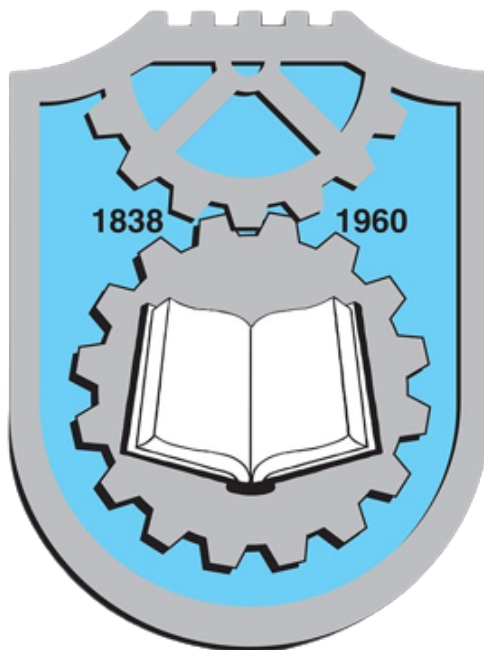


Универзитет у Крагујевцу
Факултет инжењерских наука



Неуронске мреже

Документација

Четврти домаћи задатак
Предвиђање цена акција на берзи

Професор:

Проф. др Весна Ранковић

Студент:

Каришић Ђорђе 393/2023

___ . ___ . 2023.

Садржај

1	Увод	2
2	Анализа скупа података	3
3	Преглед софтверског решења	8
3.1	Визуелизација	8
3.2	Претпроцесирање	10
3.3	Креирање, тренирање и тестирање модела мреже	11
4	Резултати	14

1 Увод

Предвиђање цена акција на тржишту акција је изазован задатак, како због сложености података тако и због сложености самих зависности између њих. Подаци, често представљени графичким форматом, значајно зависе временски једни од других.

Цена у тренутку t може зависити од цене у тренутку $t - 1$, али и $t - 2$ и тако даље до неког крајњег $t - \alpha$ тренутка, али може бити и независна, при чему се претпоставља јак екстерни утицај на тржиште.

Како је потребно моделирати сложене временске зависности, предлаже се употреба дубоког учега, конкретно модела неуронских мрежа – мрежа са дуготрајном краткорочном меморијом.

Овај тип мрежа, уз помоћ своје специфичне и једноставне структуре може моделирати веома дугорочне зависности у подацима и тиме адекватно апроксимирати функцију рачунања излаза на основу улаза чије су вредности временске серије. Једноставније речено, на основу цене претходних $N - 1$ дана, овај тип модела може прецизно и брзо одредити цену у N дану.

Такав модел на улазу треба имати адекватне податке, уређене, скалиране и “чисте” од било каквих аномалија и нежељених вредности – потребно их је адекватно процесирати.

Како би се пружио увид о сложености, структури и везама међу подацима из скупа података потребно је, такође, визуелизовати скуп података, специфично за сврху коју ће обављати. Како се у овом случају ради о подацима који представљају цену акције у датом временском тренутку на тржишту акција, потребно их је испитати у таквом окружењу.

Скуп података употребљен у оквиру овог рада је CAC40 Dataset (preprocessed), скуп временски зависних података сакупљених са Француског тржишта акцијама за 40 најрелевантнијих компанија – од којих је за овај рад одабрана Peugeot.

