

Hladina cholesterolu dle prostředí

- Naimportujeme knihovny.

```
In [1]: import pandas as pd
import matplotlib
import plt
import numpy as np
from scipy.stats import expon, norm, uniform
from scipy.stats import t, f
```

- Vypočet čísla datasetu.

```
In [2]: K = 18
n = 11
M = (((K+L)*47) % 11)+1.
print('M = ', M)

M = 11.0
```

- Pro 11 bude dataset ex0222 - hladina cholesterolu dle prostředí.

Úloha 1

Načítáme datový soubor a rozdělíme sledovanou proměnnou na příslušné dvě pozorované skupiny. Stručně popište data a zkoumaný problém. Pro každou skupinu zvlášť odhadněte střední hodnotu, rozptyl a medián příslušného rozdělení.

- Dataset ex0222 reprezentuje výsledky průzkumu ve kterém zkoumali hladinu cholesterolu u guatemalských indiánů bydlících na venkově a ve městě.
- Sloupec Cholesterol reprezentuje hladinu cholesterolu v mg/l.
- Sloupec Group může být nebo Urban nebo Rural.

- Načteme data pomocí funkce read_csv() a podíváme se na data pomocí display().

```
In [3]: df = pd.read_csv('my_data.csv', sep=';')
display(df)
```

Cholesterol	Group
133	Urban
134	Urban
155	Urban
170	Urban
175	Urban
...	...
204	Rural
220	Rural
223	Rural
226	Rural
231	Rural
...	...
204	Rural
220	Rural
223	Rural
226	Rural
231	Rural

94 rows x 2 columns

- Budeme mít dvě pozorované skupiny Rural a Urban. Celkem máme 94 záznamu. Podíváme se na počet záznamu v obou skupinách.

```
In [4]: df['Group'].value_counts()

Out[4]:
Rural    49
Urban    45
Name: Group, dtype: int64
```

- Počet pozorování pro skupinu Rural je 49 a pro skupinu Urban je 45
- Rozdělíme původní dataset podle příznaku Group na dvě skupiny Rural a Urban.
- Pomocí funkce info() ověříme, že sloupec Cholesterol má číselnou reprezentaci (není typu object).

```
In [5]: df_Rural = df[ df['Group']=='Rural']
df_Urban = df[ df['Group']=='Urban']

df_Rural.info()
df_Urban.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 49 entries, 45 to 93
Data columns (total 2 columns):
# Column Non-Null Count Dtype
---
0 Cholesterol 49 non-null int64
1 Group 49 non-null object
dtypes: int64(1), object(1)
memory usage: 1.1+ KB

<class 'pandas.core.frame.DataFrame'>
Int64Index: 45 entries, 0 to 44
Data columns (total 2 columns):
# Column Non-Null Count Dtype
---
0 Cholesterol 45 non-null int64
1 Group 45 non-null object
dtypes: int64(1), object(1)
memory usage: 1.1+ KB
```

Vytvoříme pomocné funkce pro výpočet odhadu střední hodnoty, rozptylu a mediánu

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

- Střední hodnotu odhadneme pomocí výběrového průměru.

```
In [6]: def SampleMean(data):
return data['Cholesterol'].mean()

def SampleVariance(data):
return data['Cholesterol'].var()

def SampleMedian(data):
return data['Cholesterol'].median()
```

- Využijeme naše funkce na odhad střední hodnoty, rozptylu a mediánu pro skupinu Rural

```
In [9]: print("Odhad střední hodnoty 'Cholesterolu' v Rural: ", SampleMean(df_Rural))
print("Odhad rozptylu 'Cholesterolu' v Rural: ", round(SampleVariance(df_Rural), 2))
print("Odhad mediány 'Cholesterolu' v Rural: ", SampleMedian(df_Rural))

Odhad střední hodnoty 'Cholesterolu' v Rural: 157.0
Odhad rozptylu 'Cholesterolu' v Rural: 1006.46
Odhad mediány 'Cholesterolu' v Rural: 152.0
```

- Využijeme naše funkce na odhad střední hodnoty, rozptylu a mediánu pro skupinu Urban

```
In [10]: print("Odhad střední hodnoty 'Cholesterolu' v Urban: ", round(SampleMean(df_Urban), 2))
print("Odhad rozptylu 'Cholesterolu' v Urban: ", round(SampleVariance(df_Urban), 2))
print("Odhad mediány 'Cholesterolu' v Urban: ", SampleMedian(df_Urban))

Odhad střední hodnoty 'Cholesterolu' v Urban: 216.87
Odhad rozptylu 'Cholesterolu' v Urban: 1593.62
Odhad mediány 'Cholesterolu' v Urban: 206.0
```

Úloha 2

Pro každou skupinu zvlášť odhadněte hustotu a distribuční funkci pomocí histogramu a empirické distribuční funkce.

