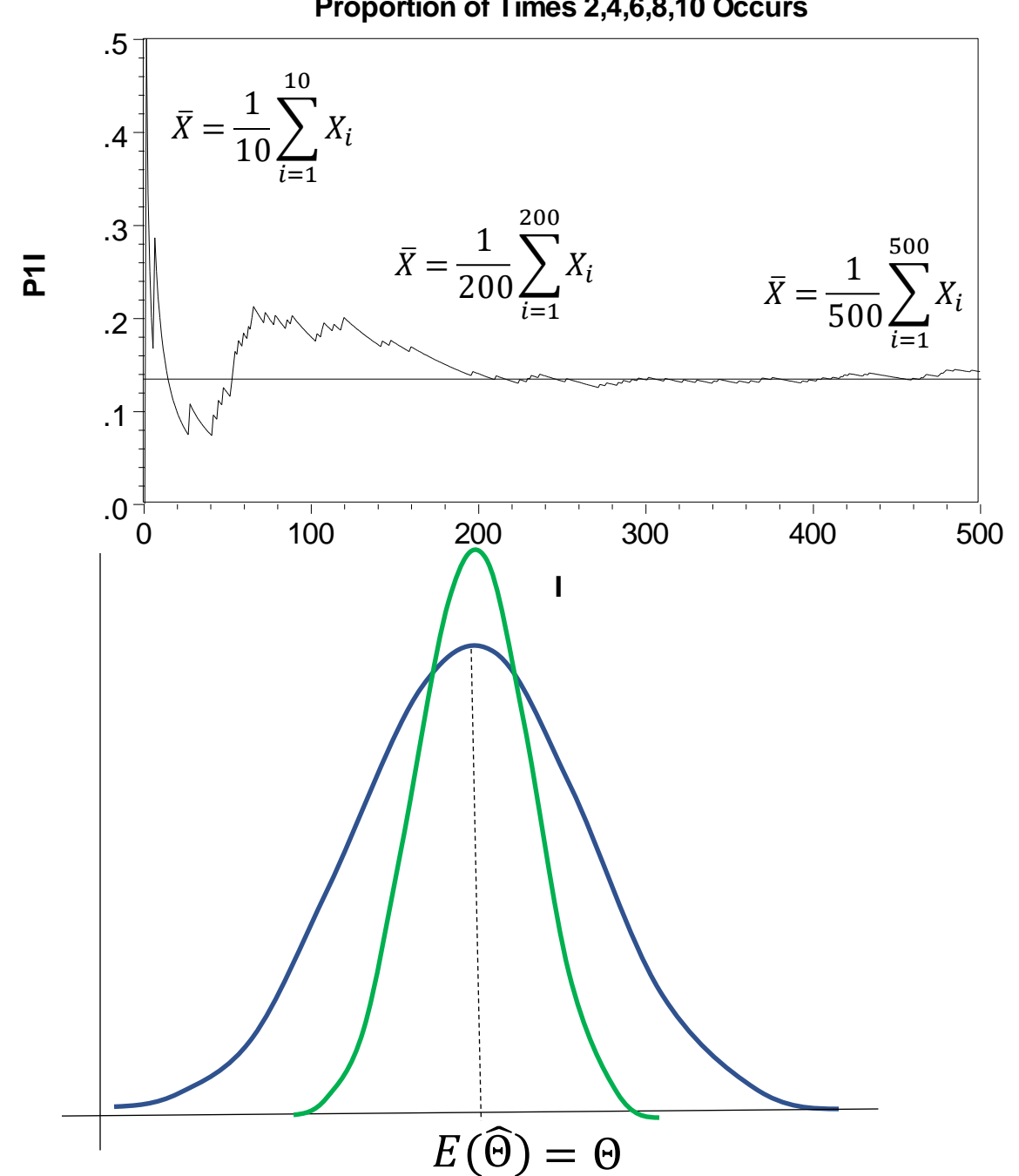
S rostoucím počtem pozorování (n)