2 Анализа скупа података

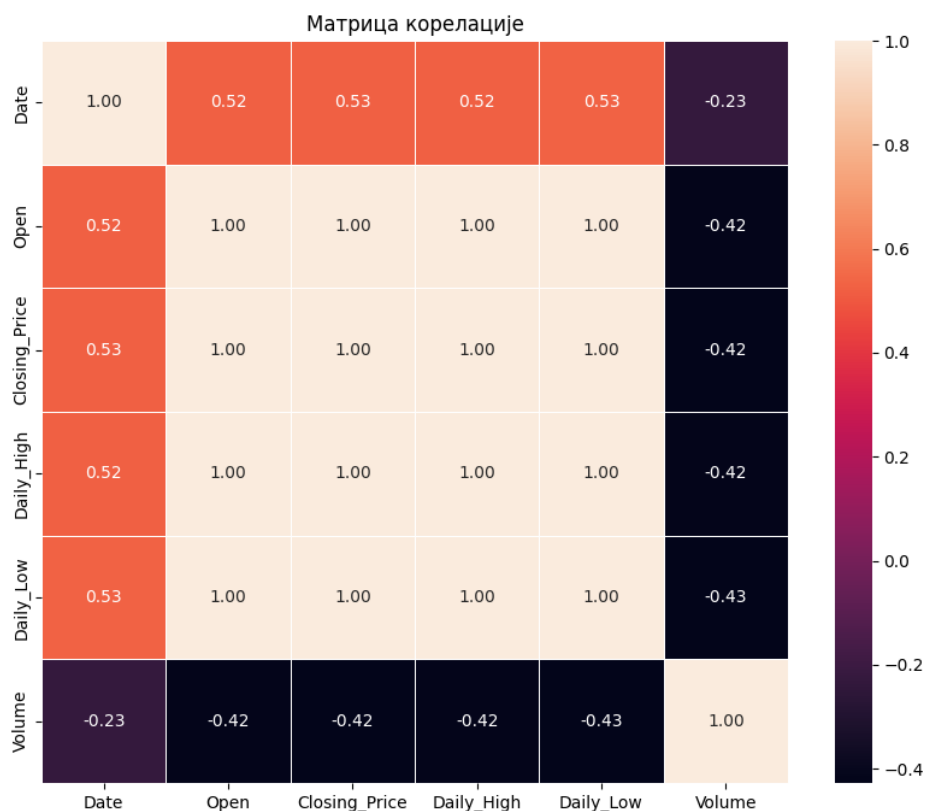
Скуп података састоји се из инстанци које карактеришу наредне карактеристике:

- **Date** – Дан за који је вршен прорачун
- **Daily Low** – Најнижа вредност акције у току дана
- **Daily High** – Највећа вредност акције у току дана
- **Closing Price** – Цена акције на крају дана
- **Open** – Цена акције на почетку дана
- **Volume** – Број продатих акција за дати дан – Обим

Циљ овог рада је прецизно предвидети карактеристику **Closing Price**. Ова карактеристика ће играти улогу излаза мреже.

Како би се видело да ли је могуће на основу осталих карактеристика предвидети излаз, потребно је посматрати колико те карактеристике међусобно зависе, као и колико излаз зависи од њих.

Слика 1 приказује корелациону матрицу ових карактеристика.



