

ÚLOHA 1 – Bayesovské odhady – 4. body

a) Konjugované apriorní a aposteriorní rozdělení, prediktivní rozdělení [2 body] Předpokládáme, že počet připojení na internetovou síť za 1 ms je popsán náhodnou veličinou s Poissonovým rozdělením s parametrem λ , tj. $X \sim \text{Poi}(\lambda)$. O parametru λ máme následující expertní odhad: každých 5 ms by mělo nastat 10 připojení. Pozorovali jsme připojení po dobu 100 ms. Pozorování a počtu připojení za každou 1 ms jsou uvedené v souboru `measurements.csv` ve sloupci „`uloha_1_a`“. Vaším zadáním je z této expertní informace určit konjugované apriorní rozdělení a parametru Poissonova rozdělení a na základě pozorování určit aposteriorní rozdělení. Dále určete apriorní a aposteriorní prediktivní rozdělení parametru λ . Požadovaný výstup: porovnejte je.

- Do jednoho obrázku vykreslte apriorní a aposteriorní hustotu parametru Poissonova rozdělení λ .
- Do jednoho obrázku vykreslte apriorní a aposteriorní prediktivní hustotu pozorování x za jeden časový interval.
- Sestavte 95% interval spolehlivosti pro parametru λ z apriorního a aposteriorního rozdělení a porovnejte je.
- Vyberte si dva aposteriorní bodové odhady parametru λ , porovnejte je a ekomentujte jejich výběr.
- Vyberte si jeden apriorní a jeden aposteriorní bodový odhad počtu pozorování a

```
In [145.]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import norm, truncnorm, gamma, poisson, nbinom, pearsonr

data = pd.read_excel("Projekt-2_Data.xlsx", engine='openpyxl')

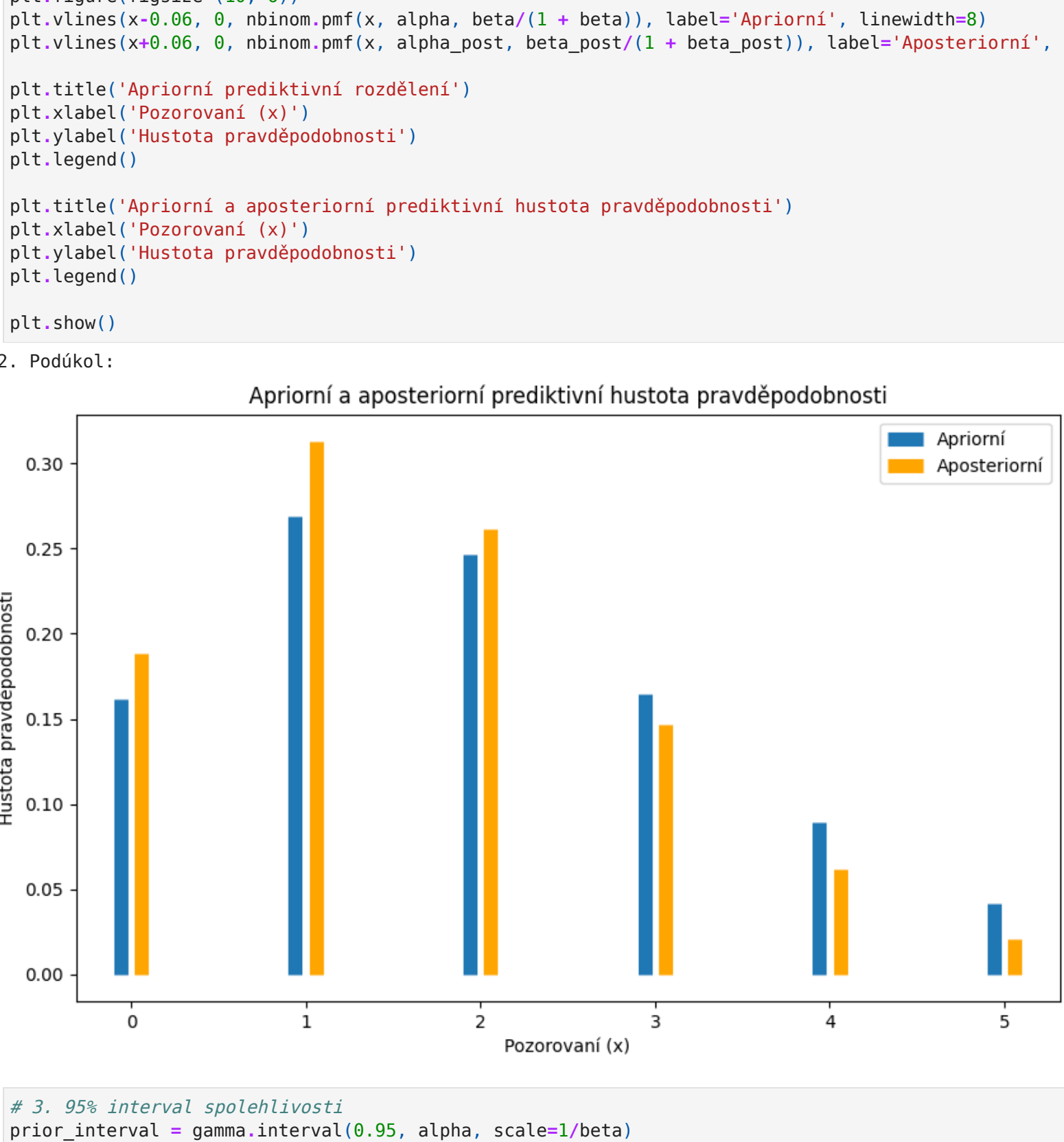
df = pd.DataFrame(data, columns=['uloha_1_a'])[:180]

# Apriorní rozdělení
alpha = 10
beta = 5

# Apriorní rozdělení
alpha_post = alpha + df.sum()
beta_post = beta + len(df)

# 1. Vykreslení apriorního a aposteriorního rozdělení  $\lambda$ 
print("1. Podúkol:")
x = np.linspace(0, 10, 1000)
prior_distribution = gamma.pdf(x, alpha, scale=1/beta)
posterior_distribution = gamma.pdf(x, alpha_post, scale=1/beta_post)
plt.figure(figsize=(10, 6))
plt.plot(x, prior_distribution, label='Apriorní rozdělení', color='blue')
plt.plot(x, posterior_distribution, label='Aposteriorní rozdělení', color='red')
plt.title('Apriorní a Aposteriorní rozdělení parametru  $\lambda$ ')
plt.xlabel('lambda')
plt.ylabel('Hustota pravděpodobnosti')
plt.legend()
plt.show()
```