konverguje hodnota odhadu

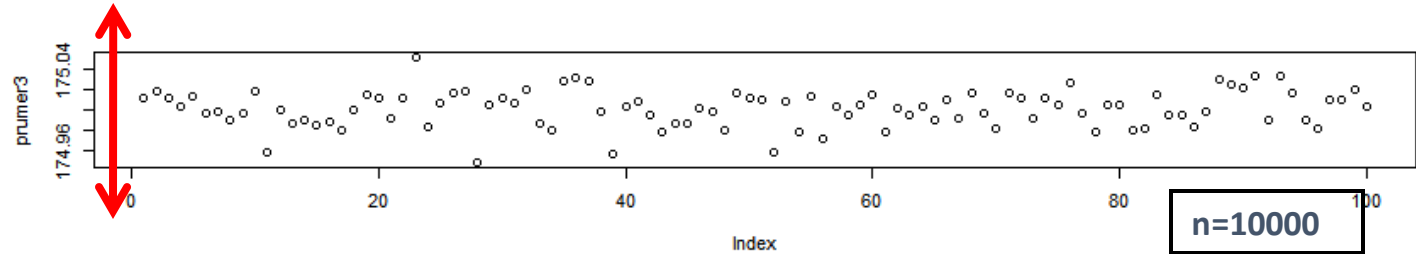
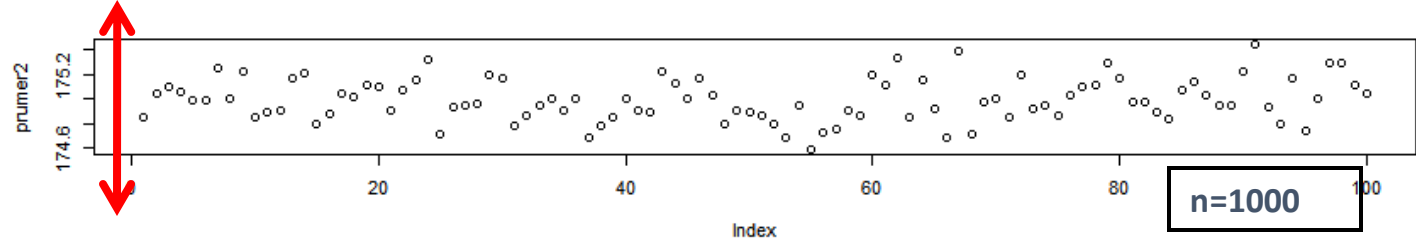
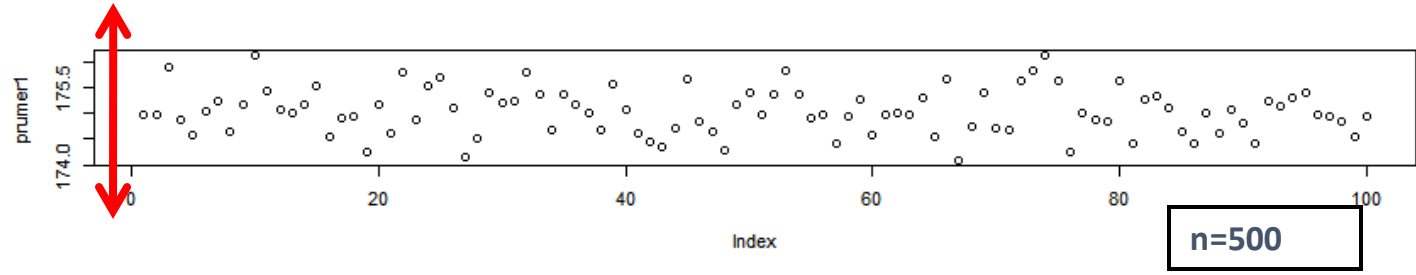
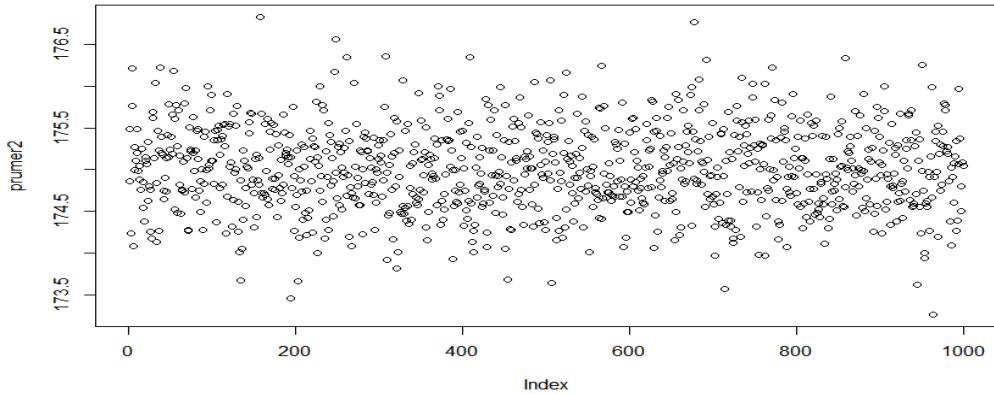
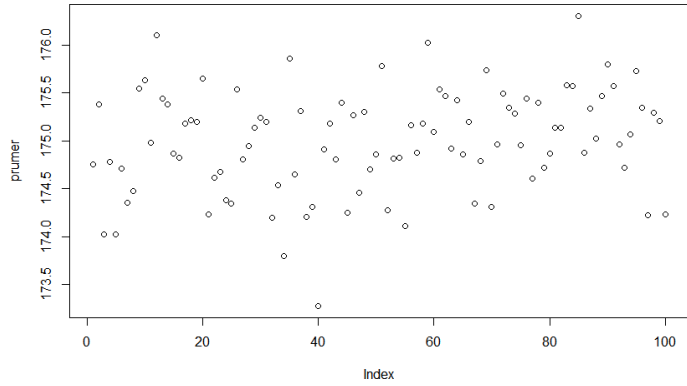
ke skutečné hodnotě parametru základního souboru