- Odhad hustoty uděláme pomocí histogramu. Použijeme na to metodu pandas plot.hist() a nastavíme parametr density=True aby platila normalizační podmínka (škáluje na interval <0,1).
- Odhad distribuční funkci uděláme pomocí empirické distribuční funkce. Použijeme na to taky metodu pandas plot.hist() a navíc nastavíme parametr cumulative=True.

```
In [11]: df_Rural.plot.hist(bins=15, grid=True, title="Odhad hustoty pro Rural", density=True, edgecolor='b', alpha=0.5)
df_Rural.plot.hist(bins=15, grid=True, title="Odhad distribuční f-ci pro Rural", density=True, cumulative=True,
```

```
Out[11]: <AxesSubplot:title='center': 'Odhad distribuční f-ci pro Rural', ylabel='Frequency'>
```

```
In [12]: df_Urban.plot.hist(bins=15, legend=False, grid=True, title="Odhad hustoty pro Urban", density=True, edgecolor='b')
df_Urban.plot.hist(bins=15, legend=False, grid=True, title="Odhad distribuční f-ci pro Urban", density=True, cu
```

```
Out[12]: <AxesSubplot:title='center': 'Odhad distribuční f-ci pro Urban', ylabel='Frequency'>
```

Úloha 3

Pro každou skupinu zvlášť najděte nejbližší rozdělení: Odhadněte parametry normálního, exponenciálního a rovnoměrného rozdělení. Zanepte příslušné hustoty s odhadnutými parametry do grafu histogramu. Diskutujte, které z rozdělení odpovídá pozorovaným datům nejlépe.

Použijeme metodu Maximální věrohodnosti pro odhad parametrů normalního, exponenciálního a rovnoměrného rozdělení.

- Pro normální rozdělení odhadneme střední hodnotu pomocí výběrového průměru a rozptyl pomocí vzorečku z přednášky:

$$\hat{\sigma}_n^2 = \hat{\tau}_n = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

- Máme na to pomocnou funkci EstimateMean_Std()

```
In [13]: # Normální rozdělení
def EstimateMean_Std(data):
    my_sum = 0
    mean = SampleMean(data)
    for i in data['Cholesterol']:
        my_sum += (i-mean)**2
    my_var = my_sum/data['Cholesterol'].count()
    my_std = np.sqrt(my_var)
    return mean, my_std

# Exp. rozdělení
my_lambda = 1/mean
print("Exponenciální rozdělení: ")
print("\tlambda: ", round(my_lambda, 4), "\n")
y_exp = expon.pdf(x, scale = mean)
ax.plot(x, y_exp, color = 'g', label = "Exponenciální rozdělení")

#Rovnoměrné rozdělení
a = data['Cholesterol'].min()
b = data['Cholesterol'].max()
print("Rovnoměrné rozdělení: ")
print("\ta: ", a)
print("\tb: ", b)
y_unif = uniform.pdf(x, a, b - a)
ax.plot(x, y_unif, color = 'orange', label = "Rovnoměrné rozdělení")
ax.legend()
```

