

# Bitcoin Price Prediction Project with ARIMA, XGBoost and LSTM

## LSTM

This notebook is part of Bitcoin Price Prediction Project with LSTM model. Complete project consists of three notebooks with sections as below.

1. Abstract
2. Data description # (01\_BTC\_Price\_Prediction - ARIMA.ipynb notebook)
3. Gathering data # (01\_BTC\_Price\_Prediction - ARIMA.ipynb notebook)
4. Data preprocessing
  - 5.1. ARIMA => Introduction
  - 5.2. ARIMA Model => function definition
  - 5.3. ARIMA Model => Forecasting
    - complete data
    - data start from 2018-01-01
    - data start from 2021-01-01
  - 5.4. ARIMA Model => summary
6. XGBoost # (02\_BTC\_Price\_Prediction - XGBoost.ipynb notebook)
  - 6.1. XGBRegressor => Introduction
  - 6.2. XGBRegressor => function definition
  - 6.3. XGBRegressor => Forecasting
    - complete data
    - data start from 2018-01-01
    - data start from 2021-01-01
  - 6.4. XGBRegressor => summary
7. LSTM Model # (this notebook)
  - 7.1. LSTM => Introduction
  - 7.2. LSTM => function definition
  - 7.3. LSTM => Forecasting
    - complete data
    - data start from 2018-01-01
    - data start from 2021-01-01
  - 7.4. LSTM => summary
8. Project summary # (this notebook)

```
In [1]: import pandas as pd
import numpy as np
from numpy import array

from datetime import datetime, date

import matplotlib.pyplot as plt

from tqdm import tqdm
```

## 4. Data preprocessing

```
In [2]: # restore data backup
btc_df_gathered = pd.read_csv("BTC_Preprocess_Backup.csv", sep = ",", skipinitialspace=True)
btc_df_gathered.head()
```

Out[2]:

	time_open	time_close	open	high	low	close	market_cap	volume
0	2010-07-17T00:00:00Z	2010-07-17T23:59:59Z	0.04951	0.04951	0.04951	0.04951	NaN	NaN
1	2010-07-18T00:00:00Z	2010-07-18T23:59:59Z	0.04951	0.04951	0.04951	0.04951	NaN	NaN
2	2010-07-19T00:00:00Z	2010-07-19T23:59:59Z	0.08584	0.08584	0.08584	0.08584	NaN	NaN
3	2010-07-20T00:00:00Z	2010-07-20T23:59:59Z	0.08080	0.08080	0.08080	0.08080	NaN	NaN
4	2010-07-21T00:00:00Z	2010-07-21T23:59:59Z	0.07474	0.07474	0.07474	0.07474	NaN	NaN

```
In [3]: # copy data to new df
btc_df_1 = btc_df_gathered
```

```
In [4]: # check datatypes
btc_df_1.dtypes
```

Out[5]:

	time_open	time_close	open	high	low	close	market_cap	volume
time_open	object							
time_close	object							
open	float64							
high	float64							
low	float64							
close	float64							
market_cap	float64							
volume	float64							
dtype:	object							

```
In [5]: # change df index to date
btc_df_1.index = pd.to_datetime(btc_df_1["time_close"]).dt.date

# drop "time_open", "time_close" columns
btc_df_1.drop(["time_open", "time_close"], axis=1, inplace=True)
btc_df_1.head()
```

Out[5]:

	open	high	low	close	market_cap	volume
time_close						
2010-07-17	0.04951	0.04951	0.04951	0.04951	NaN	NaN
2010-07-18	0.04951	0.04951	0.04951	0.04951	NaN	NaN
2010-07-19	0.08584	0.08584	0.08584	0.08584	NaN	NaN
2010-07-20	0.08080	0.08080	0.08080	0.08080	NaN	NaN
2010-07-21	0.07474	0.07474	0.07474	0.07474	NaN	NaN

```
In [6]: # show df shape
btc_df_1.shape
```

Out[6]: (4276, 6)

```
In [7]: btc_df_1.describe()
```

Out[7]:

	open	high	low	close	market_cap	volume
count	4276.000000	4276.000000	4276.000000	4276.000000	3260000e+03	3.017000e+03
mean	7620.722003	7818.648308	7411.105750	7632.136514	1.835636e+11	1.205864e+10
std	14287.143745	14651.680612	13896.018691	14100.996427	2.947438e+11	1.710047e+10
min	0.049510	0.049510	0.049510	0.049510	7.784112e+08	2.857830e+06
25%	111.217500	114.591000	107.450000	111.500000	6.510601e+09	5.567090e+07
50%	638.137000	656.522000	627.096000	638.393000	6.466303e+10	3.682395e+09
75%	8191.225324	8321.614645	7952.276240	8196.622730	1.752686e+11	1.885380e+10
max	67589.872823	68692.137037	66457.970074	67589.768671	1.275266e+12	1.363702e+11

## 7. LSTM Model

### 7.1. LSTM => Introduction

This notebook was created to predict the Bitcoin price (in USD) using the LSTM model. The model has been programmed in a function to perform a forecast with different data ranges. LSTM model contains parameters that have been empirically optimized.