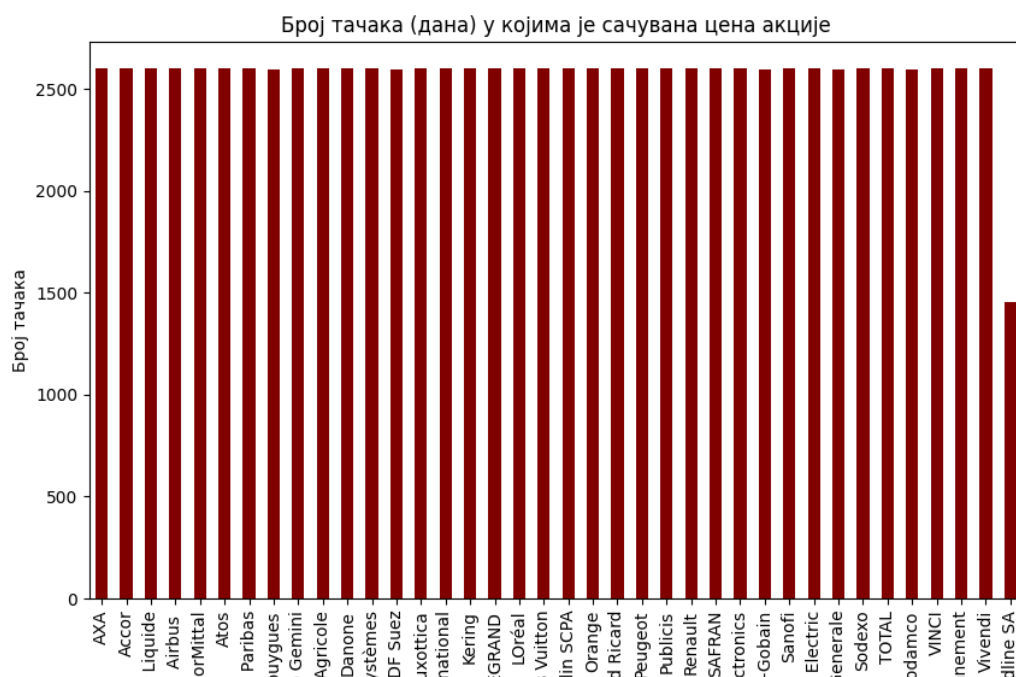
Слика 1: Корелациона матрица

Може се приметити, као што је очекивано, да су готово све карактеристике међусобно потпуно пропорционалне. Када једне расте, и остале расту и обрнуто. Ово их чини непотребним, како не доносе нове информације са собом.

У оквиру рада, изабрана је само карактеристика **Closing Price**, односно вредност те карактеристике у временским тренуцима $t - l_{seq}$ до $t - 1$, где је l_{seq} неки коначан број већи од 1, како би се одредила вредност **Closing Price** у тренутку t .

Како скуп података садржи најуспешнијих 40 француских компанија, потребно је одабрати једну за даљи развој пројекта.

Слика 2 представља број инстанци сваке компаније у оквиру скупа података.



Слика 2: Број инстанци скупа података по компанији

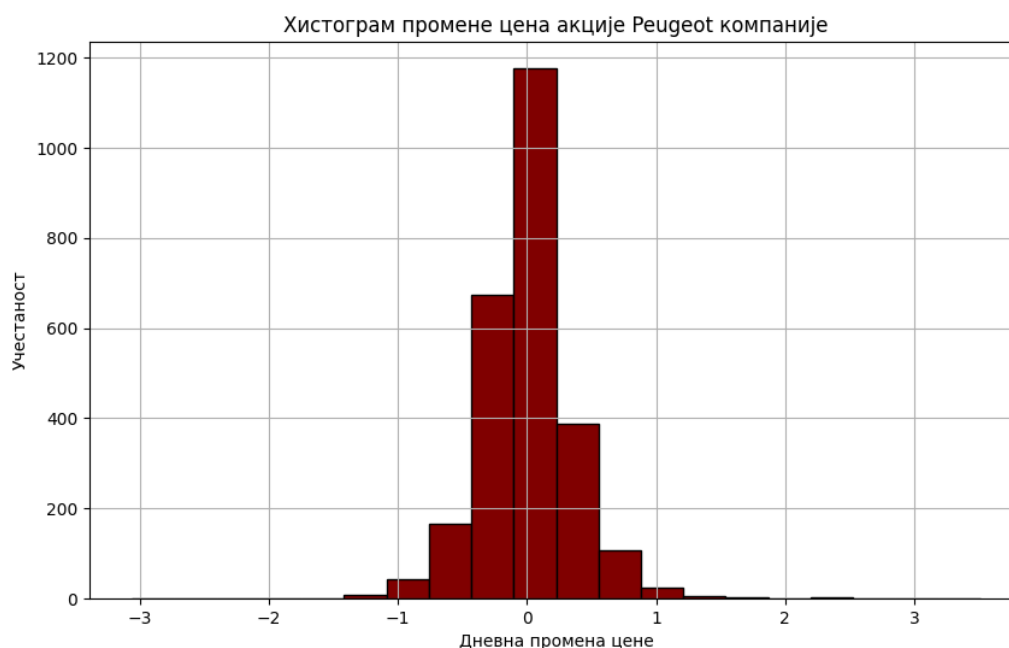
На основу 2 примећује се да све компаније, осим једне, имају идентичан број тачака или инстанци у оквиру овог скупа података, што је јако повољно јер омогућава слободан избор компаније за даљи рад.