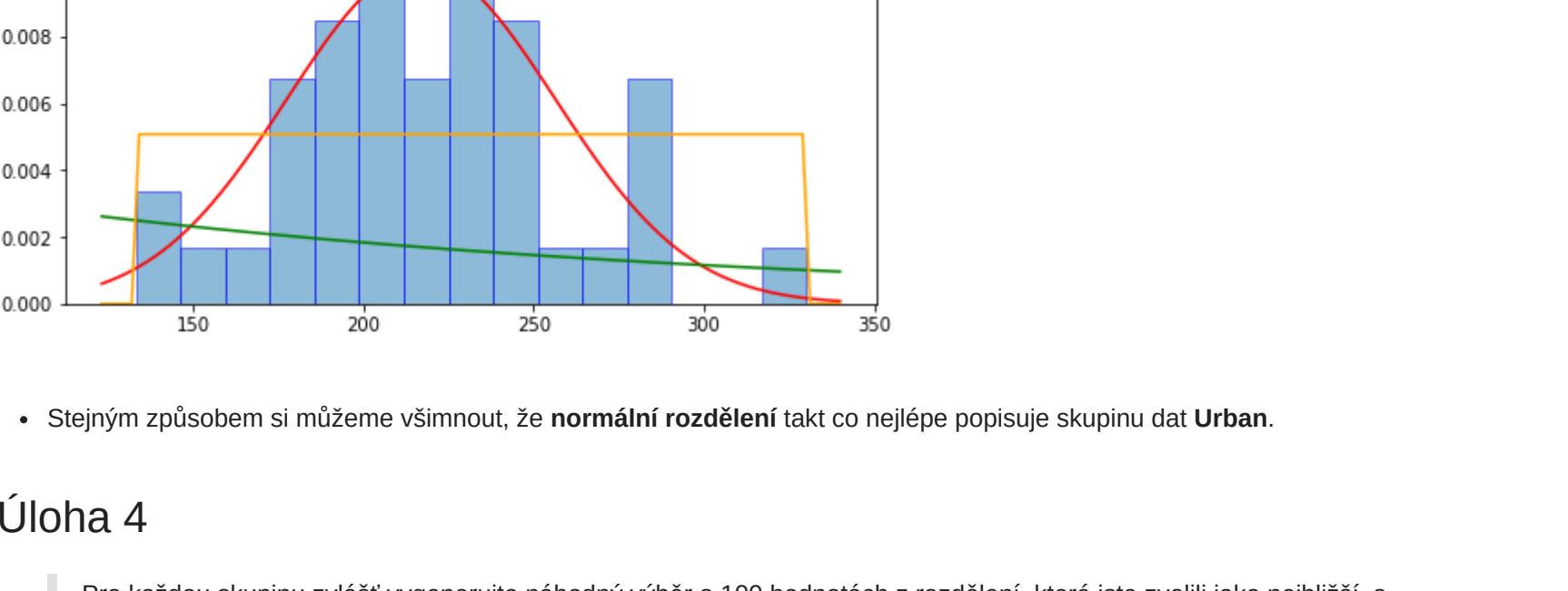
- Zavoláme funkci ExploreDeviations() pro skupinu Rural

```
In [15]: ExploreDeviations(df_Rural, 'Rural')
```

Normální rozdělení:
Maximálně věrohodný odhad střední hodnoty -> 157.0
Maximálně věrohodný odhad rozptylu -> 31.43