An additional function has been produced to prepare a sequence of a data - number of past days to look for one future day forecasting. In this notebook, LSTM model will look just one day in the history to predict next day.

The mean square error (RMSE) was used to evaluate the performance of the model - it should be as low as possible.

Training data is the first 80% of the dataset and the rest is in the test datasets.

At the end of the notebook, a summary of the task is presented along with a graph of the RMSE results of the model.

Long Short-Term Memory - in short LSTM models are type of recurrent neural network (RNN) capable of learning sequences of observations. This makes LSTM well suited for time series forecasting.

LSTM works better with normalized data. For that reason, MinMaxScaler has been used in range from 0 to 1. Reshape data for purpose of LSTM has been done as well.

After fitting, MinMaxScaler provides "inverse\_transform" method to receive more "readable" results.

EarlyStopping was used and penalty is set to 12(0.00001) in both data variants.

```
In [8]: from sklearn.metrics import mean_squared_error

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense

from tensorflow.keras.callbacks import EarlyStopping
from tensorflow.keras.regularizers import L2
```

### 7.2. LSTM => function definition

```
In [9]: # split data into sequence
def split_data(data, n_steps):

    """
    Function for determinate the sequence pattern of the data.
    n_steps is the amount of days to look back for prediction of the next day;
    n_steps = 2 ==> look back in to 2 day to predict day 3;
    n_steps = 3 ==> look back in to 3 day to predict day 4.
    """

    X = []
    y = []

    for i in range(len(data)):

        # find the end of the sequence
        end_seq = i + n_steps

        # check if we are bout of the sequence
        if end_seq > len(data)-1:
            break

        # get input and output parts of the sequence
        seq_x, seq_y = data[i:end_seq], data[end_seq]
        X.append(seq_x)
        y.append(seq_y)

    return array(X), array(y)

In [10]: def lstm_variants(df_variants: list, train_size: float, steps: int):

    """
    This function includes LSTM modeling for Bitcoin price prediction.
    Parameters have been empirically optimized.
    The function is prepared for prediction, based on a sets of data - eg. different ranges of historical data.
    """

    variant = 0 # data variant number
    rmse = [] # LSTM rmse

    # general loop fo data vizairints (data range)
    for df in tqdm(df_variants):

        # -----
        # increment for data variant number
        variant = variant + 1
        print(f"Data variant no: ", variant, "\n")
        # -----

        # -----
        # train test, split the data
        # train = data * train_size
        # test = data * (1 - train_size)
        n_steps = int((len(df) * train_size))
        train_data, test_data = df[:n_steps], df[n_steps:]
        # -----

        # split data into sequence
        X_train, y_train = make_sequence(train_data, steps)
        X_test, y_test = make_sequence(test_data, steps)

        # reshape from [samples, timesteps] into [samples, timesteps, features]
        n_features = 1
        X_train = X_train.reshape((X_train.shape[0], X_train.shape[1], n_features))
        X_test = X_test.reshape((X_test.shape[0], X_test.shape[1], n_features))
        # -----

        # define early_stopping and penalty
        reg = L2(0.00001)
        early_stopping = EarlyStopping(monitor="val_loss", # Quantity to be monitored.
                                      patience=20, # Number of epochs with no improvement after
                                      restore_best_weights=True) # Whether to restore model weights from
        # -----
        # the epoch with the best value of the monitor

        # -----
        # create model
        model=Sequential()

        # first LSTM layer
        model.add(LSTM(units=50, # Positive integer, dimensionality of the output space
                       activation="relu", # Activation function to use. Default: 'tanh'
                       input_shape=(steps, n_features), # input data shape - "None" will accept "any" given
                       kernel_regularizer=reg # Regularizer function applied to the 'kernel' weights
                       ))

        # output layer
        model.add(Dense(1)) # units: Positive integer, dimensionality of the output space.
        # Here equals 1 as 1 prediction.

        model.summary()

        # compile the model
        model.compile(optimizer="adam", # String (name of optimizer) or optimizer instance.
                     loss="mean_squared_error") # Loss function. Maybe be a string (name of loss function)

        # fit model to the training data
        history = model.fit(
            X_train, y_train, # Number of epochs to train the model.
            epochs=200, # Number of samples per gradient update. Default: 32
            batch_size=64,
            validation_data=(X_test, y_test), # Data on which to evaluate the loss and any model metrics at the
            callbacks=[early_stopping, # List of callbacks to apply during training. Here early stop
                                VerboseMode(0 = silent, 1 = progress bar
                                )
            )

        # print fitting params
        print("\n", history.params)
        # -----

        # -----
        # plot training losses
        loss = history.history["loss"]
        val_loss = history.history["val_loss"]
        epochs = range(len(loss))

        plt.figure(figsize=(8,4), dpi=100)
        plt.plot(epochs, loss, "green", label="Training loss")
        plt.plot(epochs, val_loss, "blue", label="Training loss")
        plt.title("Training and test loss")
        plt.xlabel("Epochs")
        plt.ylabel("Loss")
        plt.legend(loc="best", fontsize=8)
        plt.show()
        # -----

        # -----
        # X test reshape
        X_input = x.reshape((1, steps, n_features)) for x in X_test

        # prediction
        y_pred = (model.predict(X, verbose=0)) for x in X_input
        y_pred = y_pred[0][0] for i in range(len(y_pred))
        y_pred = pd.Series(y_pred, index=test_data[steps:].index)
        # -----

        # -----
        # plot results with train data
        plt.figure(figsize=(10,5), dpi=100)
        plt.plot(x_train_data, "green", label="Train data")
        plt.plot(test_data, "blue", label="Actual Price/ Test data")
        plt.plot(y_pred, "orange", label="Predicted Price")

        plt.title("BTC Price Prediction - with forecast")
        plt.xlabel("Time")
        plt.ylabel("BTC Price")
        plt.legend(loc="best", fontsize=8)
        plt.show()

        # plot results predictions and test data - zoom-in results
        plt.figure(figsize=(10,5), dpi=100)
        plt.plot(test_data, "blue", label="Actual Price/ Test data")
        plt.plot(y_pred, "orange", label="Predicted Price")

        plt.title("BTC Price Prediction - with forecast")
        plt.xlabel("Time")
        plt.ylabel("BTC Price")
        plt.legend(loc="best", fontsize=8)
        plt.show()
        # -----

        # -----
        # print performance
        rmse_1 = mean_squared_error(test_data[steps:].values, y_pred.values, squared=False)
        print(model._class_name_, "variant:", variant, ":", "RMSE = %.3f" % rmse_1, "\n")
        rmse.append(float(format(rmse_1, ".3f")))
        # -----

    # end of general loop for data vizairints (data range)

    return rmse
```