Свака инстанца представља вредности карактеристика за један дан једне компаније на тржишту акцијама

За сврху израде овог пројекта одабрана је компанија **Peugeot**, и њена историја на француском тржишту акција приказана у 2601 инстанци (информације дате за 2601 дан).

Како се овај рад бави предвиђањем вредности у датом тренутку на основу вредности из претходних тренутака, потребно је посматрати како се крећу промене вредности **Closing Price** карактеристике из дана у дан за изабрану компанију, са акцентом на разлику вредности карактеристике између две суседне инстанце.

Слика 3 представља визуелни приказ промене вредности **Closing Price** карактеристике међу суседним инстанцама.



Слика 3: Промена вредности **Closing Price** карактеристике

Са слике 3 може се приметити да су разлике у вредностима **Closing Price** минималне и готово да прате нормалну расподелу, што олакшава моделу дугорочне траткотрајне меморије да ефикасније врши предвиђање.

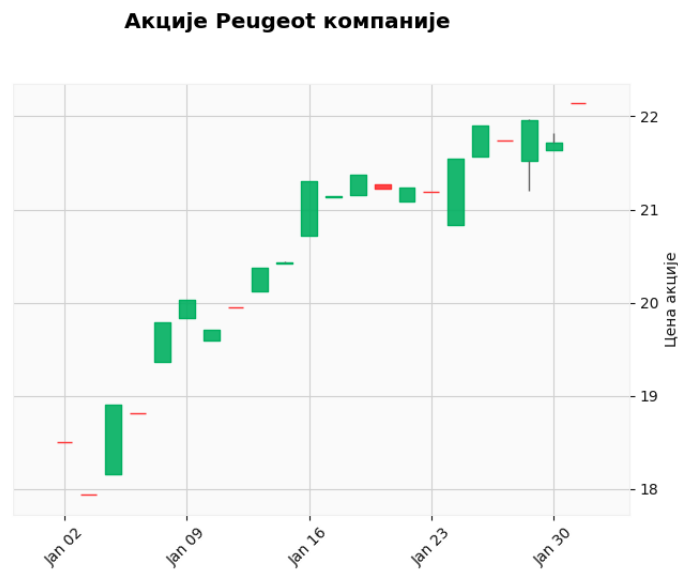
Како се ради о тржишту акција, поучно би било приказати прикладном визуелизацијом потом сагледати скуп података.

Такав приказ подразумева представљање података помоћу **Candlestick** графика. Представљање података помоћу **Candlestick** графика је метода у финансијама која се користи за визуелизацију цена на тржишту у одређеном временском периоду. Ова врста графика састоји се из “свећа”, где свака свећа представља цену атрибута (у овом случају **Closing Price**) у одређеном временском оквиру (у овом случају дан), док крајеви линија представљају највишу и најмању цену у току дана. Боја свеће указује на то да ли је цена порасла (зелена) или опала (црвена) у односу на вредност из претходног временског тренутка. Сликаом 4 дат је пример свеће.

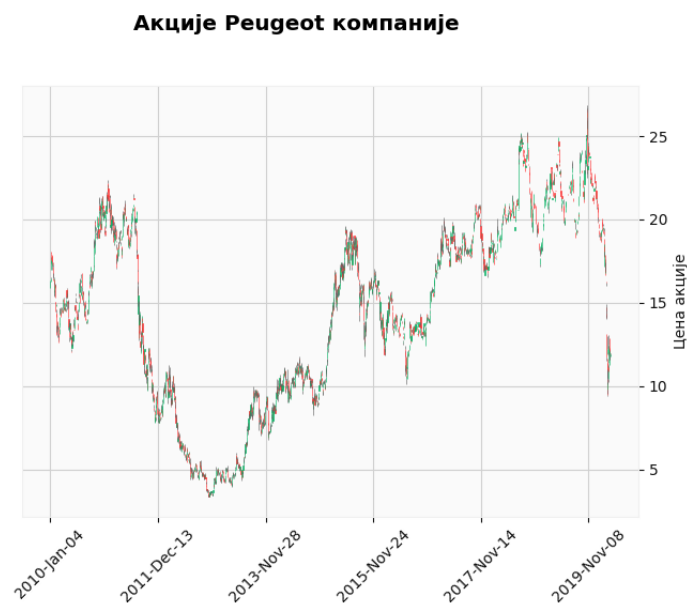


Слика 4: Пример свеће у **Candlestick** графику

Сликама 5 и 6 дат је графички приказ података помоћу **Candlestick** графика.