Exponenciální rozdělení:
lambda : 0.0064

Rovnoměrné rozdělení:
a : 95
b : 231



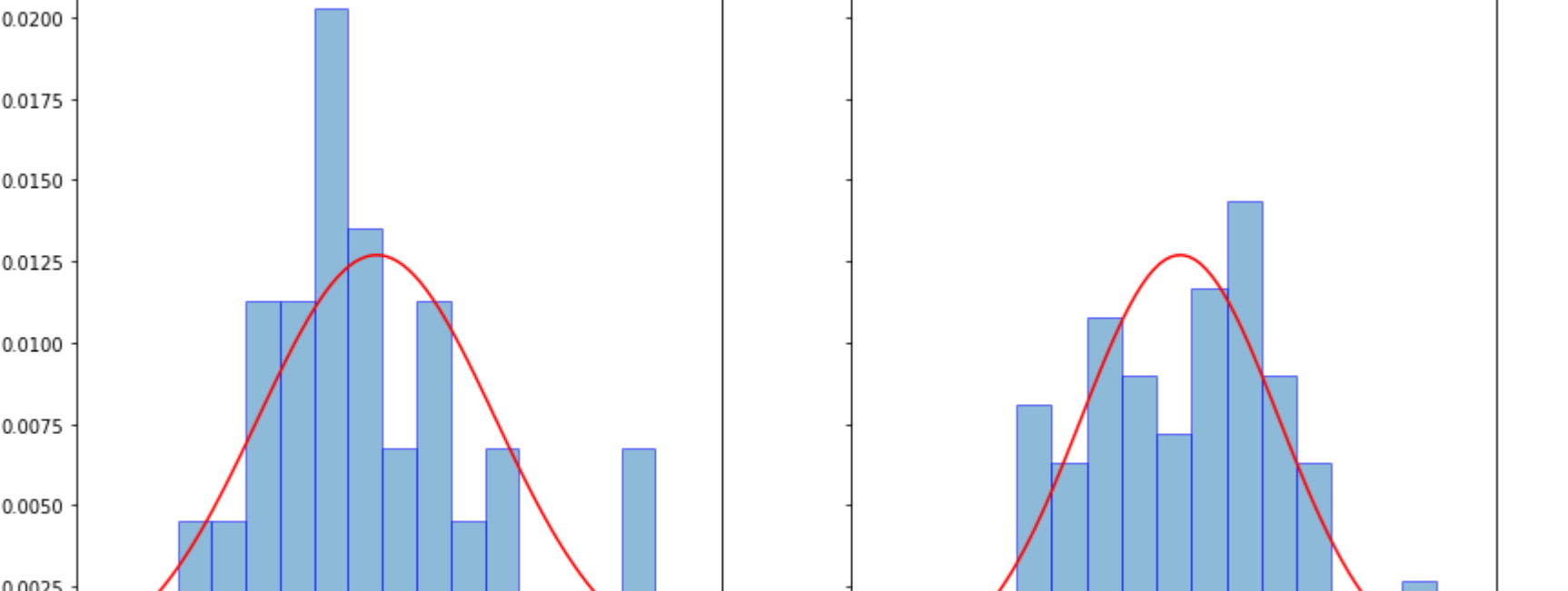
- Z grafu si můžeme všimnout, že normální rozdělení co nejlépe popisuje skupinu dat Rural.
- Ted zavoláme funkci ExploreDeviations() pro skupinu Urban

```
In [16]: ExploreDeviations(df_Urban, 'Urban')
```

Normální rozdělení:
Maximálně věrohodný odhad střední hodnoty -> 216.87
Maximálně věrohodný odhad rozptylu -> 39.47

Exponenciální rozdělení:
lambda : 0.0046

Rovnoměrné rozdělení:
a : 139
b : 330



- Stejným způsobem si můžeme všimnout, že normální rozdělení takt co nejlépe popisuje skupinu dat Urban.

Úloha 4

Pro každou skupinu zvlášť vygenerujte náhodný výběr o 100 hodnotách z rozdělení, které jste zvolili jako nejbližší s parametry odhadnutými v předchozím bodě. Porovnejte histogram simulovaných hodnot s pozorovanými daty.

- Ve funkci Simulation() generujeme náhodný výběr o 100 hodnotách z normálního rozdělení, jehož parametry používáme z předchozího bodu.
- Pro zobrazení simulovaných dat postupujeme stejně jako pro původní data v předchozím bodě.
- Dále kreslíme původní histogram pro normální rozdělení a histogram pro simulovaná data.