1. Podúkol:



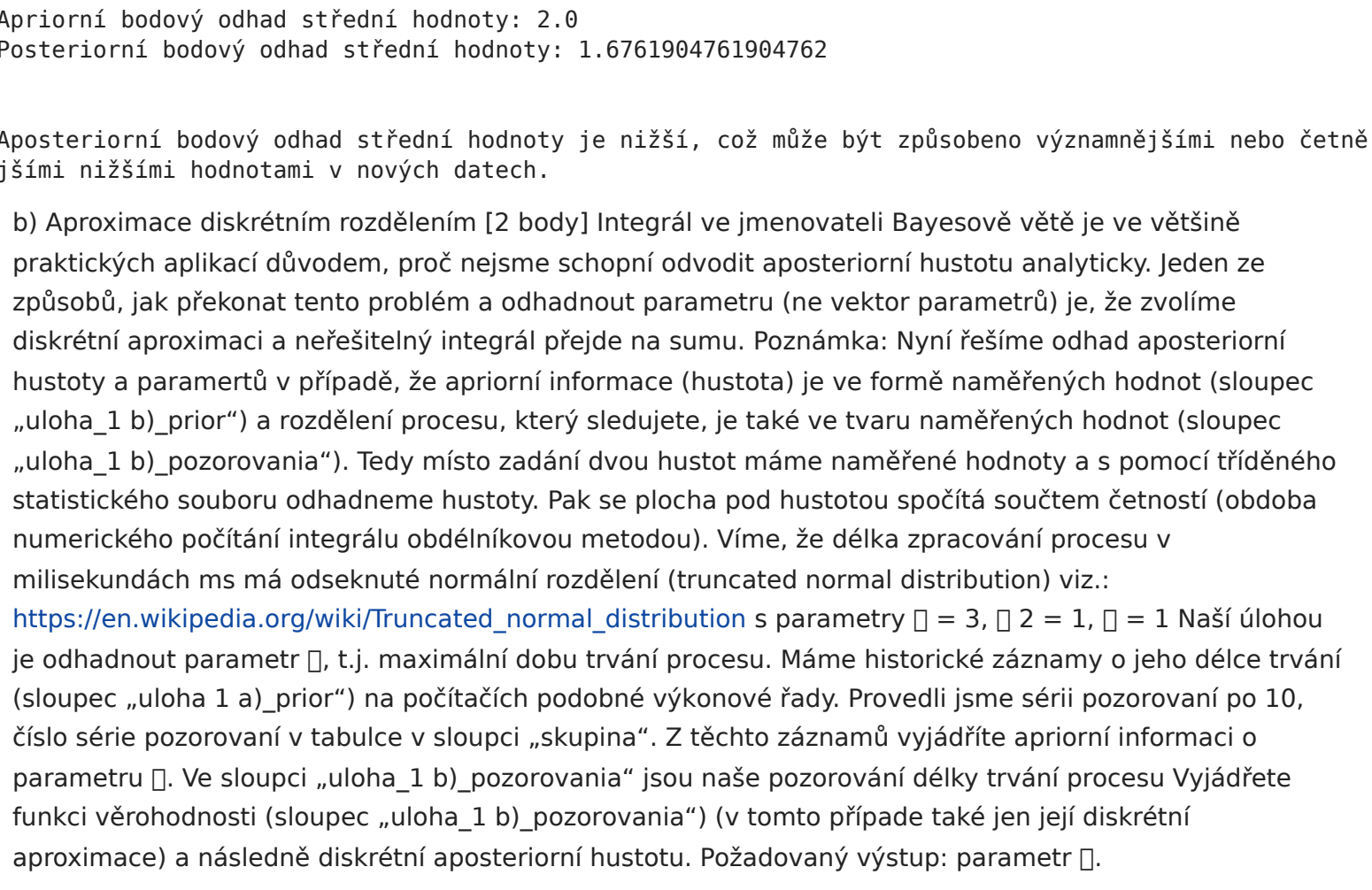
```
In [145.]: # 2. Výpočet a vykreslení apriorního a aposteriorního prediktivního rozdělení
print("2. Podúkol:")
pred_prior = poisson.pmf(x, alpha/beta)
pred_post = poisson.pmf(x, alpha_post/beta_post)

x = np.arange(0, 6)
plt.figure(figsize=(10, 6))
plt.vlines(x=[0.06, 0, nbinom.pmf(x, alpha, beta/(1 + beta))], label='Apriorní', linewidth=8)
plt.vlines(x=[0.06, 0, nbinom.pmf(x, alpha_post, beta_post/(1 + beta_post))], label='Aposteriorní', color='red', linewidth=8)

plt.title('Apriorní prediktivní rozdělení')
plt.xlabel('Pozorování (x)')
plt.ylabel('Hustota pravděpodobnosti')
plt.legend()

plt.title('Apriorní a aposteriorní prediktivní hustota pravděpodobnosti')
plt.xlabel('Pozorování (x)')
plt.ylabel('Hustota pravděpodobnosti')
plt.legend()
plt.show()
```