### 7.3. LSTM => Forecasting

```
In [11]: # copy full dataset
df_var_1 = btc_df_1["close"]

# copy data starting from 2018-01-01 and replace NaN values to 0
df_var_2 = btc_df_1["close"].iloc[2725:]

# copy data starting from 2021-01-01 and replace NaN values to 0
df_var_3 = btc_df_1["close"].iloc[3821:]

data_variants = [
    df_var_1,
    df_var_2,
    df_var_3
]
```

```
In [12]: # Forecasting
lstm_variants_rmse = lstm_variants(df_variants=data_variants, train_size=0.8, steps=1)
```



```

0%) | 0/3 [00:00<?,
Data variant no: 1
Model: "sequential"
Layer (type)                Output
-----
lstm (LSTM)                  (None)
dense (Dense)                (None)

Total params: 10,451
Trainable params: 10,451
Non-trainable params: 0

Epoch 1/200
54/54 [=====]
Epoch 2/200
54/54 [=====]
Epoch 3/200
54/54 [=====]
Epoch 4/200
54/54 [=====]
Epoch 5/200

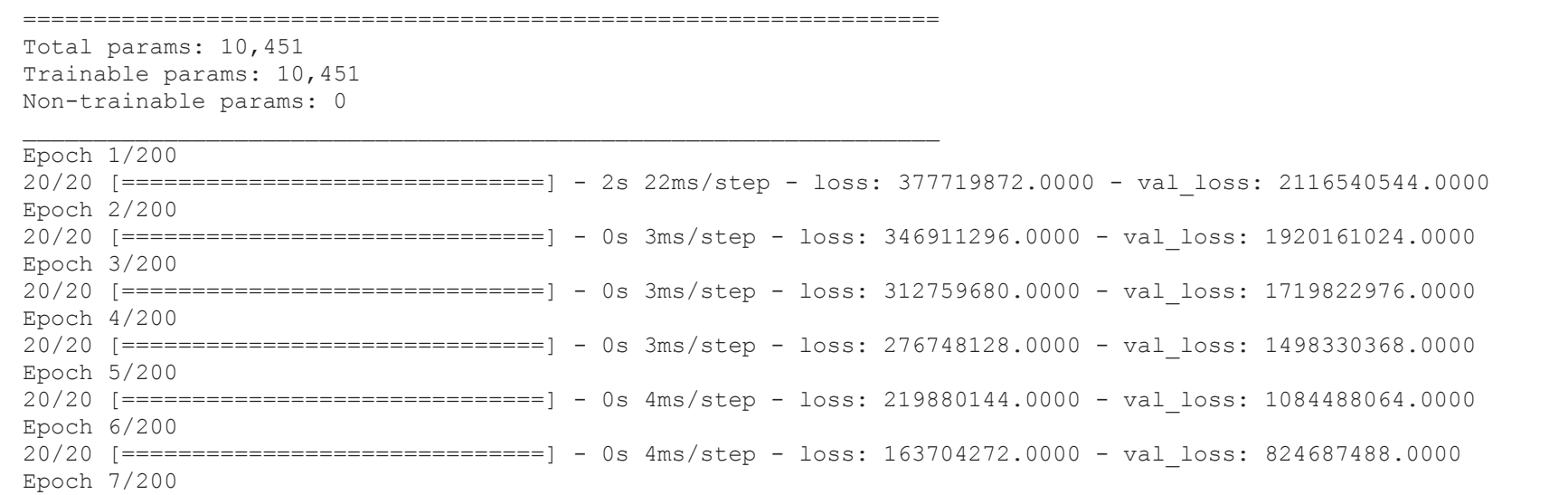
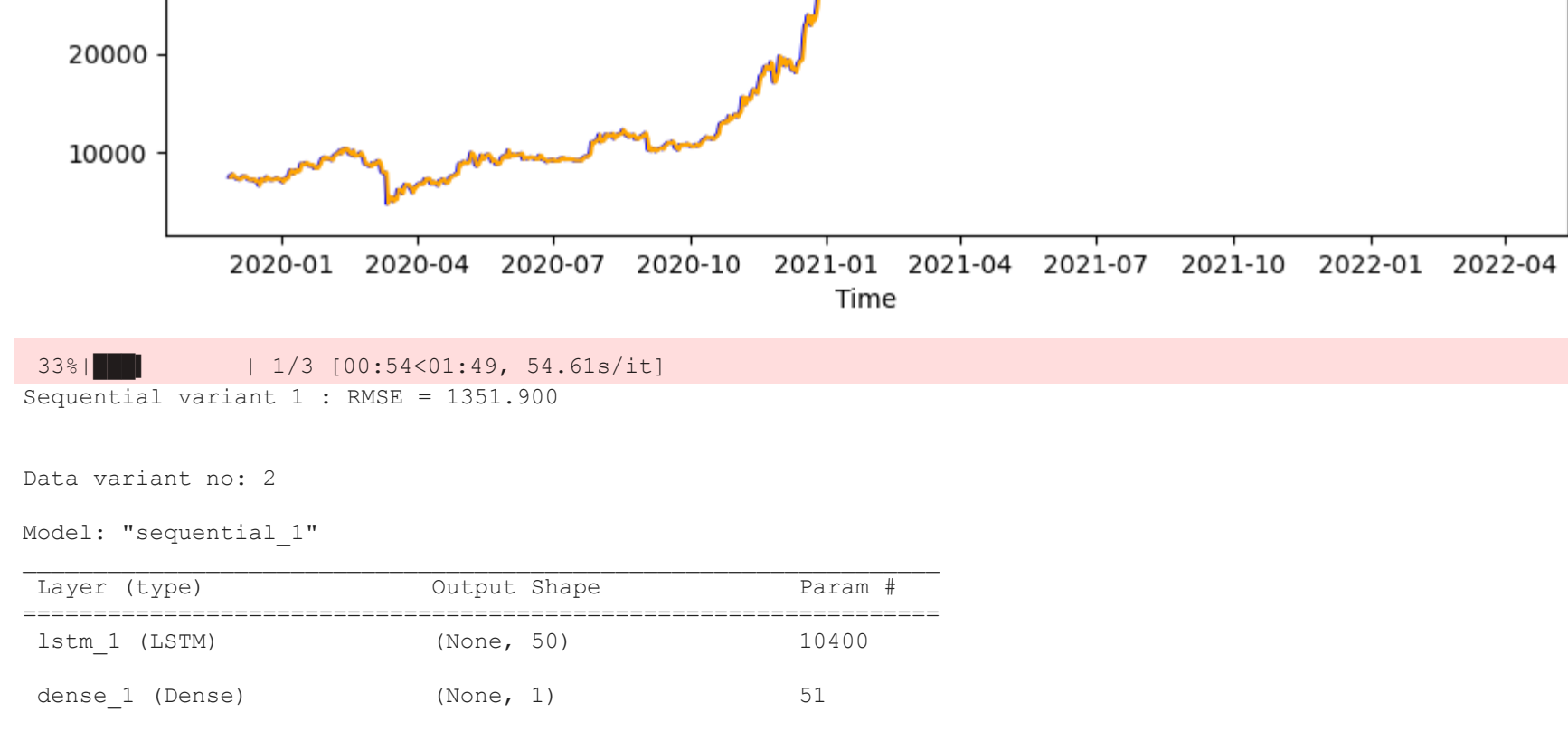
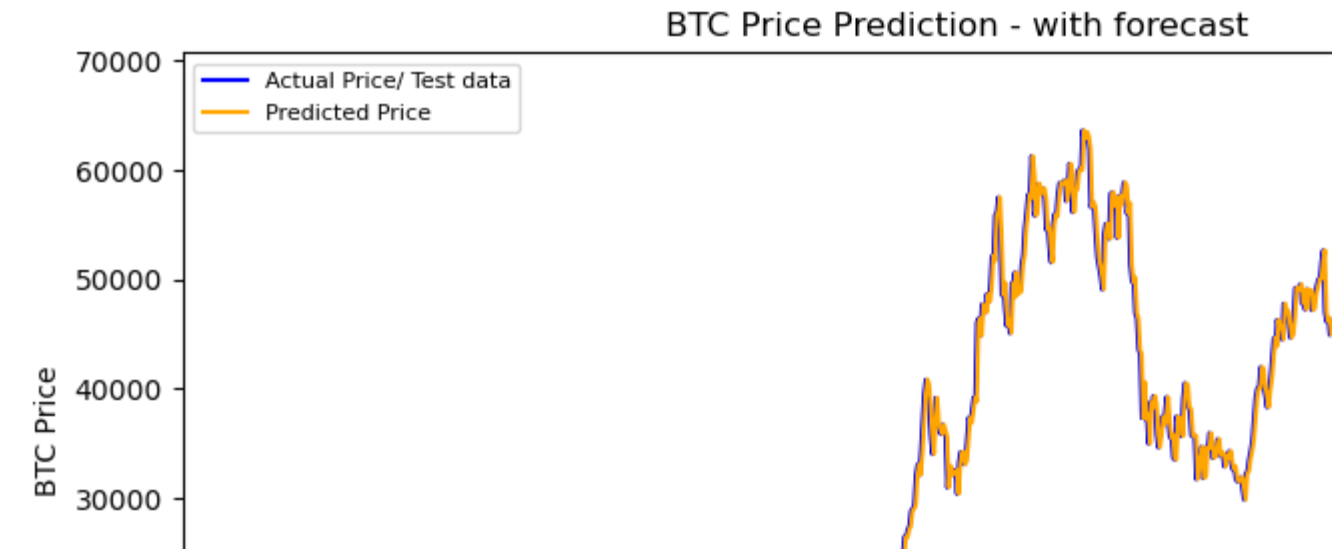
```

