

# homework

December 13, 2018

## 1 case0222 - Cholesterol In Urban And Rural Guatemalans

<https://www.rdocumentation.org/packages/Sleuth2/versions/2.0-4/topics/ex0222>

### 1.0.1 Popis dat

Dataset pochází ze studie provedené na guatemalských indiánech. Míra cholesterolu byla změna celkem 94 jedincem a byl zaznamenán jejich pvod. Bylo nameno 49 pozorování na venkov a 45 ve mst.

### 1.0.2 Formát

Dataframe obsahuje 94 pozorování na následujících 2 promnných:

- **Cholesterol** - Mnoství cholesterolu v krvi lovka (v mg/l).
- **Group** - Promnná obsahující hodnoty "Rural" a "Urban" oznaující, jestli je subjekt z venkova, nebo z msta.

### 1.0.3 Zdroj

Ramsey, F.L. and Schafer, D.W. (2002). The Statistical Sleuth: A Course in Methods of Data Analysis (2nd ed), Duxbury.

```
In [2]: library(Sleuth2)
        str(ex0222)
```

```
'data.frame':      94 obs. of  2 variables:
 $ Cholesterol: num  133 134 155 170 175 179 181 184 188 189 ...
 $ Group      : Factor w/ 2 levels "Rural","Urban": 2 2 2 2 2 2 2 2 2 2 ...
```

- [x] (1b) Nacte datový soubor a rozdlte sledovanou promnnou na písluné dv pozorované skupiny. Data strun popíte. Pro kadu skupinu zvlá odhadnte stední hodnotu, rozptyl a medián písluného rozdlení.

```
In [3]: rural <- subset(ex0222, Group=="Rural", Cholesterol, drop=TRUE)
        urban <- subset(ex0222, Group=="Urban", Cholesterol, drop=TRUE)
```

*Subjekty z venkova:*

```
In [4]: cat("Rural area indians:\n")
        cat("EX =", mean(rural), "\n")
        cat("varX =", var(rural), "\n")
        cat("median =", median(rural))
```

```
Rural area indians:
EX = 157
varX = 1008.458
median = 152
```

*Subjekty z msta:*

```
In [6]: cat("Urban area indians:\n")
        cat("EX =", mean(urban), "\n")
        cat("varX =", var(urban), "\n")
        cat("median =", median(urban))
```

```
Urban area indians:
EX = 216.8667
varX = 1593.618
median = 206
```

- [x] (1b) Pro každou skupinu zvlá odhadněte hustotu a distribuní funkci pomocí histogramu a empirické distribuní funkce.