2. Podúkol:



```
In [145.]: # 3. 95% interval spolehlivosti
prior_interval = gamma.interval(0.95, alpha, scale=1/beta)
posterior_interval = gamma.interval(0.95, alpha_post, scale=1/beta_post)

print("3. Podúkol:")
print(f"Apriorní 95% interval spolehlivosti: {prior_interval}")
print(f"Aposteriorní 95% interval spolehlivosti: {posterior_interval}")
print("Aposteriorní interval je užší než apriorní interval, což naznačuje, že nová data vedou k větš")
print("\n")

# 4. Dva aposteriorní bodové odhady parametru lambda
posterior_mean = alpha_post / beta_post
posterior_median = gamma.median(alpha_post, scale=1/beta_post)

print("4. Podúkol:")
print(f"Aposteriorní bodový odhad lambda pro střední hodnoty: {posterior_mean.iloc[0]}")
print(f"Aposteriorní bodový odhad lambda pro střední hodnoty: {posterior_mean.iloc[0]}")
print("\n")
print("Tyto hodnoty jsem vybral, protože pokud je obě známe, tak nám mohou dát jistý obrázek o rozl")
print("To, že jsou hodnoty mediánu a střední hodnoty velice podobné může naznačovat, že data mají symetrick")

# 5. Jeden apriorní a jeden aposteriorní bodový odhad počtu pozorování
prior_mean = alpha / beta
posterior_mean = alpha_post.iloc[0]/beta_post

print("5. Podúkol:")
print(f"Apriorní bodový odhad střední hodnoty: {prior_mean}")
print(f"Aposteriorní bodový odhad střední hodnoty: {posterior_mean}")
print("\n")
print("Aposteriorní bodový odhad střední hodnoty je nižší, což může být způsobeno významnějšími ne")

# 3. Podúkol:
Apriorní 95% interval spolehlivosti: (0.959877392264868, 3.4169669828383)
Aposteriorní 95% interval spolehlivosti: (1.4376938284869922, 1.9327287471868797)
```

Aposteriorní interval je užší než apriorní interval, což naznačuje, že nová data vedou k větší jistě v odhadu parametru λ .

4. Podúkol:
Aposteriorní bodový odhad λ pro střední hodnoty: 1.6761994761994762
Aposteriorní bodový odhad λ mediánu: 1.6736169441241727

Tyto hodnoty jsem vybral, protože pokud je obě známe, tak nám mohou dát jistý obrázek o rozložení d at, což dále komentuju.
To, že jsou hodnoty mediánu a střední hodnoty velice podobné může naznačovat, že data mají symetrick tvar a nejsou výrazně zkreslená do jednoho nebo druhého směru.

- Apriorní bodový odhad střední hodnoty: 2.8
Posteriorní bodový odhad střední hodnoty: 1.6761994761994762

Aposteriorní bodový odhad střední hodnoty je nižší, což může být způsobeno významnějšími nebo četnějšími nižšími hodnotami v nových datech.