Слика 5: Вредност акције **Peugeot** компаније у једном месецу



Слика 6: Вредност акције **Peugeot** компаније 2010.–2020.

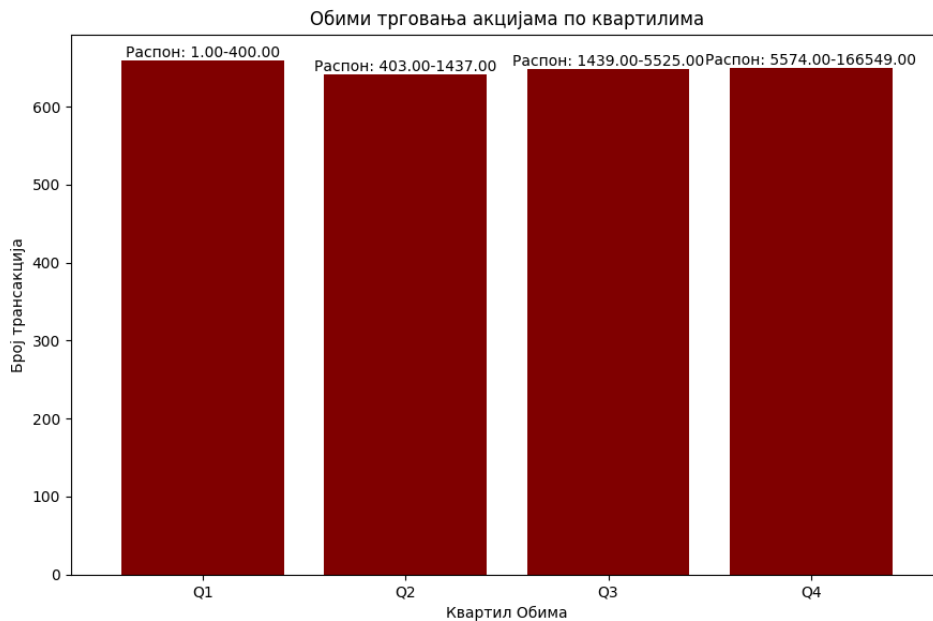
Са слика 5 и 6 може се приметити да, иако су промене у вредностима између суседних инстанци минималне (слика 3), долази до великих флукуација у вредности **Closing Price** карактеристике.

Једноставнији приказ вредности **Closing Price** карактеристике у функцији времена дат је сликом 7.



Слика 7: Једноставнији приказ промене вредности акција

Претходно је поменуто број продатих акција за дати дан као важан фактор у моделирању временске зависности у овом проблему, и иако је та карактеристика уклоњена, посматрањем њених својства могу се добити важне информације. Слика 8 дат је обим трговања акцијом **Peugeot** компаније подељен на четири кватрала идентичне величине.



Слика 8: Једноставнији приказ промене вредности акција

Обим трговања акцијом одабране компаније се може поделити у четири кватила идентичне величине бирањем распона датих на слици 8. На основу резултата те поделе може се приметити да се трговање акцијама одабране компаније обављало у најмањој могућој, као и у значајно великој количини у неким временским тренуцима, што може бити веома важан фактор. Међутим, са корелационе матрице (слика 1) примећује се да је обим заправо обрнуто пропорционалан већини осталих карактеристика, што значи да га је могуће занемарити за сада.