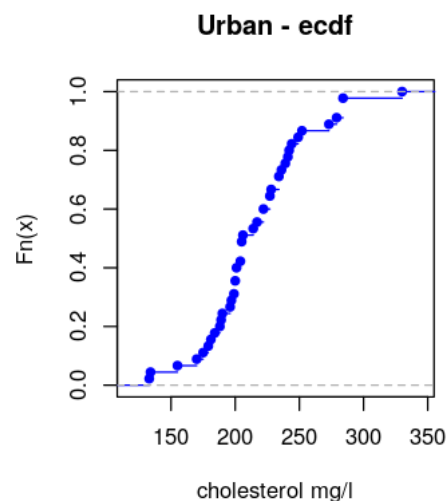
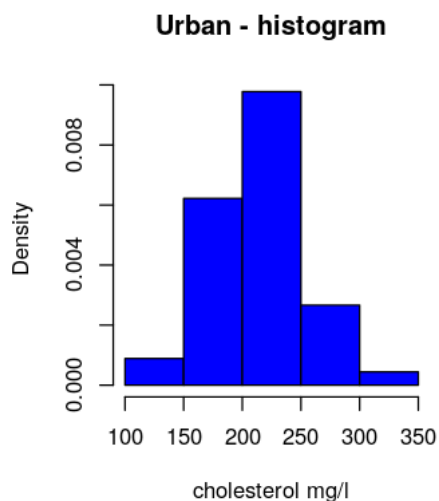
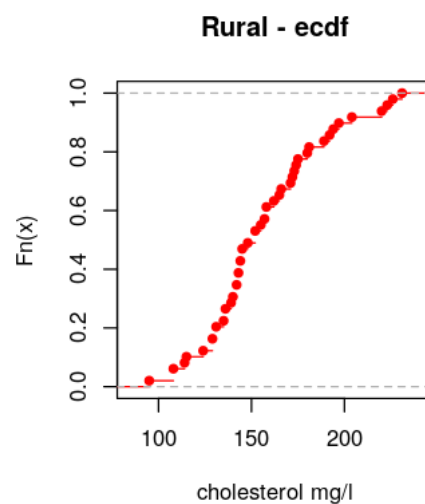
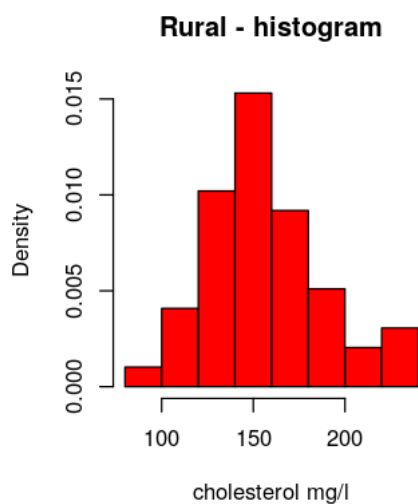
```
In [18]: par(mfrow = c(2, 2), pty="s")
```

```
# rural
```

```
hist(rural, col="red", main="Rural - histogram", probability=T, xlab="cholesterol mg/l")
plot.ecdf(rural, col="red", main="Rural - ecdf", xlab="cholesterol mg/l")
```

```
# urban
```

```
hist(urban, col="blue", main = "Urban - histogram", probability=T, xlab="cholesterol mg/l")
plot.ecdf(urban, col="blue", main="Urban - ecdf", xlab="cholesterol mg/l")
```



- [ ] (3b) Pro každou skupinu zvlášť najdte nejbližší rozdělení: Odhadněte parametry normálního, exponenciálního a rovnoměrného rozdělení. Zanešte příslušné hustoty s odhadnutými parametry do grafu histogramu. Diskutujte, které z rozdělení odpovídá pozorovaným datům nejlépe.

**Odhady rozdělení** Pro provedení odhadu jsou využity funkce `mean()` a `sd()` zabudované do standardní knihovny jazyka R. Odhad je proveden shodně i pro množinu urban.

### Normální rozdělení

```
EX = mean(rural)
s = sd(rural)
```

## Exponenciální rozdělení

```
lambda = 1/mean(rural)
```

## Uniformní rozdělení

```
a = min(rural)
```

```
b = max(rural)
```

```
In [16]: x <- seq(min(rural), max(rural), length=100)
```

```
# hodnoty pro jednotlivá rozložení
```

```
y_norm <- dnorm(x, mean=mean(rural), sd=sd(rural))
```

```
y_exp <- dexp(x, 1/mean(rural))
```

```
y_unif <- dunif(x, min=min(rural), max=max(rural))
```

```
hist(rural, probability=T, main="Rural - odhady rozdělení", xlab="cholesterol", ylab=
```