b) Aproximace diskrétního rozdělení [2 body] Integrál ve jmenovateli Bayesové věty je ve většině praktických aplikací důvodem, proč nejsme schopni odvodit aposteriorní hustotu analyticky. Jeden ze způsobů, jak překonat tento problém a odhadnout parametru (ne vektor parametru) je, že zvolíme diskrétní aproximaci a neřešíme integrál přejde na sumu. Poznámka: Nyní řešíme odhad aposteriorní hustoty a parametru v případě, že apriorní informace (hustota) je ve formě naměřených hodnot (sloupec „`uloha_1_b_prior`“) a rozdělení procesu, který sledujeme, je také ve tvaru naměřených hodnot (sloupec „`uloha_1_b_pozorovani`“). Tedy místo zadání dvou hustot máme naměřené hodnoty a s pomocí třídného statistického souboru odhadneme hodnoty. Pak se plocha pod hustotou spočítá součtem četností (obdobna numerického počítání integrálu obdélníkovou metodou). Víme, že délka zpracování procesu v milisekundách má omezenou normální rozdělení (truncated normal distribution) viz: https://en.wikipedia.org/wiki/Truncated_normal_distribution s parametrem $\mu = 3$, $\sigma^2 = 1$. Naši úlohu je odhadnout parametru λ , tj. maximální dobu trvání procesu. Máme historické záznamy o jeho délce trvání (sloupec „`uloha_1_a_prior`“) na počítacích podobné výkonové řady. Provedli jsme sérii pozorování po 10, číslo série pozorování v tabulce v sloupci „skupina“. Z těchto záznamů vyjádříte apriorní informaci o parametru λ . Ve sloupci „`uloha_1_b_pozorovani`“ jsou naše pozorování délky trvání procesu. Vyjádříte funkci věrohodnosti (sloupec „`uloha_1_b_pozorovani`“) (v tomto případě také jen její diskrétní aproximace) a následně diskrétní aposteriorní hustotu. Požadovaný výstup: parametru λ .

- Do jednoho grafu vykreslte apriorní, aposteriorní hustotu a funkci věrohodnosti. Funkci věrohodnosti normujte tak, aby její součet byl 1 kvůli porovnatelnosti v obrázku.
- Z aposteriorní hustoty určete 95% interval spolehlivosti (konfidenční interval) pro
- Vyberte dva bodové odhady parametru λ a spočítejte je.

```
In [145.]: ##### APRIORNI FUNKCE #####

df_prior = pd.DataFrame(data, columns=['uloha_1_b_prior', 'skupina'])
grouped_data_prior = df_prior.groupby('skupina')['uloha_1_b_prior'].apply(list).reset_index(name='prior_values')

# Funkce pro nalezení maximální hodnoty v listu
def find_max_in_list(list):
    return max(list)

# Nalezení maximální hodnoty v každé skupině
max_prior_values = grouped_data_prior['prior_values'].apply(find_max_in_list)

# Počet intervalů, na které se rozdělí data
num_intervals = 50

hist, bin_edges = np.histogram(max_prior_values, bins=num_intervals, density=True)

prior_values = hist
bin_centers = (bin_edges[:-1] + bin_edges[1:]) / 2
plt.figure(figsize=(10, 6))
plt.bar(bin_centers, prior_values, width=(bin_edges[1] - bin_edges[0]), align='center', edgecolor='black')

##### VEROHODNOSTNI FUNKCE #####
df_observ = pd.DataFrame(data, columns=['uloha_1_b_pozorování'])[:100]
mu=3
sigma = 1
a=1
a_normalized = (a - mu) / sigma

likelihood_values = [truncnorm.pdf(df_observ, a=a_normalized, b=(b - mu) / sigma, loc=mu, scale=sigma)
likelihood_values = [np.prod(x) for x in likelihood_values]
normalization_factor = np.trapz(likelihood_values, bin_centers)
likelihood_values = likelihood_values / normalization_factor

plt.bar(bin_centers, likelihood_values, width=(bin_edges[1] - bin_edges[0]), align='center', edgecolor='black')

##### APOSTERIORNI FUNKCE #####

# Výpočet normalizačního faktoru
normalization_factor = np.trapz([likelihood_values[x]*prior_values[x] for x in range(len(likelihood_values))])

# Výpočet aposteriorní funkce
aposterior_values = [(likelihood_values[x]*prior_values[x]/normalization_factor) for x in range(len(likelihood_values))]

plt.bar(bin_centers, aposterior_values, width=(bin_edges[1] - bin_edges[0]), align='center', edgecolor='black')

print("1. Podúkol:")
plt.xlabel('Hustota pravděpodobnosti')
plt.ylabel('Hustota pravděpodobnosti')
plt.title('Apriorní a aposteriorní funkce hustoty a věrohodnostní funkce')
plt.legend()
plt.show()

# Kumulativní distribuční funkce (CDF)
cdf = np.cumsum(aposterior_values) / np.sum(aposterior_values)
# Určení hranice intervalu pomocí CDF
lower_bound = bin_centers[np.argmax(cdf >= 0.025)]
upper_bound = bin_centers[np.argmax(cdf >= 0.975)]

print("2. Podúkol:")
print("95% intervalový odhad pro parametru b: ({}, {})", format(lower_bound, upper_bound))

likelihood_val = np.sum(bin_centers * aposterior_values) / np.sum(aposterior_values)
median_val = bin_centers[np.argmax(cdf >= 0.5)]

print("\n")
print("Bodové odhad střední hodnoty: ", median_val)
print("Bodové odhad mediánu: ", median_val)
```