3 Преглед софтверског решења

Преглед софтверског решења подељен је у више секција:

- Визуелизија
- Претпроцесирање
- Креирање, тренирање и тестирање модела мреже

3.1 Визуелизација

Зависности подпрограма за визуелизацију, функција за исцртавање графичког приказа модела, функција за исцртавање графика за приказ броја инстанци по компанији и функција за исцртавање матрице корелације дати су исечком 1, док су исечком 2 дате остале функције за исцртавање.

```
1 from matplotlib import pyplot as plt
2 import pandas as pd
3 import mplfinance as mpf
4 import seaborn as sns
5 from preprocess import load_data
6 import datetime as dt
7 from tensorflow.keras.models import load_model
8 import visualkeras
9 from PIL import ImageFont
10
11 def plot_lstm_model(path):
12     model = load_model(path)
13     font = ImageFont.truetype("arial.ttf", 12)
14     visualkeras.layered_view(model, legend=True, to_file='gen/model_architecture.png', font=font)
15 def plot_komp_cnt(name_counts):
16     name_counts.plot(kind='bar', figsize=(10, 6), color='maroon')
17     plt.title('Број тачака (дана) у којима је сачувана цена акције')
18     plt.xlabel('Име компаније')
19     plt.ylabel('Број тачака')
20     plt.savefig("gen/brdana_komp.png")
21     plt.close()
22 def plot_corr(df):
23     cm = df.corr()
24     plt.figure(figsize=(10, 8))
25     sns.heatmap(cm, annot=True, cmap='rocket', fmt=".2f", linewidths=0.5)
26     plt.title('Матрица корелације')
27     plt.savefig("gen/cm.png")
28     plt.close()
```

Listing 1: Зависности и функције за исцртавање

```

1 def plot_closing_price(df):
2     plt.figure(figsize=(12, 6))
3     plt.plot(df['Date'], df['Closing_Price'], label='Цена затварања акција', color='maroon')
4     plt.title('Peugeot - цена затварања акција у функцији времена')
5     plt.xlabel('Датум')
6     plt.ylabel('Цена затварања акција')
7     plt.legend()
8     plt.savefig("gen/cena_zatv_akc.png")
9 def plot_daily_price_changes(df):
10    plt.figure(figsize=(10, 6))
11    df['Daily_Price_Change'] = df['Closing_Price'].diff()
12    df['Daily_Price_Change'].hist(bins=20, edgecolor='black', color='maroon')
13    plt.title('Хистограм промене цена акције Peugeot компаније')
14    plt.xlabel('Дневна промена цене')
15    plt.ylabel('Учестаност')
16    plt.savefig("gen/promena_cene.png")
17 def plot_candlestick_daytoday(df):
18    sd, ed = pd.Timestamp(2019,1,1), pd.Timestamp(2019,1,31)
19    df_month = df[(df['Date'] >= sd) & (df['Date'] <= ed)]
20    ohlc = df_month[['Date', 'Open', 'Daily_High', 'Daily_Low', 'Closing_Price']]
21    .set_index('Date')
22    ohlc.columns = ['Open', 'High', 'Low', 'Close']
23    mpf.plot(ohlc, type='candle', style='yahoo',
24    title='Акције Peugeot компаније', ylabel='Цена акције',
25    show_nontrading=False, savefig="gen/marketplot_jan.png")
26 def plot_volume(df):
27    df['Volume_Quartiles'] = pd.qcut(df['Volume'], q=[0, 0.25, 0.5, 0.75, 1.0],\
28    labels=['Q1', 'Q2', 'Q3', 'Q4'])
29    volume_quartile_counts = df['Volume_Quartiles'].value_counts()
30    quartile_ranges = df.groupby('Volume_Quartiles')['Volume'].agg(['min', 'max'])
31    plt.figure(figsize=(9, 6))
32    volume_quartile_counts.sort_index().plot(kind='bar',\
33    color='maroon', width=0.8)
34    bars = plt.bar(volume_quartile_counts.index, volume_quartile_counts.sort_index(),\
35    color='maroon', width=0.8)
36    for bar, (q, (min_val, max_val)) in zip(bars, quartile_ranges.iterrows()):
37        plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height(),
38        f'Распон: {min_val:.2f}-{max_val:.2f}', ha='center', va='bottom')
39    plt.title('Обими трговања акцијама по квантилима')
40    plt.xlabel('Квантил Обима')
41    plt.ylabel('Број трансакција')
42    plt.tight_layout()
43    plt.savefig("gen/volume_quartiles_plot.png")