Například u časových řad, nebudou dané estimátory poskytovat nevychýlený odhad.

Budou však konzistentní.



Populační průměr (střední hodnota) je 175 cm



Provedu 100-krát výběr o velikosti 500 lidí
a výpočtu pro každý výběr průměru

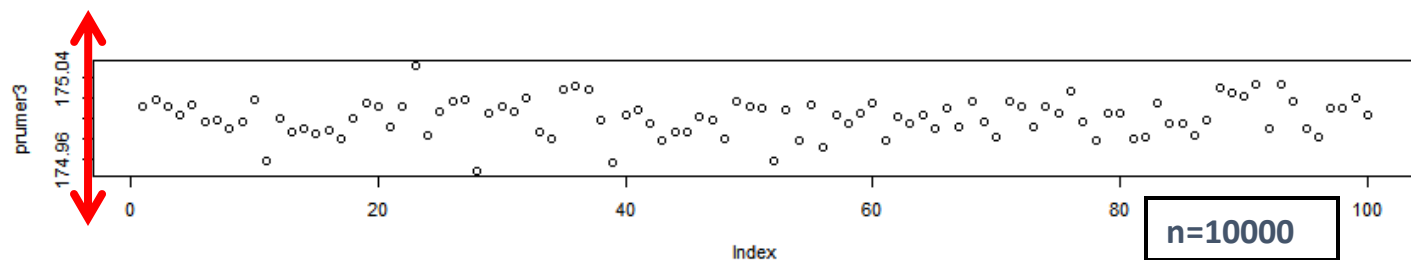
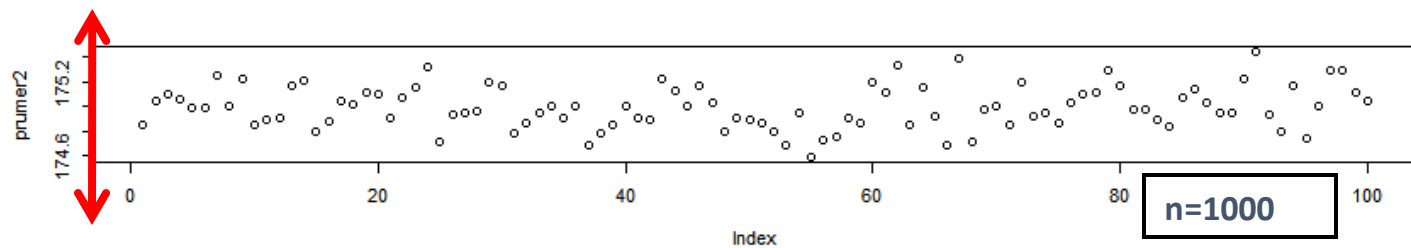
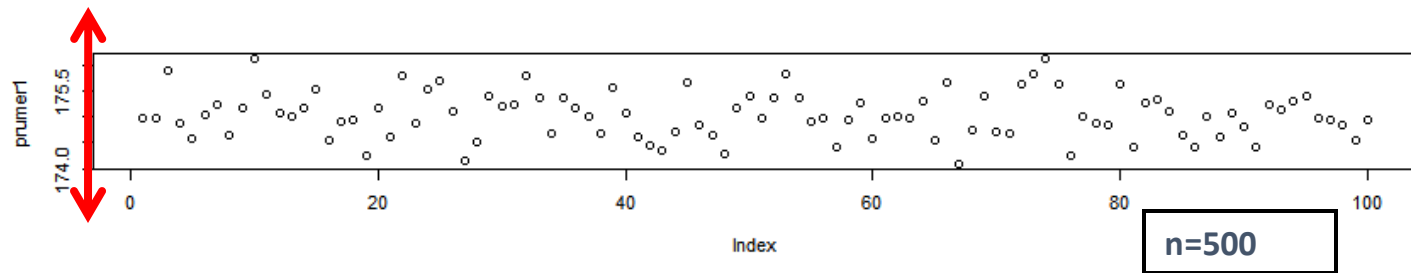
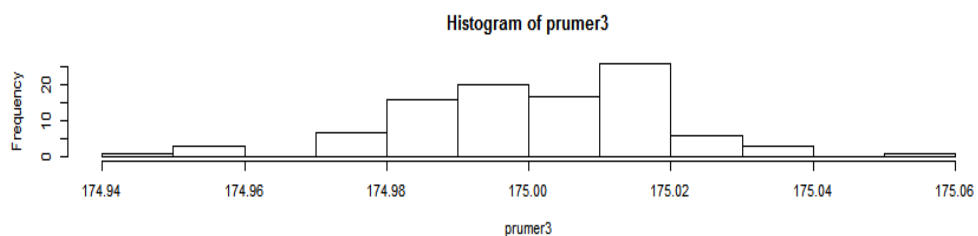
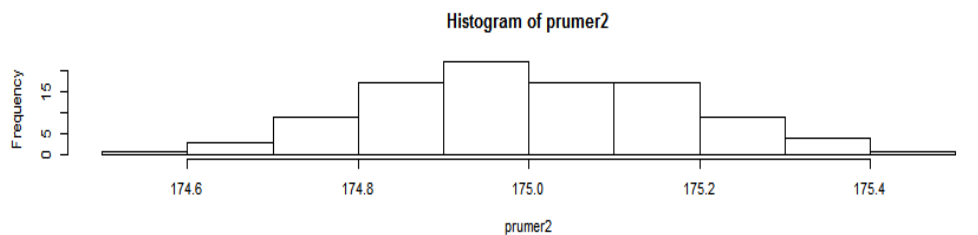
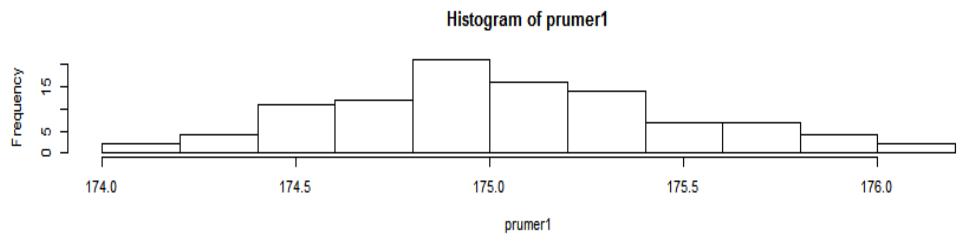
Provedu 1000-krát výběr o velikosti 500 lidí

Provedu 100-krát výběr o velikosti 500 lidí

a výpočtu pro každý výběr průměru

Provedu 100-krát výběr o velikosti 1000 lidí

Provedu 100-krát výběr o velikosti 10 000 lid



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

$n \rightarrow \infty$

$$E(\bar{X}) = \frac{n}{n-1} \mu = \mu$$

Daná metoda poskytuje
Sice zkreslené, ale konzistentní odhady

$$E(\bar{X}) = E\left(\frac{1}{n-1} \sum_{i=1}^n X_i\right)$$

$$E(\bar{X}) = \frac{1}{n-1} E\left(\sum_{i=1}^n X_i\right)$$

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

$$E(\bar{X}) = \frac{1}{n-1} \sum_{i=1}^n \mu$$

$$E(\bar{X}) = \frac{1}{n-1} n \cdot \mu = \frac{n}{n-1} \mu$$

$$E(\bar{X}) \neq \mu$$

Proč to řešíme?

Někdy nenajdeme nezkreslený odhad

Chceme, aby byl alespoň konzistentní

Zkreslený odhad

Vydatný(eficientní, efficient) odhad

Získám 2 NEZKRESLENÉ odhady parametru beta b_1 a \widehat{b}_1

Tedy střední hodnota obou odhadů je shodná

Který odhad použít?