1. Podúkol:



- Podúkol:
95% intervalový odhad pro parametru b: (5.693712028182375, 7.008910628347767)
- Podúkol:
Bodové odhad střední hodnoty: 6.052771319832352
Bodové odhad mediánu: 5.956751748215453

ÚLOHA 2 – Regrese – 8. bodů

Úkoly a požadovaný výstup: považujte nyní kvadratický model (všechny interakce druhého řádu a všechny druhé

- Pomocí zpětné eliminace určete vhodný regresní model. Za výchozí „plný“ model

mocniny, které dávají smysl). Zapíše rovnici Vašeho finálního modelu. Diskutujte splnění předpokladů lineární regrese a základní regresní diagnostiky. Pokud (až během regresního modelování) identifikujete některé „extrémně odlehle hodnoty“ můžete ty „nejodlehlejší“ hodnoty, po alespoň krátkém zdůvodnění, vyřadit. (4. bodů) 2) Pomocí Vašeho výsledného modelu identifikujte, pro které nastavení parametrů má odrazem neproblematictější hodnotu. (1. bod) 3) Odhadněte hodnotu odezvy uživatele s Windows při průměrném nastavení ostatních parametrů a předtiskněte konfidenční interval a predikční interval pro toto nastavení. (1. bod) 4) Na základě jakýchkoli vypočtených charakteristik argumentujte, zdali je Váš model „vhodný“ pro další použití. (2. body)

```
In [146.]: import statsmodels.api as sm
from statsmodels.stats.outliers_influence import variance_inflation_factor

df = pd.read_excel("Projekt-2_Data.xlsx", sheet_name="Uloha 2")

# Odhadení závislosti InteractingPct a ScrollingPct pomocí korelační matice ###
X = df.iloc[:, 1:]

### Prevedení kategoriálních dat do jednotlivých sloupců ###
mat=pd.get_dummies(df.drop_first=True)
mat=mat.astype(float)

X = sm.add_constant(mat[['ActiveUsers', 'InteractingPct', 'OSType_MacOS', 'OSType_Windows', 'OSType_iOS']])

In [146.]: ### VYTVOŘENÍ ÚPLNEHO KVADRATICKÉHO MODELU ###
X['ActiveUsers*ActiveUsers'] = X['ActiveUsers'] * X['ActiveUsers']
X['InteractingPct*InteractingPct'] = X['InteractingPct'] * X['InteractingPct']
X['ActiveUsers*InteractingPct'] = X['ActiveUsers'] * X['InteractingPct']
X['ActiveUsers*OSType_MacOS'] = X['ActiveUsers'] * X['OSType_MacOS']
X['ActiveUsers*OSType_Windows'] = X['ActiveUsers'] * X['OSType_Windows']
X['ActiveUsers*OSType_iOS'] = X['ActiveUsers'] * X['OSType_iOS']
X['InteractingPct*OSType_MacOS'] = X['InteractingPct'] * X['OSType_MacOS']
X['InteractingPct*OSType_Windows'] = X['InteractingPct'] * X['OSType_Windows']