Epoch 6/200  
54/54 [=====] - 0s 3ms/step - loss: 44296.5078 - val\_loss: 1859150.2500  
Epoch 7/200  
54/54 [=====] - 0s 7ms/step - loss: 44130.6172 - val\_loss: 1824081.3750  
Epoch 8/200  
54/54 [=====] - 0s 3ms/step - loss: 44232.6016 - val\_loss: 1834264.8750  
Epoch 9/200  
54/54 [=====] - 0s 3ms/step - loss: 44537.0586 - val\_loss: 1833163.7500  
Epoch 10/200  
54/54 [=====] - 0s 3ms/step - loss: 44216.3438 - val\_loss: 1838281.7500  
Epoch 11/200  
54/54 [=====] - 0s 3ms/step - loss: 44077.2393 - val\_loss: 1828905.0000  
Epoch 12/200  
54/54 [=====] - 0s 3ms/step - loss: 44174.2305 - val\_loss: 1847508.3750  
Epoch 13/200  
54/54 [=====] - 0s 3ms/step - loss: 44360.3672 - val\_loss: 1861474.7500  
Epoch 14/200  
54/54 [=====] - 0s 3ms/step - loss: 44463.0820 - val\_loss: 1840830.1250  
Epoch 15/200  
54/54 [=====] - 0s 3ms/step - loss: 44347.4414 - val\_loss: 1827875.5000  
Epoch 16/200  
54/54 [=====] - 0s 3ms/step - loss: 44444.6250 - val\_loss: 1827778.6250  
Epoch 17/200  
54/54 [=====] - 0s 3ms/step - loss: 44109.7852 - val\_loss: 1902953.0000  
Epoch 18/200  
54/54 [=====] - 0s 3ms/step - loss: 46300.9727 - val\_loss: 1917391.7500  
Epoch 19/200  
54/54 [=====] - 0s 3ms/step - loss: 45036.0273 - val\_loss: 1827631.8750  
Epoch 20/200  
54/54 [=====] - 0s 3ms/step - loss: 44983.6484 - val\_loss: 1830603.5000  
Epoch 21/200  
54/54 [=====] - 0s 3ms/step - loss: 44838.0547 - val\_loss: 1914473.7500  
Epoch 22/200  
54/54 [=====] - 0s 3ms/step - loss: 45255.1641 - val\_loss: 1867231.2500  
Epoch 23/200  
54/54 [=====] - 0s 3ms/step - loss: 44412.6250 - val\_loss: 1847387.3750  
Epoch 24/200  
54/54 [=====] - 0s 3ms/step - loss: 44365.0703 - val\_loss: 1882596.0000  
Epoch 25/200  
54/54 [=====] - 0s 3ms/step - loss: 44282.9805 - val\_loss: 1904663.7500  
Epoch 26/200  
54/54 [=====] - 0s 3ms/step - loss: 44120.9414 - val\_loss: 1830839.1250  
Epoch 27/200  
54/54 [=====] - 0s 3ms/step - loss: 44863.4375 - val\_loss: 1837792.0000  
Epoch 28/200  
54/54 [=====] - 0s 3ms/step - loss: 44502.9297 - val\_loss: 1828243.0000  
Epoch 29/200  
54/54 [=====] - 0s 3ms/step - loss: 44607.2930 - val\_loss: 1828883.1250  
Epoch 30/200  
54/54 [=====] - 0s 3ms/step - loss: 44290.0352 - val\_loss: 2124370.2500  
Epoch 31/200  
54/54 [=====] - 0s 3ms/step - loss: 45847.9023 - val\_loss: 1918996.0000  
Epoch 32/200  
54/54 [=====] - 0s 3ms/step - loss: 44760.1523 - val\_loss: 1981174.6250  
Epoch 33/200  
54/54 [=====] - 0s 3ms/step - loss: 45051.8789 - val\_loss: 1865748.2500  
Epoch 34/200  
54/54 [=====] - 0s 3ms/step - loss: 45004.6992 - val\_loss: 1848342.6250  
Epoch 35/200  
54/54 [=====] - 0s 3ms/step - loss: 44696.7773 - val\_loss: 1847893.5000  
Epoch 36/200  
54/54 [=====] - 0s 3ms/step - loss: 44523.6602 - val\_loss: 1882308.7500  
Epoch 37/200  
54/54 [=====] - 0s 3ms/step - loss: 44621.5352 - val\_loss: 1838324.7500  
Epoch 38/200  
54/54 [=====] - 0s 3ms/step - loss: 44392.3663 - val\_loss: 1839152.5000  
Epoch 39/200  
54/54 [=====] - 0s 3ms/step - loss: 44208.2891 - val\_loss: 1883848.5000