```

Listing 2: Остале функције за исцртавање

3.2 Претпроцесирање

Исечком 3 дата је функција за читавање и “чишћење” скупа података.

```
1 def load_data(start,end):
2     df = pd.read_csv("preprocessed_CAC40.csv")
3     df = df.interpolate(method='linear')
4     df['Date'] = pd.to_datetime(df['Date'])
5     df.drop(df.filter(regex="Unname"),axis=1, inplace=True)
6     df = df[df["Name"] == "Peugeot"]
7     df = df.loc[(df['Date'] > start) & (df['Date'] < end), :]
8     df['Volume'] = pd.to_numeric(df['Volume'].str.replace(', ', ''), errors='raise')
9     df['Volume'] = df['Volume'].ffill()
10    df = df.sort_values('Date')
11    return df
```

Listing 3: Функција за читавање и “чишћење” скупа података

Ова функција отклања проблеме код којих су неке карактеристике инстанци без вредности, на више различитих начина.

Исечком 4 дата је функција за креирање секвенци.

```
1 def create_sequences(arr, seq_len):
2     X_train, y_train = [], []
3     for x in range(seq_len, len(arr)):
4         X_train.append(arr[x - seq_len : x, -1])
5         y_train.append(arr[x, -1])
6     X_train, y_train = np.array(X_train), np.array(y_train)
7     return X_train, y_train
```

Listing 4: Функција за креирање секвенци

На основу параметра временског корака, генеришу се низови улазних података. За сваку инстанцу везујемо претходних N инстанци како би се моделирала временска зависност.

Поред ових функција, врши се и скалирање. Скалирање се врши у главном програму, како је постојала потреба за тиме због зависности појединих објеката. Исечком 5 дата је функција скалирања (која се налази у главном програму).

```
1 scaler = MinMaxScaler(feature_range=(0, 1))
2 def preprocess_data(data, scaler, cols):
3     temp = data.copy(deep=True)
4     temp[cols] = scaler.fit_transform(temp[cols])
5     return temp[cols].values
```

Listing 5: Функција за скалирање

3.3 Креирање, тренирање и тестирање модела мреже

Зависности су дате исечком 6.

```
1 from matplotlib import pyplot as plt
2 import pandas as pd
3 import numpy as np
4 from tensorflow.keras.models import Sequential
5 from tensorflow.keras.layers import LSTM, Dense, Dropout, BatchNormalization
6 from preprocess import load_data, create_sequences
7 from tensorflow.keras.callbacks import ModelCheckpoint, EarlyStopping, ReduceLROnPlateau, TensorBoard
8 from sklearn.preprocessing import MinMaxScaler
9 import datetime as dt
```

Listing 6: Зависности главног програма

Процес креирања модела мреже обавља функција дата исечком 7.

```
1 def build_lstm_model(input_shape=(60, 1)): #vremenski korak = 60
2     model = Sequential()
3     model.add(LSTM(50, return_sequences=True, input_shape=input_shape))
4     model.add(BatchNormalization())
5     model.add(Dropout(0.2))
6     model.add(LSTM(units=50, return_sequences=True))
7     model.add(BatchNormalization())
8     model.add(Dropout(0.2))
9     model.add(LSTM(units=50))
10    model.add(Dense(1))
11    model.compile(optimizer='adam', loss='mean_squared_error')
12    return model
```