X['InteractingPct*OSType_iOS'] = X['InteractingPct'] * X['OSType_iOS']

y = df['Ping [ms]'].astype(int)

orig_X = X.copy()

# Normalizace hodnot
X.iloc[:, 1:] = (X.iloc[:, 1:] - X.iloc[:, 1:].min()) / (X.iloc[:, 1:].max() - X.iloc[:, 1:].min())

In [146.]: model = sm.OLS(y, X).fit()

# Metoda zdlší eliminace
for col, pval in zip(X.columns, model.pvalues):
    if col == 'const':
        continue
    else:
        if pval > 0.05:
            X = X.drop(col, axis=1)
            model = sm.OLS(y, X).fit()

In [146.]: ### VIF (Multikolinearita) ###
vif = pd.Series([variance_inflation_factor(X.values, i) for i in range(X.shape[1])], index=X.columns)
vif_df = vif.to_frame()

# Nastavení názvu sloupce
vif_df.columns = ['VIF']

print(vif_df)
print("\n")
X = X.drop('ActiveUsers*ActiveUsers', axis=1)
vif = pd.Series([variance_inflation_factor(X.values, i) for i in range(X.shape[1])], index=X.columns)
vif_df = vif.to_frame()
vif_df.columns = ['VIF']

print("Hodnoty VIF po odstranění hodnoty ActiveUsers*ActiveUsers:")
print(vif_df)

const          32.761712
ActiveUsers    24.981312
InteractingPct  5.567126
OSType_Windows 5.551936
ActiveUsers*ActiveUsers  22.204978
ActiveUsers*InteractingPct  8.566088
ActiveUsers*OSType_MacOS  1.656643
ActiveUsers*OSType_Windows  6.474746
ActiveUsers*OSType_iOS      1.559683

Hodnoty VIF po odstranění hodnoty ActiveUsers*ActiveUsers:
const          21.405888
ActiveUsers    4.567019
InteractingPct  5.567126
OSType_Windows 5.546736
ActiveUsers*InteractingPct  8.558677
ActiveUsers*OSType_MacOS   1.636529
ActiveUsers*OSType_Windows  6.464752
ActiveUsers*OSType_iOS     1.559671
```

```
In [146.]: results = model
residuals = results.resid

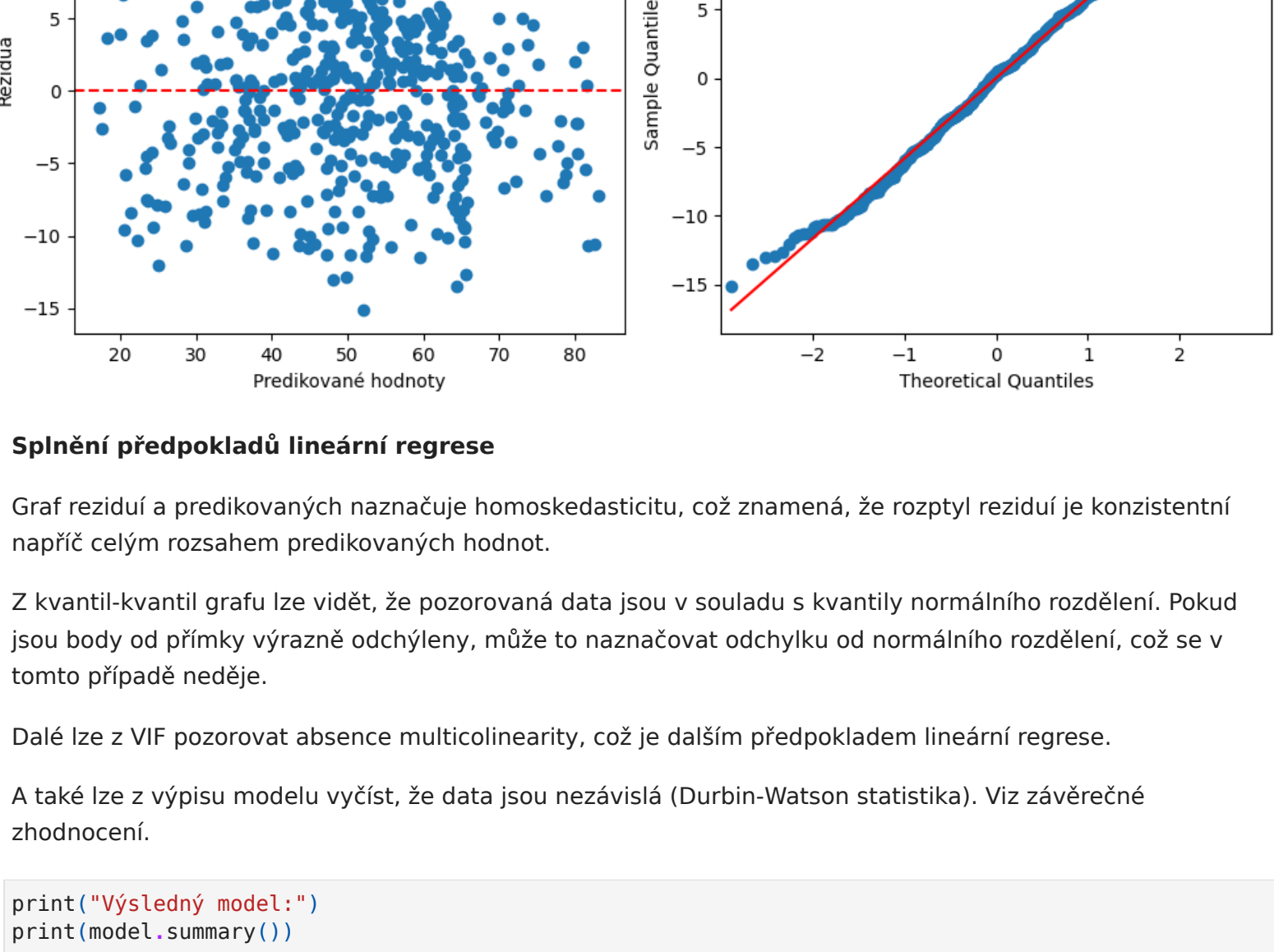
def plot_residuals(results):
    fig, axs = plt.subplots(1, 2, figsize=(10, 5))

    # Residua a predikované hodnoty
    axs[0].scatter(results.fittedvalues, results.resid)
    axs[0].set_xlabel('Predikované hodnoty')
    axs[0].set_title('Residua a predikované hodnoty')
    axs[0].set_ylabel('Residua')
    axs[0].axhline(y=0, color='red', linestyle='--')

    # Porovnání kvantilů pozorovaných dat s kvantily teoretického normálního rozdělení
    sm.qqplot(results.resid, line='s', ax=axs[1])
    axs[1].set_title('Kvantil-Kvantil graf')

    plt.tight_layout()
    plt.show()

plot_residuals(results)
```