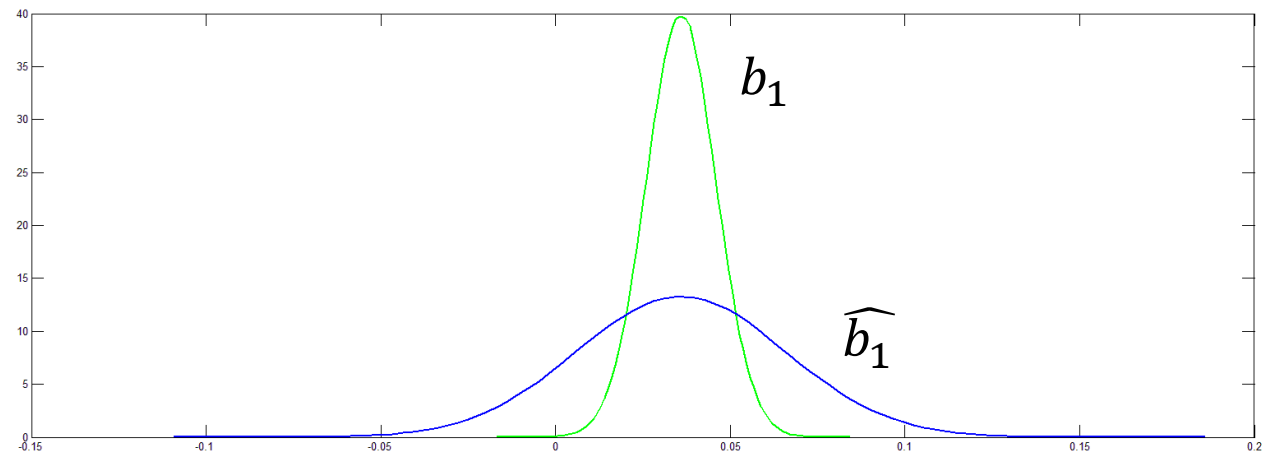
Bude nás zajímat také rozptyl odhadu

Pro \widehat{b}_1 - hodnoty dosahují i záporných hodnot – značný rozptyl

b_1 je vydatnější než \widehat{b}_1

Existuje více nezkraslených statistik – my hledáme tu nejvydatnější

S nejmenším rozptylem!!!



**NECHCI PO VÁS DŮKAZY NEZKRESLENOSTI, ALE MUSÍTE SI UVĚDOMIT, JAK DŮLEŽITÉ JSOU PŘEDPOKLADY.
NEJEN V EKONOMETRII, ALE I V EKONOMII. POKUD NEJSOU SPLNĚNY PŘEDPOKLADY MODELU, ESTIMÁTORU ATD.
TAK VLASTNĚ NEVÍTE, JAK SE DANÝ MODEL, ESTIMÁTOR CHOVÁ. MŮŽETE TAK BÝT ÚPLNĚ MIMO REALITU**

Metoda nejmenších čtverců

Výběrový soubor – „proženeme nějakou funkcí“ (estimator) – získáme odhad (bodový, intervalový)

Výška studentů VŠE

Estimátor

Odhadneme střední hodnotu výšky dospělých lidí

$$\bar{X} = \frac{1}{n} \sum X_i$$

Vztah mezi C a YD

$$b = (X'X)^{-1}X'y$$

Odhadneme parametry podmíněné střední hodnoty

Ok odhadujeme parametry podmíněné střední hodnoty, ale co to znamená?