(verbose=1, 'epoch': 200, 'steps': 54)

**1e8 Training and test loss**

**BTC Price Prediction - with forecast**



```
20/20 [-----]
Epoch 8/200
20/20 [-----]
Epoch 9/200
```

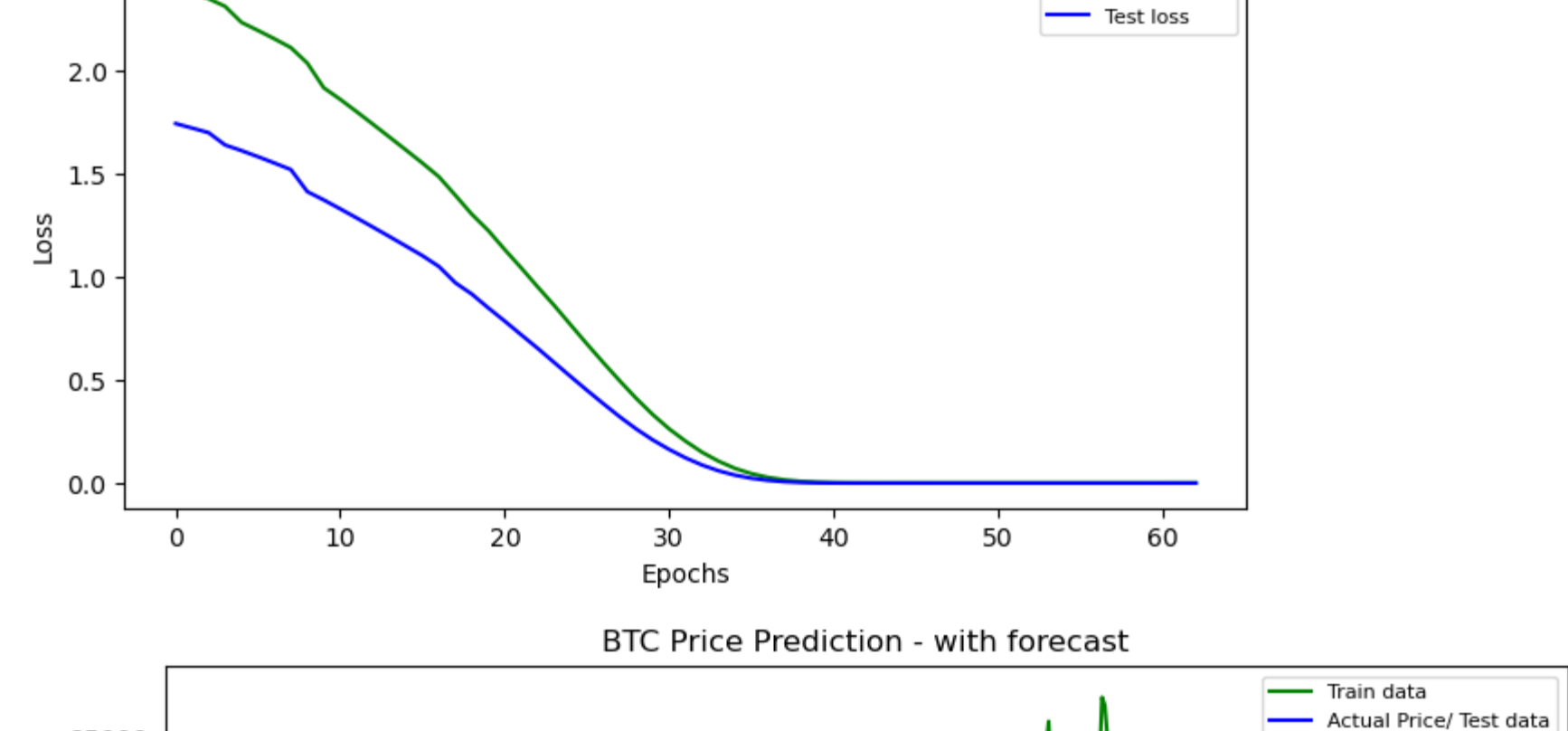
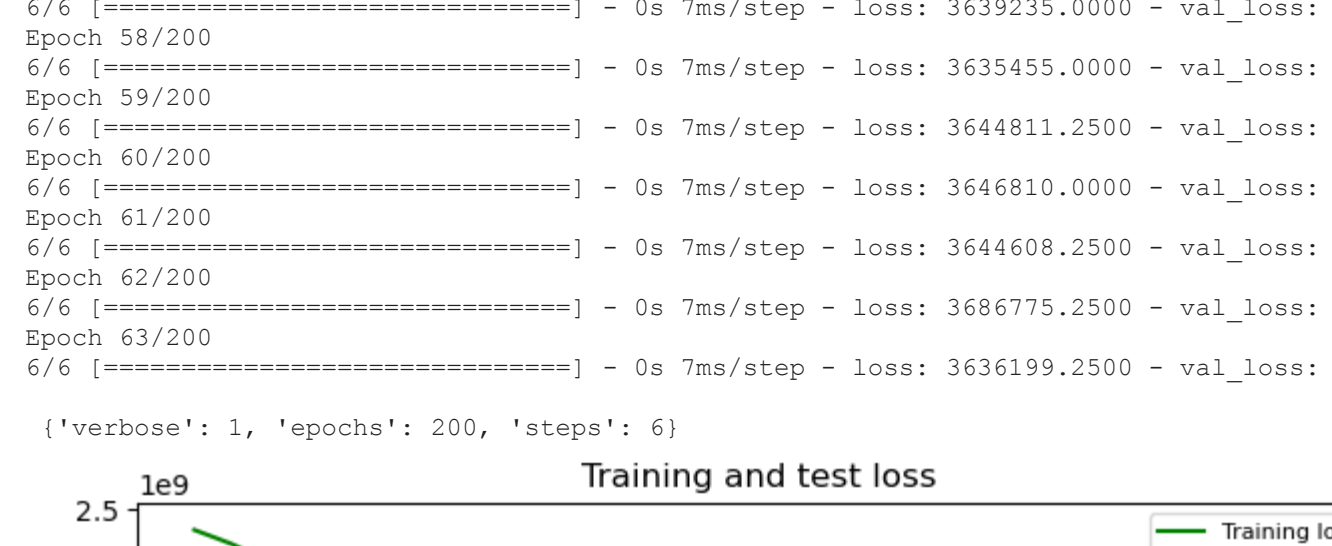
```
20/20 [.....] - 0s
Epoch 10/200
20/20 [.....] - 0s
Epoch 11/200
```

```
Epoch 12/200
20/20 [-----]
Epoch 13/200
20/20 [-----]
```