Z grafů jsou patrné dvě odlehle hodnoty, které tedy vymažu, jelikož by přiřlí (negativně) ovlivnily výsledný model.

```
In [146.]: # Denormalizace (navrácení původních hodnot)
data_denorm = orig_X
results = model

# Vlny odlehlelých hodnot (viz democení)
influence = results.get_influence()
leverage = influence.hat_matrix_diag()
cooks_d = influence.cooks_distance()
standardized_residuals = influence.resid_studentized_internal
studentized_residuals = influence.resid_studentized_external

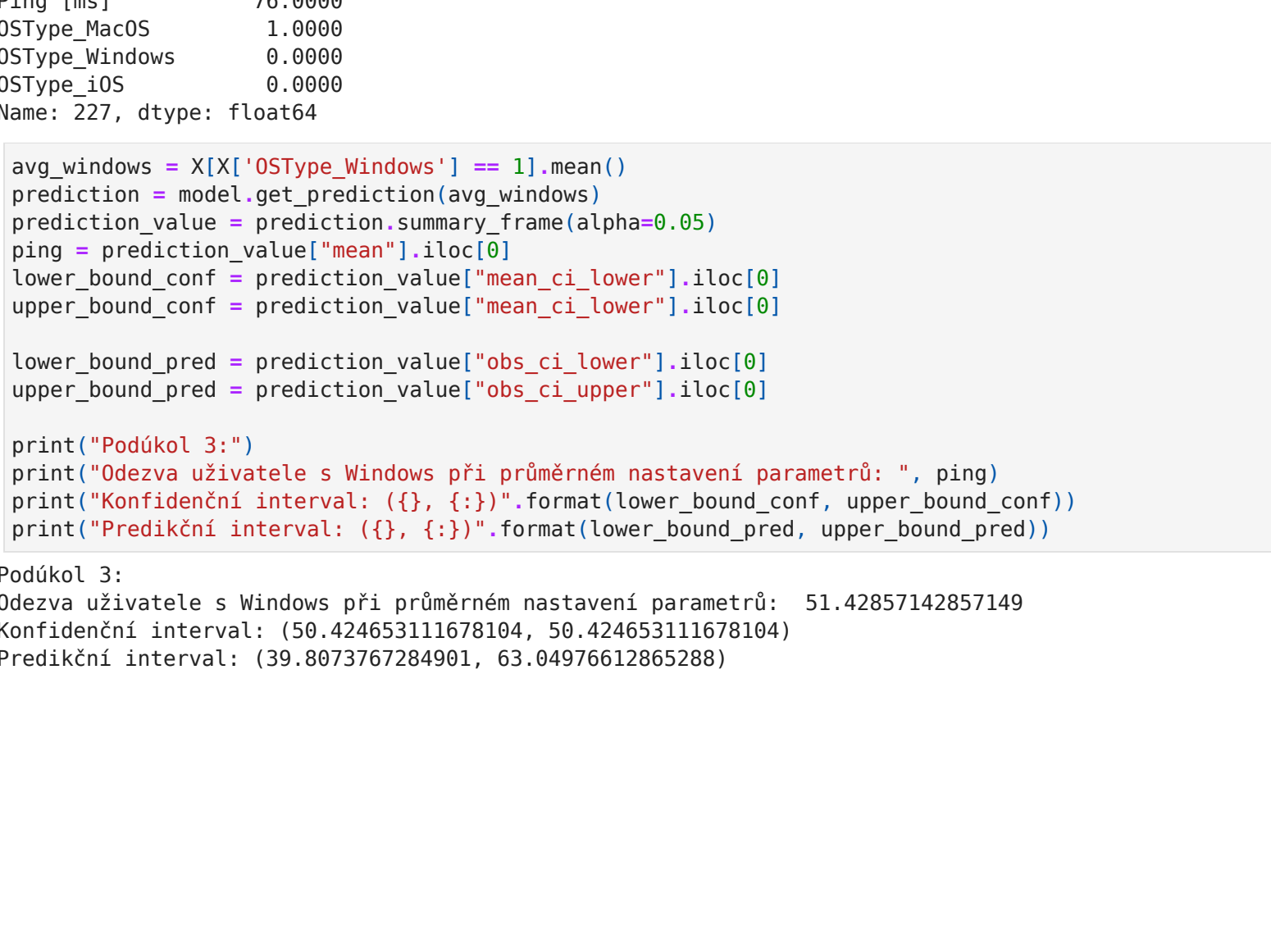
outl_stats_df = pd.DataFrame({
    'Leverage': leverage,
    'Standardized Residuals': influence.resid_studentized_internal,
    'Studentized Residuals': influence.resid_studentized_external,
    'Cook's Distance': influence.cooks_distance(),
    'Cook's Distance p-value': influence.cooks_distance(),
    'index': data_denorm.index
})
(influence.hat_matrix_diag > 3 * len(results.params) / data_denorm.shape[0]) |
(np.abs(influence.resid_studentized_internal) > 2) |
(influence.cooks_distance[1] < 0.05)

print(outl_stats_df)

# Odstranění dvou odlehlelých hodnot
y_drop(index=[255, 476], inplace=True)
X_drop(index=[255, 476], inplace=True)

In [146.]: model = sm.OLS(y, X).fit()

# Zobrazení upravených reziduí
plot_residuals(model)
```



Splnění předpokladů lineární regrese

Graf reziduí a predikovaných naznačuje homoskedasticitu, což znamená, že rozptyl reziduí je konzistentní napříč celým rozsahem predikovaných hodnot.