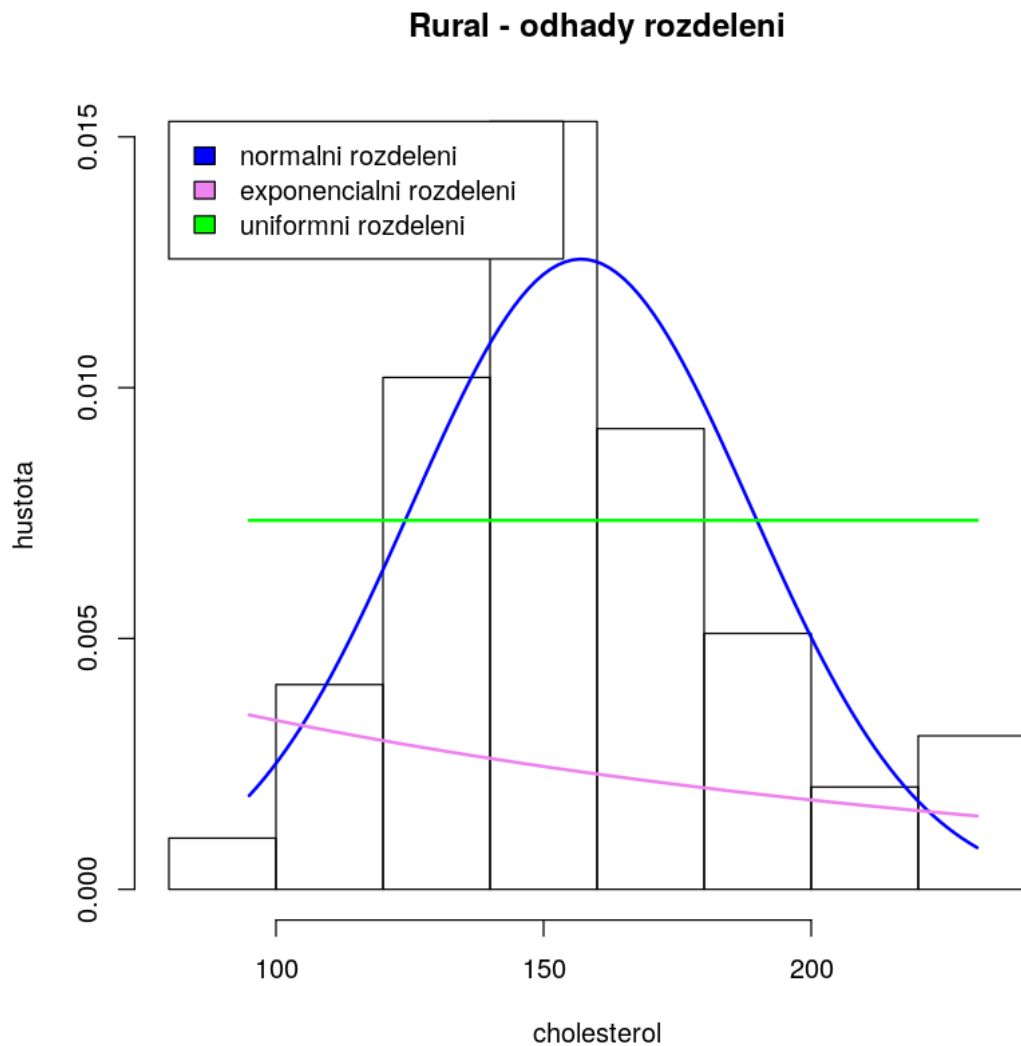
```
lines(x, y_norm, col="blue", lwd=2)
```

```
lines(x, y_exp, col="violet", lwd=2)
```

```
lines(x, y_unif, col="green", lwd=2)
```

```
legend("topleft", inset=0.037, fill=c("blue","violet","green"),
```

```
legend=c("normalní rozdělení", "exponenciální rozdělení", "uniformní rozdělení")
```



```
In [64]: x <- seq(min(urban), max(urban), length=40)
```

```
y_norm <- dnorm(x, mean=mean(urban), sd=sd(urban))
```

```
y_exp <- dexp(x, 1/mean(urban))
```

```
y_unif <- dunif(x, min=min(urban), max=max(urban))
```