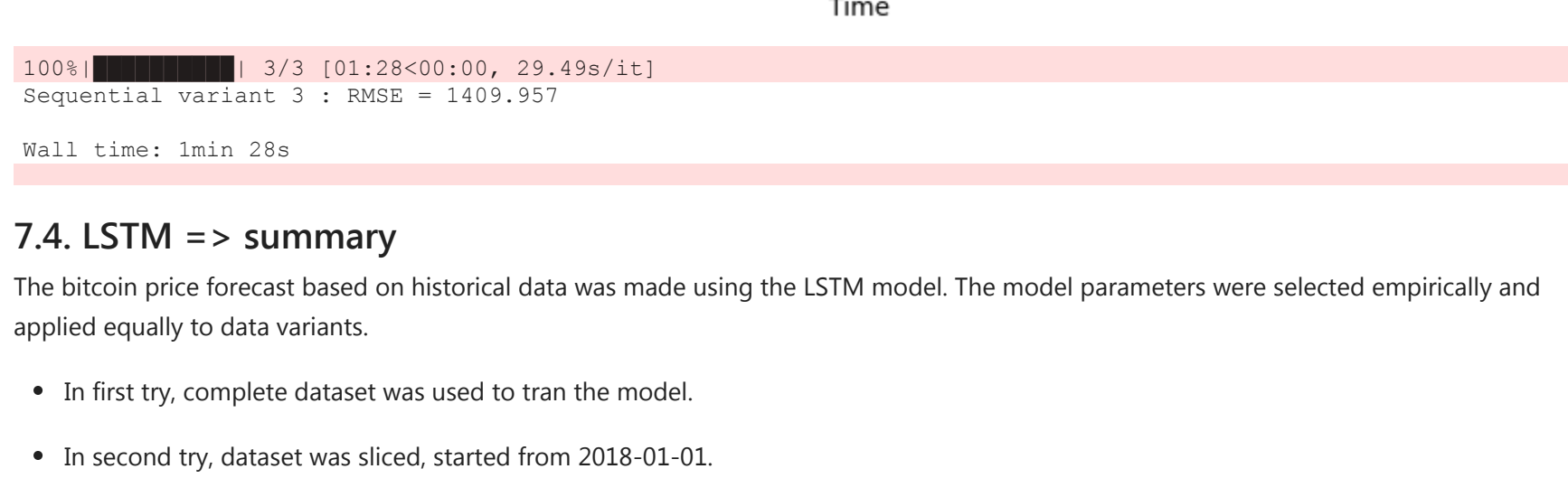
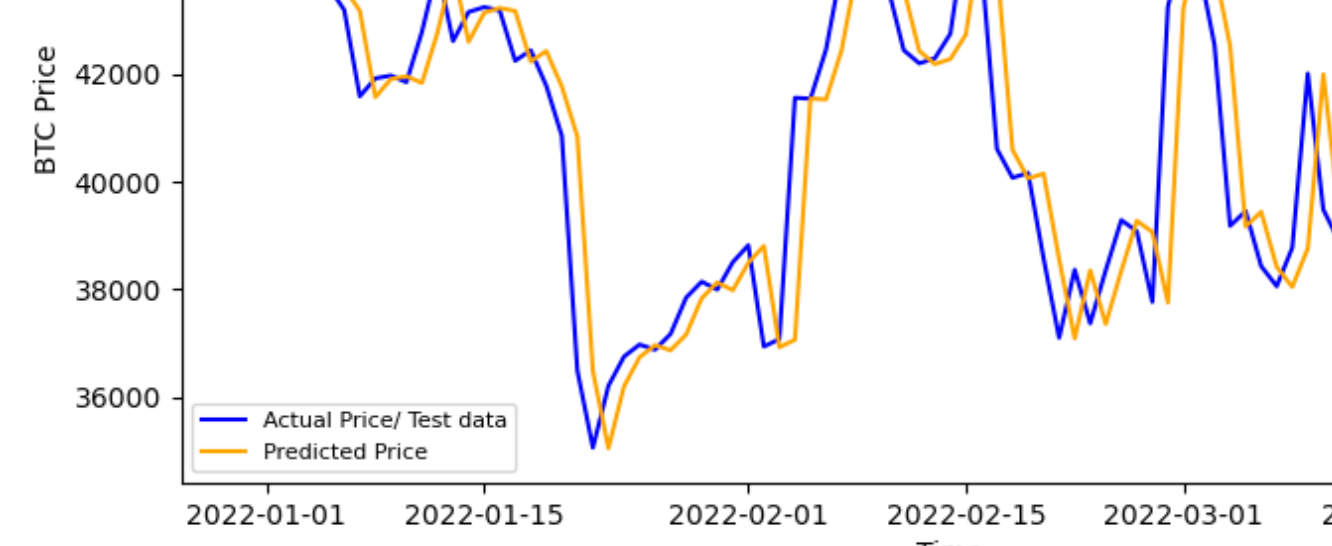
```

20/20 =====] - 0s 3ms/step - loss: 701
Epoch 14/200
20/20 =====] - 0s 3ms/step - loss: 684
Epoch 15/200
20/20 =====] - 0s 3ms/step - loss: 832
Epoch 16/200
20/20 =====] - 0s 3ms/step - loss: 691
Epoch 17/200
20/20 =====] - 0s 4ms/step - loss: 688
Epoch 18/200
20/20 =====] - 0s 3ms/step - loss: 693
Epoch 19/200
20/20 =====] - 0s 3ms/step - loss: 693
Epoch 20/200
20/20 =====] - 0s 3ms/step - loss: 688