Listing 7: Функција за креирање мреже

Функције за тренирање мреже и припрему тестног скупа су дата исечком 8.

```
1 def train_model(model, X_train, y_train, epochs=100, batch_size=32, validation_split=0.15):
2     checkpointer = ModelCheckpoint(filepath='best_weights.h5',
3         monitor='val_loss',
4         verbose=1,
5         save_best_only=True,
6         mode='min')
7     early_stopping = EarlyStopping(monitor='val_loss',
8         patience=15,
9         restore_best_weights=True)
10    reduce_lr = ReduceLROnPlateau(monitor='val_loss',
11        factor=0.2,
12        patience=3,
13        min_lr=1e-5)
14    tensorboard = TensorBoard(log_dir='logs',
15        histogram_freq=1,
16        write_graph=True,
17        write_images=True)
18    model.fit(X_train,
19        y_train,
20        epochs=epochs,
21        batch_size=batch_size,
22        validation_split=validation_split,
23        callbacks=[checkerpointer, early_stopping, reduce_lr, tensorboard])
24
25 def prepare_test_data(train_data, test_data, scaler, seq_length=60):
26     actual_prices = test_data['Closing_Price'].values
27     total_dataset = pd.concat((train_data['Closing_Price'], test_data['Closing_Price']), axis=0)
28     model_inputs = total_dataset[len(total_dataset) - len(test_data) - seq_length:].values
29     model_inputs = model_inputs.reshape(-1, 1)
30     model_inputs = scaler.transform(model_inputs)
31     X_test, _ = create_sequences(model_inputs, seq_length)
32     X_test = np.reshape(X_test, (X_test.shape[0], X_test.shape[1], 1))
33     return actual_prices, X_test
```

Listing 8: Функције за тренинг и припрему тестног скупа података

Функција за приказ резултата мреже и главни програм приказани су исечком 9

```
1 def plot_results(actual_prices, predicted_prices, title):
2     print(np.shape(actual_prices))
3     plt.plot(actual_prices, color='black', label="Права вредност")
4     plt.plot(predicted_prices, color='green', label="Предложена вредност")
5     plt.title(title)
6     plt.xlabel("Време")
7     plt.ylabel("Peugeot цена акција")
8     plt.legend()
9     plt.savefig("gen/rezultat.png")
10    plt.close()
11
12    START_DATE = dt.datetime(2015, 1, 1)
13    END_DATE = dt.datetime(2019, 12, 30)
14    START_DATE_TEST = END_DATE
15
16    scaler = MinMaxScaler(feature_range=(0, 1))
17    cols = ['Closing_Price']
18
19    train_data = load_data(START_DATE, END_DATE)
20    temp_train = preprocess_data(train_data, scaler, cols)
21    X_train, y_train = create_sequences(temp_train, 60)
22
23    model = build_lstm_model(input_shape=(60, 1))
24    train_model(model, X_train, y_train)
25
26    test_data = load_data(START_DATE_TEST, dt.datetime(2020, 12, 30))
27    actual_prices, X_test = prepare_test_data(train_data, test_data, scaler)
28
29    predicted_prices = model.predict(X_test)
30    predicted_prices = scaler.inverse_transform(predicted_prices)
31
32    plot_results(actual_prices, predicted_prices, "Peugeot цена акција")
```

Listing 9: Функција за приказ резултата мреже и главни програм

Функција за приказ резултата мреже генерише график који приказује предвиђања мреже за 67 дана. Тренирање мреже вршило се над подацима од 1.1.2015. до 30.12.2019. (31.12.2019. не садржи икакве податке), док су за тестни скуп узети подаци од 1.1.2020. до 4.3.2020. На графику се исцртавају предвиђања као и тачне вредности, ради лакшег поређења.

4 Резултати

Улаз мреже је димензија (60,1), како је претходно постављен временски корак од вредности 60, што значи да ће мрежа за цену акције C_t у обзир узети цене акција $C_{t-1} \dots C_{t-60}$.

Тренирање мреже вршено је над 100 епоха са димензијом подскупа 32 (**batch size**), где је постављено пар правила које условљавају мрежу да прекида са обучавањем уколико не дође до побољшања, као и динамичка промена стопе учења и чување најефикаснијих тежина модела мреже.

Модел над тренинг скупом након 100 епоха остварује грешку (**MSE**) у вредности од 0.0018, док над валидационим скупом (узето 15% од оригиналног тренинг скупа) остварује грешку од 0.00199.

Како је речено у претходној секцији, мрежа се тестира над скупом података цена акција од 1.1.2020. до 4.3.2020.

Резултати које је модел остварио дати су сликом 9.



Слика 9: Резултат – поређење предвиђања мреже и правих вредности

Са слике 9 може се приметити да је мрежа остварила добар резултат, и што се тиче појединих цена акција и праћења линије тренда цене.