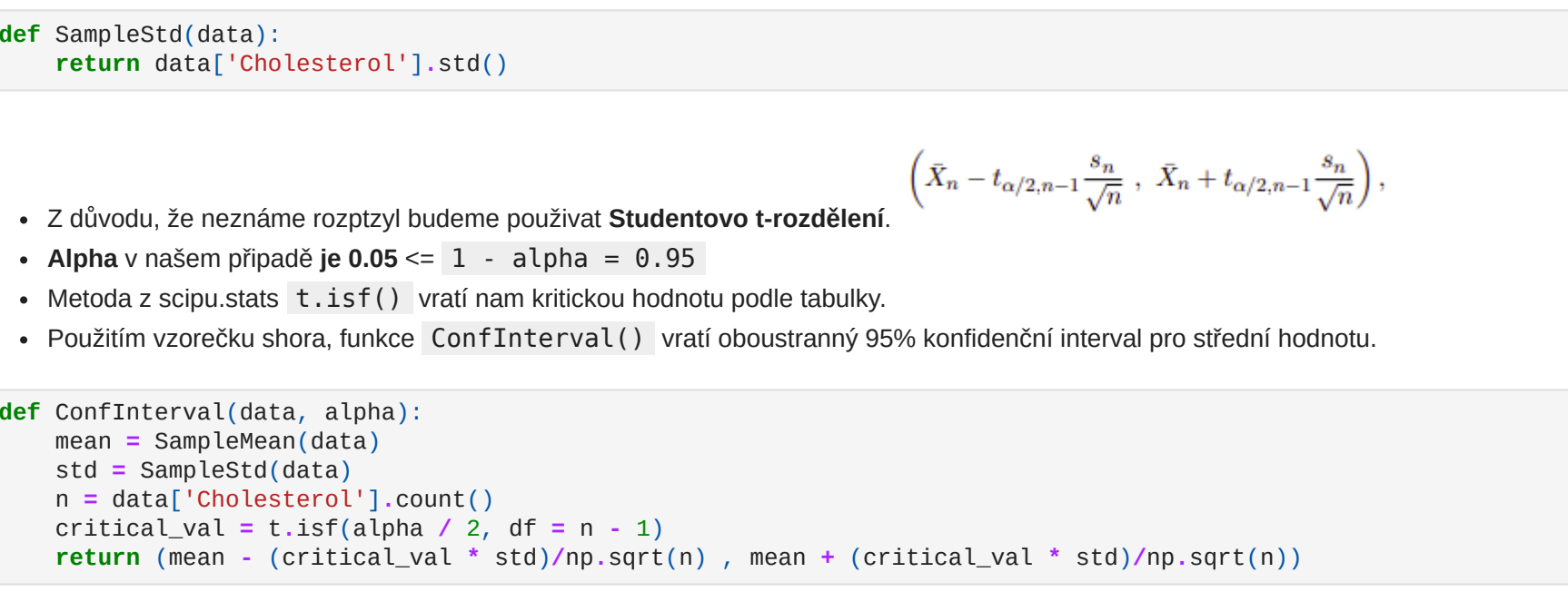
```
In [17]: def Simulation(data, name):
fig, (ax1, ax2) = plt.subplots(1, 2, figsize = (14,7), sharey = True)
print("Nulová hypotéza: K (", K, ") se rovná střední hodnotě (", round(SampleMean(data), 2), ") na hladině významnosti 5%")
sim_data = np.random.normal(mean, std, 100)
ax1.hist(data['Cholesterol'], bins=15, density=True, alpha = 0.5, edgecolor='b')
ax1.set_title("Rozdělení původních dat - " + name)

x = np.linspace(data['Cholesterol'].min()-10, data['Cholesterol'].max()+10, 100)
x_sim = np.linspace(sim_data.min()-10, sim_data.max()+10, 100)

y_norm = norm.pdf(x, mean, std)
y_norm_sim = norm.pdf(x_sim, mean, std)
ax1.plot(x, y_norm, color = 'r')
ax2.set_title("Rozdělení simulovaných dat - " + name)
ax2.hist(sim_data, bins=15, density=True, alpha = 0.5, edgecolor='b')
ax2.plot(x_sim, y_norm_sim, color = 'r')
```

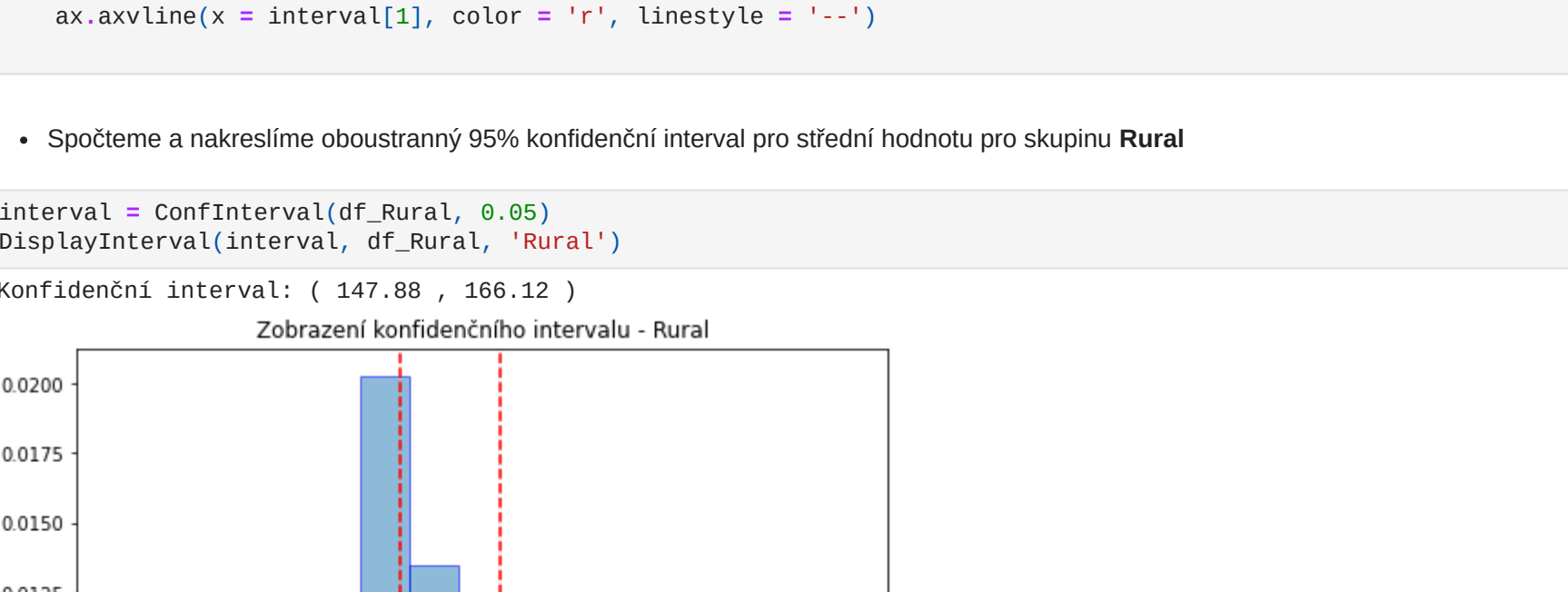
- Zavoláme funkci Simulation() pro skupinu Rural

```
In [18]: Simulation(df_Rural, 'Rural')
```