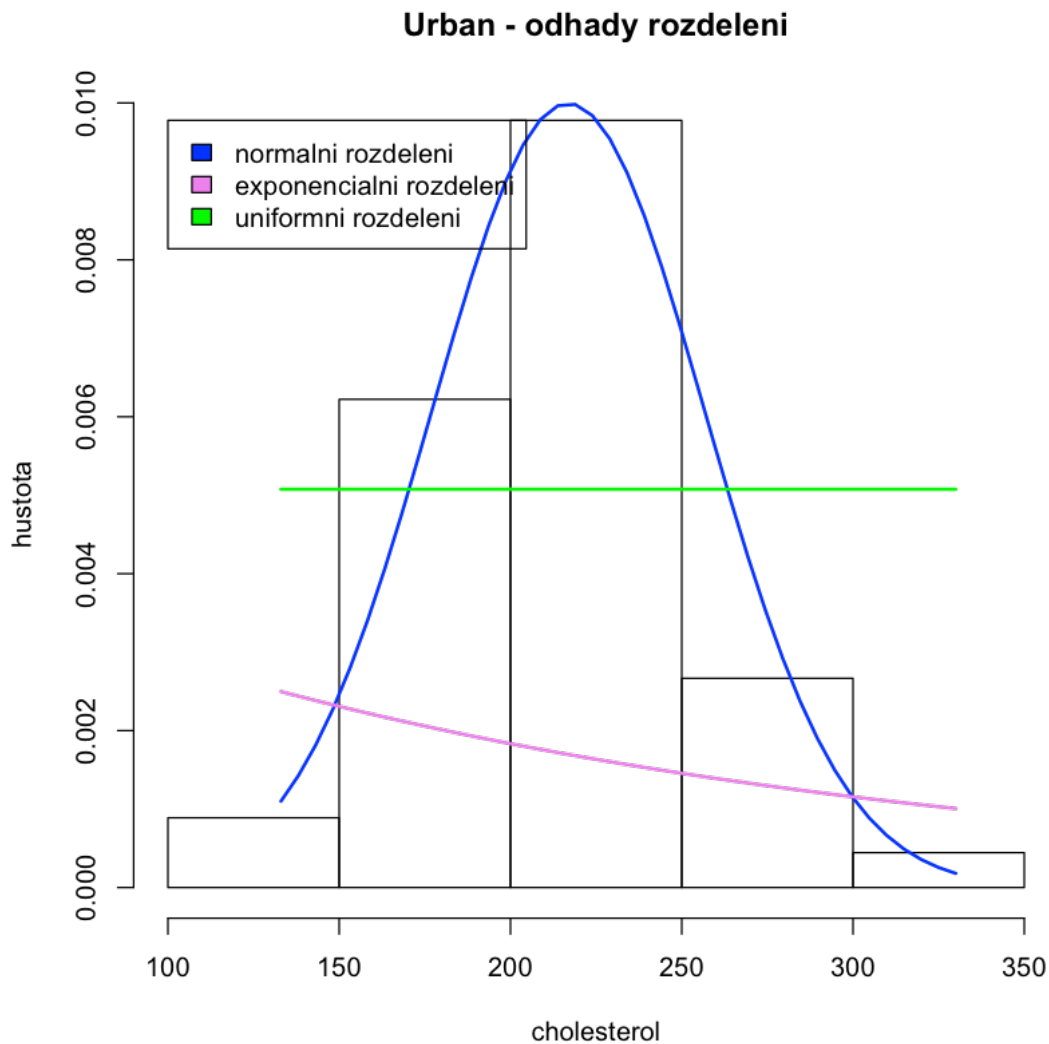
```
hist(urban, probability=T, main="Urban - odhady rozdeleni", xlab="cholesterol", ylab="hustota")
```

```
lines(x, y_norm, col="blue", lwd=2)
```

```
lines(x, y_exp, col="violet", lwd=2)
```

```
lines(x, y_unif, col="green", lwd=2)
```

```
legend("topleft", inset=0.037, fill=c("blue","violet","green"),  
      legend=c("normalni rozdeleni", "exponencialni rozdeleni", "uniformni rozdeleni"))
```