```

[illegible]

Year	Number of people in the labor force
1990	30,000
1991	31,000
1992	32,000
1993	33,000
1994	34,000
1995	33,000
1996	32,000
1997	33,000
1998	34,000
1999	34,500
2000	34,500



- In third try, dataset was sliced, started from 2021-01

In coparison of model performence of different data range

```
# the end, to period of time around four years.

# plot for the RMSE's of data variants for LSTM

LSTM_variant = [{"LSTM_RMSE_" + str(x) for x in range(len(lstm_variants_rmse))}]
summary = {}

for variant, result in zip(LSTM_variants_rmse, lstm_variants_rmse):
    summary.append(variant + " = " + str(result))

fig = plt.figure(figsize=(10,3))
ax = fig.add_axes([0,0,1,1])
ax.bar(summary, lstm_variants_rmse)
```

```
plt.xlabel("LSTM variant")
plt.ylabel("RMSE")
plt.title("RMSE of data variants for LSTM")
plt.show()
```

LSTM variant	RMSE
1	~1740
2	~1800
3	~1740

Country	Range
Canada	900
France	1000
Germany	1100

LSTM variant	RMSE
LSTM_RMSE_0	1351.9
LSTM_RMSE_1	1609.196
LSTM_RMSE_2	1409.957

Using ARIMA, XGBoost models. Each model used a the same three data ranges to generate forecasts.

The variant of the ARIMA model and LSTM model have similar approach for forecasting. Both look at close price of one past day prediction of the close price of the next day. Performance of the models, measured by RMSE is also similar in every data range.

XGBoostRegressor variant have good results but a different approach in comparison of other models. For forecasting, columns "mark-volume" has been used as well. The model is using GridSearchCV and "eval\_set" on the fit () method and L2 penalty specified by the project.

In conclusion – the LSTM was the best model used to predict Bitcoin's price. This model gave relatively acceptable results with a [mean\\_absolute\\_error](#) of 0.0049, which is the best result.