$$E(y|x) = \beta_0 + \beta_1 \cdot x$$

Vliv x_i na vývoj střední hodnoty y_i

měří změnu střední hodnoty (y) - tedy $E(y|x)$
v závislosti na změně x

$$C = -1800 + 0.75YD + e$$

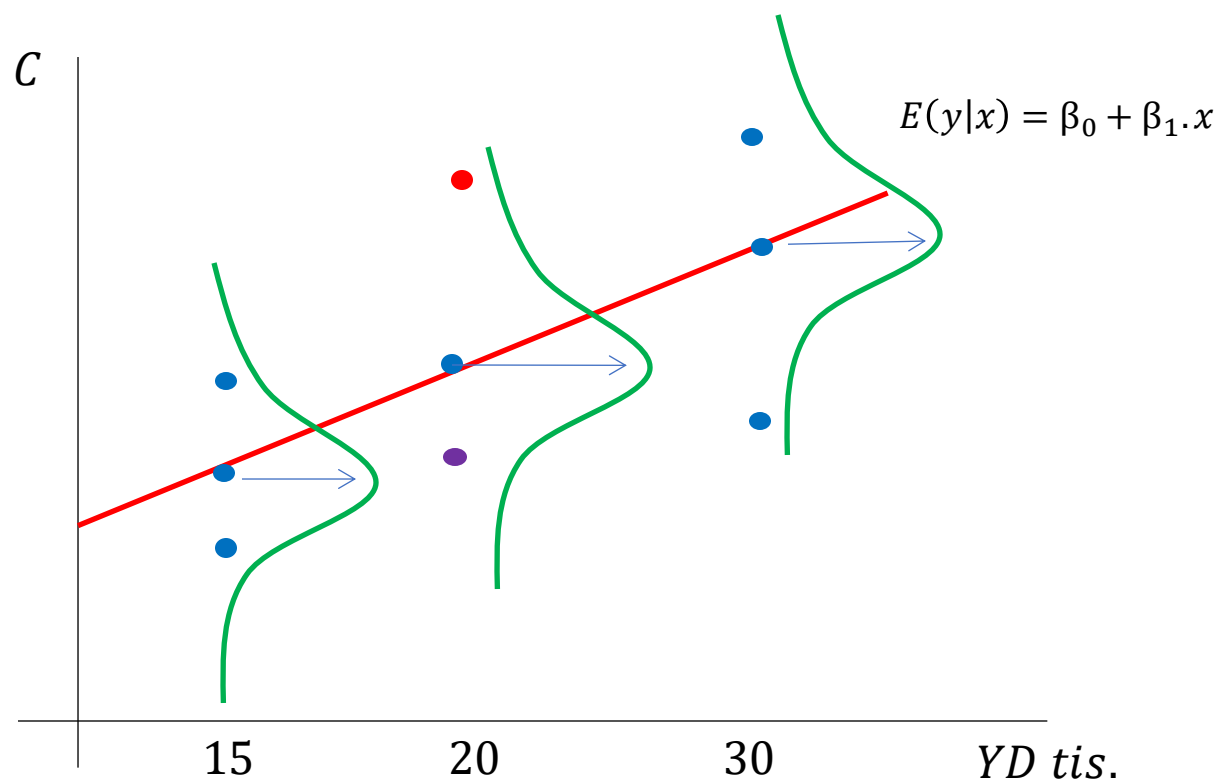
Markéta má plat 18 000 a její spotřeba je 15 000

Měla by však být

$$-1800 + 0.75 \times 18000 = 11\,700$$

Ne?

Myšlenka „v průměru“



Přímková regrese

$$y = \beta_0 + \beta_1 \cdot x + \epsilon$$

$$\hat{y} = b_0 + b_1 \cdot x$$

$$Q = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \dots \min$$

b_0 je odhad β_0

b_1 je odhad β_1

$Q \min$ – hledáme extrém – minimum

Tedy takové $b_{0,1}$, které budou minimalizovat funkci Q

Index i představuje i -té pozorování
Mzdu, vzdělání Natálie

Hledáme takové hodnoty $b_{0,1}$

Které budou minimalizovat vzdálenost $\sum_{i=1}^n (y_i - \hat{y}_i)^2$

$$b = (X'X)^{-1}X'y$$

$$y = X\beta + \epsilon$$

$$b = (X'X)^{-1}X'(X\beta + \epsilon)$$

$$b = (X'X)^{-1}X'X\beta + (X'X)^{-1}X'\epsilon$$

$$b = \beta + (X'X)^{-1}X'\epsilon$$

$$E(b) = \beta + E[(X'X)^{-1}X'\epsilon]$$

$$E(b) = \beta + (X'X)^{-1}X'E(\epsilon)$$

pokud $E(\epsilon) = 0$

$E(b) = \beta$ – nezkreslený odhad

$$\text{Var}(b) = \text{Var}(\beta + (X'X)^{-1}X'\epsilon)$$

$$\text{Var}(b) = \text{Var}[(X'X)^{-1}X'\epsilon]$$

$$\text{Var}(AX) = A\text{Var}(X)A'$$

$$\text{Var}(b) = (X'X)^{-1}X'\text{Var}(\epsilon)\left((X'X)^{-1}X'\right)'$$

$$\text{Var}(b) = (X'X)^{-1}X'\text{Var}(\epsilon)X(X'X)^{-1}$$

$$\text{Var}(b) = (X'X)^{-1}X'\sigma^2IX(X'X)^{-1}$$

$$(X'X)^{-1}X'X = I$$

σ^2 – konstanta, lze vytknout před

$$\text{Var}(b) = \sigma^2(X'X)^{-1}$$

Odvodíme nyní Gauss-Markov předpoklady

Abychom získali nezkreslený odhad:

Model je lineární v parametrech- musí platit, jinak bychom měli nelineární metodu nejmenších čtverců

$$E(b) = \beta + (X'X)^{-1}X'E(\epsilon)$$

$$\text{pokud } E(\epsilon) = 0$$

$$E(b) = \beta - \text{nezkreslený odhad}$$

A máme předpoklad!

Dále máme matici $(X'X)^{-1}$.

Abychom mohli vytvořit inverzní matici, musí platit, že čtvercová matice $(X'X)$ o velikosti $n \times n$ je tzv. regulární matice

Tedy hodnota matice je rovna n .

Řádky i sloupce jsou lineárně nezávislé

Kdyby tato podmínka neplatila, tak nejsme schopni spočítat $b = (X'X)^{-1}X'y$

Navíc chceme mít BLUE odhad – nejlepší s nejmenším rozptylem

$$\text{Var}(b) = (X'X)^{-1}X'\text{Var}(\epsilon)X(X'X)^{-1}$$

$$\text{Var}(b) = (X'X)^{-1}X'\sigma^2IX(X'X)^{-1}$$

$$\text{Var}(b) = \sigma^2(X'X)^{-1}$$

Pokud $\text{Var}(\epsilon) \neq \sigma^2I$, tak potom ani

$$\text{Var}(b) \neq \sigma^2(X'X)^{-1}$$

OLS již nebude BLUE. Bude existovat nějaká dodatečná informace o tvaru $\text{Var}(\epsilon)$, kterou se kterou bude pracovat jiná metoda.

Pamatujete jaká?

Proč nás tak zajímá rozptyl a co jej ovlivňuje?

$$\text{Var}(b) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

$$\text{Var}(b) = \widehat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1}$$

$$\text{Var}(b_j) = \frac{\sigma^2}{\sum (x_j - \bar{x}_j)^2 (1 - R_j^2)}$$

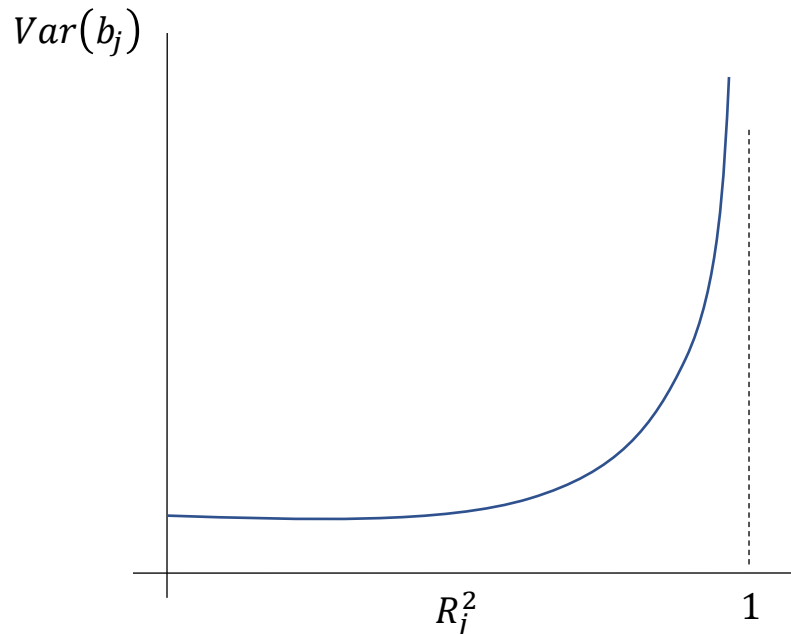
$$\text{Var}(b_j) = \frac{\widehat{\sigma}^2}{\sum (x_j - \bar{x}_j)^2 (1 - R_j^2)}$$

Pozor R_j^2 není R^2 z výstupu programu!!