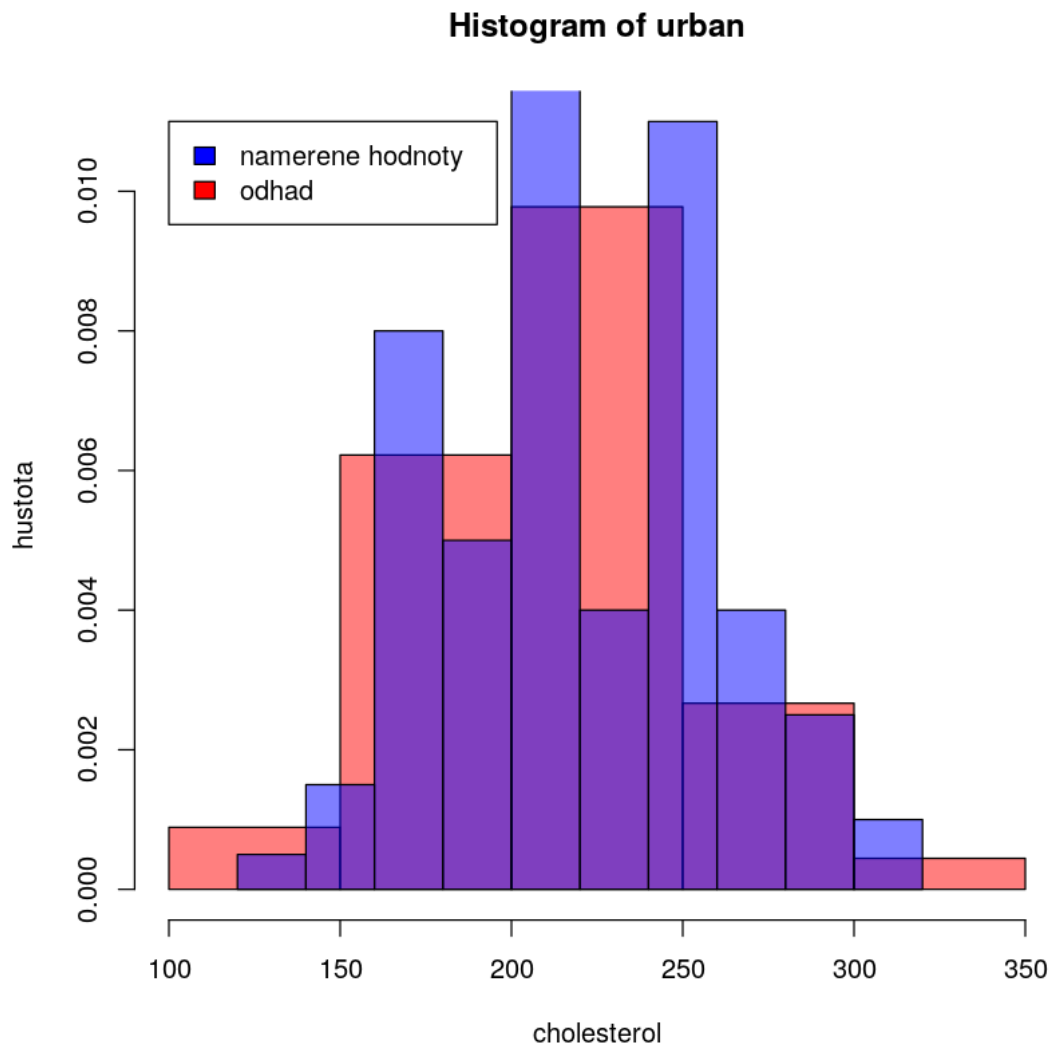
- [ ] (1b) Pro každou skupinu zvlášť vygenerujte náhodný výběr o 100 hodnotách z rozdělení, které jste zvolili jako nejbližší, s parametry odhadnutými v předchozím bod. Porovnejte histogram simulovaných hodnot s pozorovanými daty.

Na náš dataset se nejvíce hodí normální rozdělení. Parametry jsme odhadli pomocí knihovnických funkcí `mean()` a `sd()`. Samotný náhodný výběr jsme vygenerovali následujícím příkazem:

```
rnorm(100, mean=mean(urban), sd=sd(urban))
```