- Z porovnání histogramů můžeme jistě říct, že data skupiny Rural odpovídají normálnímu rozdělení.
- Ted zavoláme funkci Simulation() pro skupinu Urban

```
In [19]: Simulation(df_Urban, 'Urban')
```



- Z porovnání histogramů můžeme jistě říct, že data skupiny Urban taky odpovídají normálnímu rozdělení.

Úloha 5

Pro každou skupinu zvlášť spočítejte oboustranný 95% konfidenční interval pro střední hodnotu.

- Vytvoříme pomocnou funkci pro výpočet výběrové směrodatné odchylky, využijeme metody pandas std(), která funguje podle vzorečku

$$s_n = \sqrt{s_n^2}$$

```
In [20]: def SampleStd(data):
return data['Cholesterol'].std()
```

- Z důvodu, že neznáme rozptyl budeme používat Studentovo t-rozdělení.
- Alpha v našem případě je 0.05 <= 1 - alpha = 0.95
- Metoda z sccpu.stats t.isf() vrátí nam kritickou hodnotu z tabulky.
- Použitím vzorečku shora, funkce ConfInterval() vrátí oboustranný 95% konfidenční interval pro střední hodnotu.

```
In [21]: def ConfInterval(data, alpha):
mean = SampleMean(data)
std = SampleStd(data)
n = data['Cholesterol'].count()
critical_val = t.isf(alpha / 2, df = n - 1)
return (mean - critical_val * std)/np.sqrt(n), mean + (critical_val * std)/np.sqrt(n))
```