Z kvantil-kvantil grafu lze vidět, že pozorovaná data jsou v souladu s kvantily normálního rozdělení. Pokud jsou body od přímky výrazně odchyleny, může to naznačovat odchylku od normálního rozdělení, což se v tomto případě neděje.

Další lze z VIF pozorovat absence multicolinearity, což je dalším předpokladem lineární regrese.

A také lze z výpisu modelu vyčíst, že data jsou nezávislá (Durbin-Watson statistika). Viz závěrečné zhodnocení.

```
In [146.]: print("Výsledný model:")
print(model.summary())

Výsledný model:

OLS Regression Results
=====
Dep. Variable:          Ping [ms]          R-squared:                0.842
Model:                  OLS              Adjusted R-squared:        0.874
Method:                 Least Squares     F-statistic:              374.2
Date:                   Tue, 19 Dec 2023   Prob(F-statistic):        1.71e-192
Time:                   22:50:54          Log-Likelihood:           -1592.3
No. Observations:       509              AIC:                     3261.1
Df Residuals:           492              BIC:                     3234.
Df Model:               7
Covariance Type:        nonrobust

=====
coef    std err          t      P>|t|      [0.025     0.975]
-----
const                9.4638         1.230     7.693     0.000     7.047    11.881
ActiveUsers          55.5677         2.180    25.489     0.000    51.284    59.851
InteractingPct       2.1826         2.183     1.000     0.317     -2.183     2.183
OSType_Windows       8.1826         1.494     5.477     0.000     5.424    10.941
ActiveUsers*InteractingPct -32.7889         3.329   -9.825     0.000   -39.249   -26.168
ActiveUsers*OSType_MacOS  16.8209         1.212    13.874     0.000    14.439    19.283
ActiveUsers*OSType_Windows  6.4647         2.390     2.703     0.008     1.708     9.621
ActiveUsers*OSType_iOS -9.9763         1.234   -8.087     0.000   -12.400   -7.552
=====
Omnibus:            4.707      Durbin-Watson:           1.914
Prob(Omnibus):      0.095      Jarque-Bera (JB):           3.366
Skew:               0.020      Prob(JB):                  0.186
Kurtosis:           2.600      Cond. No.                  22.4
=====
```

Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

```
In [146.]: pings = model.predict()
max_ping_index = model.predict().argmax()

# Zjištění nejvyšší odezvy z predikce
print("Maximální hodnota odezvy:", pings[max_ping_index])
# Výpis parametrů s nejvyšší odezvou
print("\n")
print("Problémové nastavení nastavení parametrů:")
print(mat.iloc[max_ping_index])

Podúkol 2:
Maximální hodnota odezvy: 83.1894558968771

Problémové nastavení nastavení parametrů:
ActiveUsers      9953.0800
InteractingPct   0.6729
ScrollingPct     0.3271
Ping [ms]        1.0000
OSType_MacOS     0.0000
OSType_Windows   0.0000
OSType_iOS       0.0000
Name: 227, dtype: float64
```

```
In [147.]: avg_windows = X[X['OSType_Windows'] == 1].mean()
prediction = model.get_prediction(avg_windows)
prediction_value = prediction.summary_frame(alpha=0.05)
ping = prediction_value['mean'].iloc[0]
lower_bound_conf = prediction_value['mean_ci_lower'].iloc[0]
upper_bound_conf = prediction_value['mean_ci_upper'].iloc[0]

lower_bound_pred = prediction_value['obs_ci_lower'].iloc[0]
upper_bound_pred = prediction_value['obs_ci_upper'].iloc[0]

print("Podúkol 3:")
print("Odezva uživatele s Windows při průměrném nastavení parametrů: ", ping)
print("Konfidenční interval: ({}, {})", format(lower_bound_conf, upper_bound_conf))
print("Predikční interval: ({}, {})", format(lower_bound_pred, upper_bound_pred))

Podúkol 3:
Odezva uživatele s Windows při průměrném nastavení parametrů: 51.42857142857149
Konfidenční interval: (50.424653111678194, 58.424653111678194)
Predikční interval: (39.04973767284901, 63.04973767284901)
```


Podúkol 4:

Použitelnost, správnost mého modelu indikuje spousta ukazatelů, které můžeme vyčíst ze shrnutí modelu:

- Durbin-Watson statistika (1.914), která se blíží hodnotě 2, což značí, že hodnoty reziduí nejsou na sobě nijak závislé.
- Cond. No. (22.4) je nízké tedy matice plánu je dobře podmíněná a malá změna ve vstupech způsobí malou (úměrnou) změnu v koeficientech.
- Prob(JB) (0.186) je vyšší než hladina významnosti (0.05) a tedy nezamítáme nulovou hypotézu o normálním rozdělení reziduí.
- R-squared(0.842) mřira variability, která se blíží hodnotě 1, indikuje, že model dobře vysvětluje variabilitu závislé proměnné na základě vstupních dat.