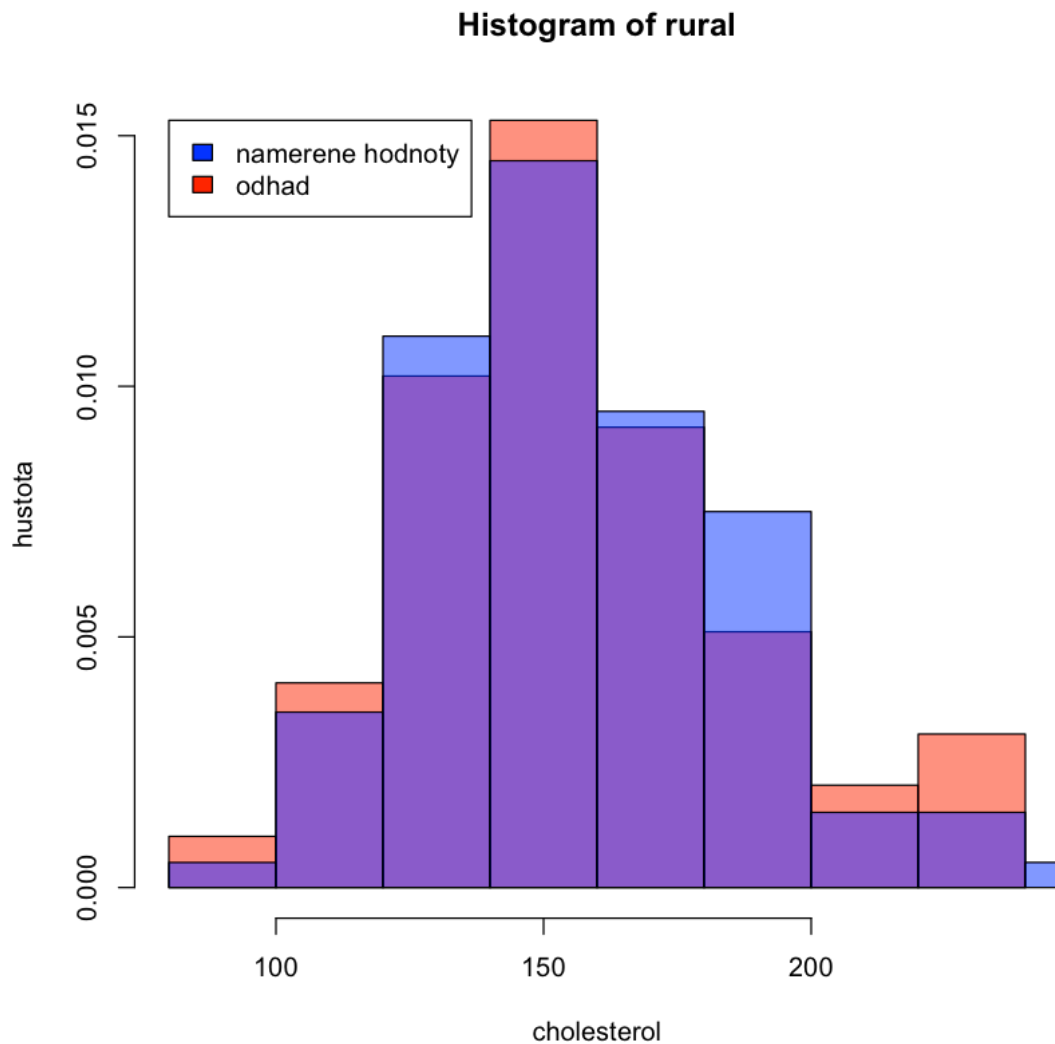
```
In [19]: y <- rnorm(100, mean=mean(urban), sd=sd(urban))
```

```
hist(urban, probability=T, col=rgb(1, 0, 0, 0.5), ylim=c(0, 0.011), xlab="cholesterol")
hist(y, probability=T, col=rgb(0, 0, 1, 0.5), add=T)
legend("topleft", inset=0.037, fill=c("blue", "red"),
      legend=c("namerene hodnoty", "odhad"))
```



```
In [66]: y <- rnorm(100, mean=mean(rural), sd=sd(rural))
```

```
hist(rural, probability=T, col=rgb(1, 0, 0, 0.5), xlab="cholesterol", ylab="hustota")
hist(y, probability=T, col=rgb(0, 0, 1, 0.5), add=T)
legend("topleft", inset=0.037, fill=c("blue", "red"),
      legend=c("namerene hodnoty", "odhad"))
```



- [ ] (1b) Pro každou skupinu zvlá spočítejte oboustranný 95% konfidenní interval pro střední hodnotu.

In [6]: `length(rural)`

49

- [ ] (1b) Pro každou skupinu zvlá otestujte na hladin významnosti 5% hypotézu, zda je střední hodnota rovná hodnot K (parametr úlohy), proti oboustranné alternativ. Můžete použít bu výsledek z předchozího bodu, nebo výstup z příslušné vestavěné funkce vašeho softwaru.

In [144]: `# TODO`



- [ ] (2b) Na hladin spolehlivosti 5% otestujte, jestli mají pozorované skupiny stejnou střední hodnotu. Typ testu a alternativy stanovte tak, aby vaše volba nejlépe korespondovala s povahou zkoumaného problému.

In [145]: # *TODO*