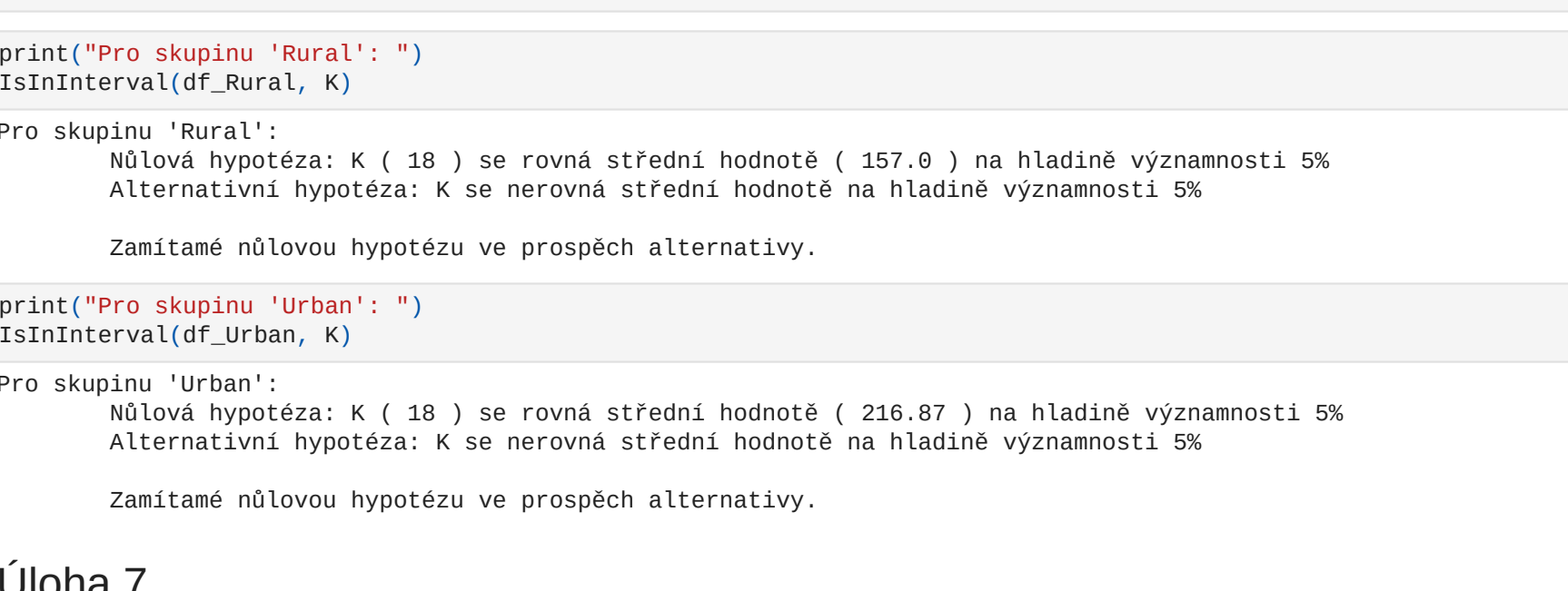
- Tobhle je funkce, která vykreslí konfidenční interval pro střední hodnotu.

```
In [22]: def DisplayInterval(interval, data, name):
print("Konfidenční interval: (", round(interval[0], 2), ",", round(interval[1], 2), ")")

f, ax = plt.subplots(figsize = (8, 6))
ax.hist(data['Cholesterol'], bins=15, density=True, alpha = 0.5, edgecolor='b')
ax.set_title("Zobrazení konfidenčního intervalu - " + name)
ax.axvline(x = interval[0], color = 'r', linestyle = '--')
ax.axvline(x = interval[1], color = 'r', linestyle = '--')
```

- Spočteme a nakreslíme oboustranný 95% konfidenční interval pro střední hodnotu pro skupinu Rural

```
In [23]: interval = ConfInterval(df_Rural, 0.05)
DisplayInterval(interval, df_Rural, 'Rural')
```



- Stejným způsobem spočteme a nakreslíme oboustranný 95% konfidenční interval pro střední hodnotu pro skupinu Urban

```
In [24]: interval = ConfInterval(df_Urban, 0.05)
DisplayInterval(interval, df_Urban, 'Urban')
```



Úloha 6

Pro každou skupinu zvlášť otestujte na hladině významnosti 5 % hypotézu, zda je střední hodnota rovná hodnotě K (parametr úlohy), proti oboustranné alternativě. Můžete použít buď výsledek z předchozího bodu, nebo výstup z příslušné vestavěné funkce vašeho softwaru.

- Pro testování na hladině významnosti 5 % hypotézu (zda je střední hodnota rovná hodnotě K), používáme interval z předchozího bodu.
- Leží-li hodnota K v konfidenčním intervalu, pak platí nulová hypotéza (střední hodnota rovná hodnotě K). V opačném případě zamítneme nulovou hypotézu ve prospěch alternativy.

```
In [25]: def IsInInterval(data, K):
print("Nulová hypotéza: K (", K, ") se rovná střední hodnotě (", round(SampleMean(data), 2), ") na hladině významnosti 5%")
print("Alternativní hypotéza: K se nerovná střední hodnotě na hladině významnosti 5%")
interval = ConfInterval(data, 0.05)
if K in interval[0]:
print("Zamítáme nulovou hypotézu.")
else:
print("Nezamítáme nulovou hypotézu ve prospěch alternativy.")
```

```
In [26]: print("Pro skupinu 'Rural': ")
IsInInterval(df_Rural, K)

Pro skupinu 'Rural':
Nulová hypotéza: K ( 18 ) se rovná střední hodnotě ( 157.0 ) na hladině významnosti 5%
Alternativní hypotéza: K se nerovná střední hodnotě na hladině významnosti 5%
Zamítáme nulovou hypotézu ve prospěch alternativy.
```

```
In [27]: print("Pro skupinu 'Urban': ")
IsInInterval(df_Urban, K)

Pro skupinu 'Urban':
Nulová hypotéza: K ( 18 ) se rovná střední hodnotě ( 216.87 ) na hladině významnosti 5%
Alternativní hypotéza: K se nerovná střední hodnotě na hladině významnosti 5%
Zamítáme nulovou hypotézu ve prospěch alternativy.
```