Čím větší bude R_j^2 - tím větší bude rozptyl odhadu

R_j^2 - R^2 z regrese x_j na všechny ostatní X + intercept

$$y = \beta_0 + \beta_1 \cdot x_1 + \beta_2 \cdot x_2 + \epsilon \quad x_1 = \gamma_0 + \gamma_2 \cdot x_2 + v \quad R_1^2$$



$$R_j^2 \rightarrow 1 \quad \text{Var}(b_j) \rightarrow \infty$$

Problém u testování hypotéz

Interpretace R^2 - kolik variability v x_1
Se nám podařilo vysvětlit pomocí x_2
Čím více - tím větší je vzájemný vztah

- Počet pozorování
- Variabilita x
- R_j^2
- Rozptyl náhodné složky σ^2

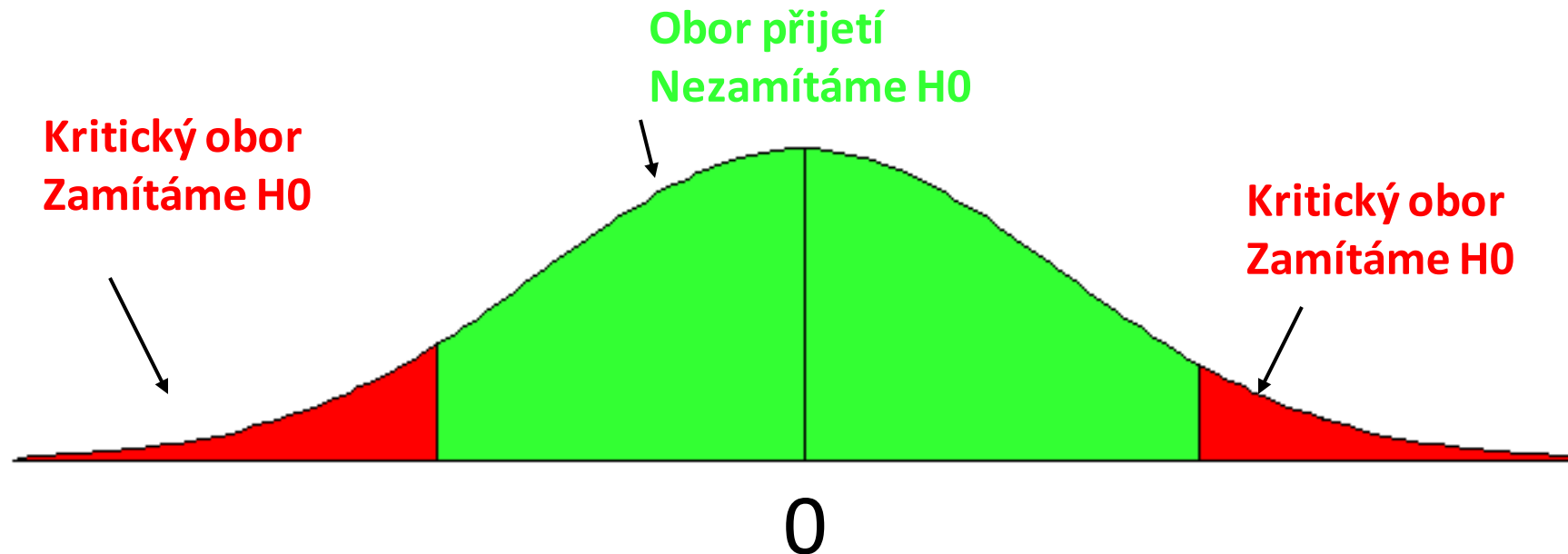
$$\text{spotřeba} = \beta_0 + \beta_1 \cdot \text{mzda} + \beta_2 \cdot \text{kap. vynos} + \beta_3 \cdot \text{urok. mira} + \varepsilon$$

$$t = \frac{b_3}{\sqrt{\text{Var}(b_3)}} \sim t(n - k - 1)$$

Jaké hodnoty svědčí pro H1?
„okrajové“ nebo blízko nuly?

$$H_0: \beta_3 = 0$$

$$H_1: \beta_3 \neq 0$$



Oživit si

Koeficient determinace vs. Upravený koeficient determinace

Log modely umět interpretovat

Dummy proměnné

Náhodná složka vs. Residuum

$$y = b_0 + b_1x + e \text{ vs. } \hat{y} = b_0 + b_1x$$

OPRAVDU CHÁPAT TESTY HYPOTÉZ!!!!



EVROPSKÁ UNIE
Evropské strukturální a investiční fondy
Operační program Výzkum, vývoj a vzdělávání



Toto dílo podléhá licenci Creative Commons
Uveďte původ – Zachovejte licenci 4.0 Mezinárodní.