Úloha 7

Na hladině významnosti 5 % otestujte, jestli mají pozorované skupiny stejnou střední hodnotu. Typ testu a alternativy stanovte tak, aby vaše volba nejlépe korespondovala s povahou zkoumaného problému.

- Pro řešení této úlohy použijeme dvouvýběrový t-test.
- Ale nejprve si musíme prozkoumat rovnost rozptylů dvou skupin.
- Nalézáme naše data pochází z normálního rozdělení, proto můžeme využít F-test pro ověření shodnosti rozptylů.
- Pomocí funkce z knihovny scipy f.cdf() provedeme F-test. Na základě výsledného p-value rozhodneme o shodnosti rozptylů.
- Kolik p-value bude větší než alpha(0.05) zamítneme nulovou hypotézu o shodnosti rozptylů.

```
In [28]: var_Rural = SampleVariance(df_Rural)
var_Urban = SampleVariance(df_Urban)
alpha = 0.05

F = var_Rural/var_Urban
df1 = df_Rural['Cholesterol'].count() - 1
df2 = df_Urban['Cholesterol'].count() - 1

p_value = f.cdf(F, df1, df2)
print("p-value F testu: ", p_value.round(3))
if p_value > alpha:
print("Zamítáme nulovou hypotézu o shodosti rozptylů")
else:
print("Nezamítáme nulovou hypotézu o shodosti rozptylů")
```

- V případě neshodnosti rozptylů, použijeme dvouvýběrový t-test pro případ nerovnosti rozptylů.

H_0	H_A	testová statistika T	kritický obor W_α
$\mu_1 = \mu_2$	$\mu_1 \neq \mu_2$	$T = \frac{\bar{X}_{n_1} - \bar{Y}_{n_2}}{s_d}$	$ T > t_{\alpha/2, n_d}$
$\mu_1 \leq \mu_2$	$\mu_1 > \mu_2$		$T > t_{\alpha, n_d}$
$\mu_1 \geq \mu_2$	$\mu_1 < \mu_2$		$T < -t_{\alpha, n_d}$

- Kde $s_d = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$
- $n_d = \frac{1}{\frac{1}{n_1-1} \left(\frac{s_1^2}{n_1} \right)^2 + \frac{1}{n_2-1} \left(\frac{s_2^2}{n_2} \right)^2}$

```
In [29]: print("Nulová hypotéza: ", SampleMean(df_Rural), " = ", round(SampleMean(df_Urban), 2))

Nulová hypotéza: 157.0 = 216.87
```

```
In [30]: n1 = df_Rural['Cholesterol'].count()
n2 = df_Urban['Cholesterol'].count()

s_d = np.sqrt(var_Rural/n1 + var_Urban/n2)
n_d = (s_d**4)/(1*(n1-1)) + (var_Rural/n1)**2 + (1*(n2-1)) * (var_Urban/n2)**2

t_critical_val = t.isf(alpha/2, df = n_d)
print("Kritická value t = ", round(t_critical_val, 5))

T = (SampleMean(df_Rural) - SampleMean(df_Urban))/s_d
print("T = ", round(abs(T), 5), "\n", round(t_critical_val, 5))

if (abs(T) > t_critical_val):
print("Nulová hypotéza: K se nerovná střední hodnotě na hladině významnosti 5%")
else:
print("Zamítáme nulovou hypotézu o shodnosti průměru dvou testovacích skupin.")
print("Nezamítáme nulovou hypotézu o shodnosti průměru dvou testovacích skupin.")
```

Kritická value t = 1.99861
|T| = 8.09041
8.09041 > 1.99861
Zamítáme nulovou hypotézu o shodnosti průměru dvou testovacích skupin.