

Configuration Manual

MSc Research Project
Data Analytics

Ram Abhilash Vasamsetti
Student ID: x22117491

School of Computing
National College of Ireland

Supervisor: Athanasios Staikopoulos

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Ram Abhilash Vasamsetti
Student ID:	x22117491
Programme:	Data Analytics
Year:	2023
Module:	MSc Research Project
Supervisor:	Athanasios Staikopoulos
Submission Due Date:	14/12/2023
Project Title:	Configuration Manual
Word Count:	682
Page Count:	9

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

ALL internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	
Date:	31st January 2024

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).	<input type="checkbox"/>
Attach a Moodle submission receipt of the online project submission , to each project (including multiple copies).	<input type="checkbox"/>
You must ensure that you retain a HARD COPY of the project , both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	<input type="checkbox"/>

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Ram Abhilash Vasamsetti
x22117491

1 Introduction

This Configuration Manual describes the hardware and software requirements and necessary configurations utilized in the research project "Using Time Series Predictive Models for Early Detection of Gambling Addiction in Problem Gamblers"

The manual is divided into 4 main sections. Section 2 gives the overview of the Research project. The section 3 highlights the Hardware and Software Prerequisites. The next Section 4 elaborates the implementation requirements where necessary libraries are discussed. The section 5 gives an overview of the dataset and its import. Section 6 briefs about the models and any necessary configurations.

2 Research Overview

This research project focuses on using Time series predictive models such as ARIMA, SARIMA and LSTM to forecast the future betting patterns of Problem Gamblers, thereby detecting Gambling Addiction. The work uses K means clustering for detection of problem gamblers and forecasting models to predict their future betting patterns.

3 System Prerequisites

3.1 Hardware Prerequisites

The implementation is carried on the windows system with the following configuration :

Processor	Intel(R) Core(TM) i5-7300HQ CPU @ 2.50GHz 2.50 GHz
Installed RAM	16.0 GB (15.9 GB usable)
System type	64-bit operating system, x64-based processor

Table 1: Operating System configuration

3.2 Software Prerequisites

Jupyter Notebook from Anaconda with Python Version: 3.11.3 has been for the implementation for this Research project.

4 Implementation Requirements

The Fig. 1 presents all the python libraries used in the research project. Any libraries missing in the machine are installed via the command "pip install [library name]".

```
In [1]: 1 import pandas as pd
        2 import seaborn as sns
        3 import matplotlib.pyplot as plt
        4 from sklearn.preprocessing import LabelEncoder
        5 import numpy as np
        6 from sklearn.cluster import KMeans
        7 from sklearn.decomposition import PCA
        8 from sklearn.metrics import silhouette_score
        9 from pandas.plotting import parallel_coordinates
       10 from ipywidgets import interact, DatePicker
       11 from statsmodels.tsa.stattools import adfuller
       12 from statsmodels.tsa.statespace.sarimax import SARIMAX
       13 from sklearn.preprocessing import MinMaxScaler
       14 from sklearn.metrics import mean_squared_error, mean_absolute_error
       15 from statsmodels.tsa.arima.model import ARIMA
       16 from keras.models import Sequential
       17 from keras.layers import Dense, LSTM
       18 from keras.optimizers import Adam
       19
```

Figure 1: Required Libraries for implementation

5 Dataset

5.1 Dataset Source

The dataset have been obtained from the Transparency Project. A Harvard Medical university initiative for encouraging research of Addictions (Division on Addiction; 2021). The dataset used has been contributed to The Transparency Project via (Gray et al.; 2012) research on Responsible Gambling. The folder contains 5 files as in Table 2

Raw Datset I.Demographics_Gray_LaPlante_PAB_2012.dat
Raw Datset II.Daily aggregates_Gray_LaPlante_PAB_2012.dat
Raw Datset III.Responsible gambling details_Gray_LaPlante_PAB_2012.dat
CodeBook_for Gray_LaPlante_PAB_2012 (Variable Definitions)
AnalyticDataset_Gray_LaPlante_PAB_2012.dat (Not Used for this Research)

Table 2: Data Source Contents

Out of all the dataset source files, Only Raw Dataset I (Demographic Information), Raw Dataset II (Daily Aggregates) and Raw Dataset III (Responsible Gambling details) are imported into the dataframes. The CodeBook file contains the legend and column definitions for each dataset used. The Analytics Dataset is a comprehensive dataset collected by the primary researcher. As per the research objectives, dedicated analytics have been performed in the data gathering and transformation phase and, thus, this file has not been used.

5.2 Dataset Import

The code as indicated in Fig. 2 is used with delimiter as `\t` . All the three raw datasets are imported. The path of the datasets need to be specified as mentioned in Fig. 2 where datasets are stored in a folder called [dataset]. Please ensure the datasets are not tampered by manually copy pasting the data.

```
1 daily_agg_df = pd.read_csv('./datasets/Raw Dataset II.Daily aggregates_Gray_LaPlante_PAB_2012.dat', delimiter='\t')
2 rg_det_df = pd.read_csv('./datasets/Raw Dataset III.Responsible gambling details_Gray_LaPlante_PAB_2012.dat', delimiter='\t')
3 demog_df = pd.read_csv('./datasets/Raw Dataset I.Demographics_Gray_LaPlante_PAB_2012.dat', delimiter='\t')
4
5
```

Figure 2: Code for importing datasets to the project

5.3 Merging of Dataset

The datasets upon minor transformations on datetime, are merged to simplify the data relation. The Fig. 3 shows the final merged dataframe containing all three dataframes. The product type == '2' indicates casino games from the codebook. This parameter can be changed to identify forecasts in other games played by different user.

Merging of Datasets

```
In [13]: 1 merged_df = daily_agg_df.merge(demog_df, on='UserID', how='outer')
          2 merged_df = merged_df.merge(rg_det_df, on='UserID', how='outer')
          3
```

```
In [14]: 1 merged_df['RG_case'].value_counts()
```

```
Out[14]: 1    811570
          0    170233
          Name: RG_case, dtype: int64
```

```
In [15]: 1 merged_df
```

```
Out[15]:
```

	UserID	ProductType	Turnover	Hold	NumberOfBets	Aggregate_Date	RG_case	CountryName	LanguageName	Gender	YearofBirth	Registration_d
0	31965	1.0	15.3388	15.3388	1	2000-05-08	1	19	8	1	1971	1999-09
1	31965	1.0	34.1594	34.1594	5	2000-05-10	1	19	8	1	1971	1999-09
2	31965	1.0	24.5419	24.5419	4	2000-05-18	1	19	8	1	1971	1999-09
3	31965	1.0	2.5309	2.5309	1	2000-05-22	1	19	8	1	1971	1999-09
4	31965	1.0	15.3387	15.3387	2	2000-05-23	1	19	8	1	1971	1999-09
...

Figure 3: Code for merging datasets

5.4 Feature Selection

The features are checked for correlation and highly correlated items are dropped by picking only the low correlated columns in the K - means clustering. The Fig. 4 shows the correlated columns. More columns can be dropped to find optimal features for the project (BUYRUKOĞLU and AKBAŞ; 2022). This paper has only excluded columns which are highly correlated.

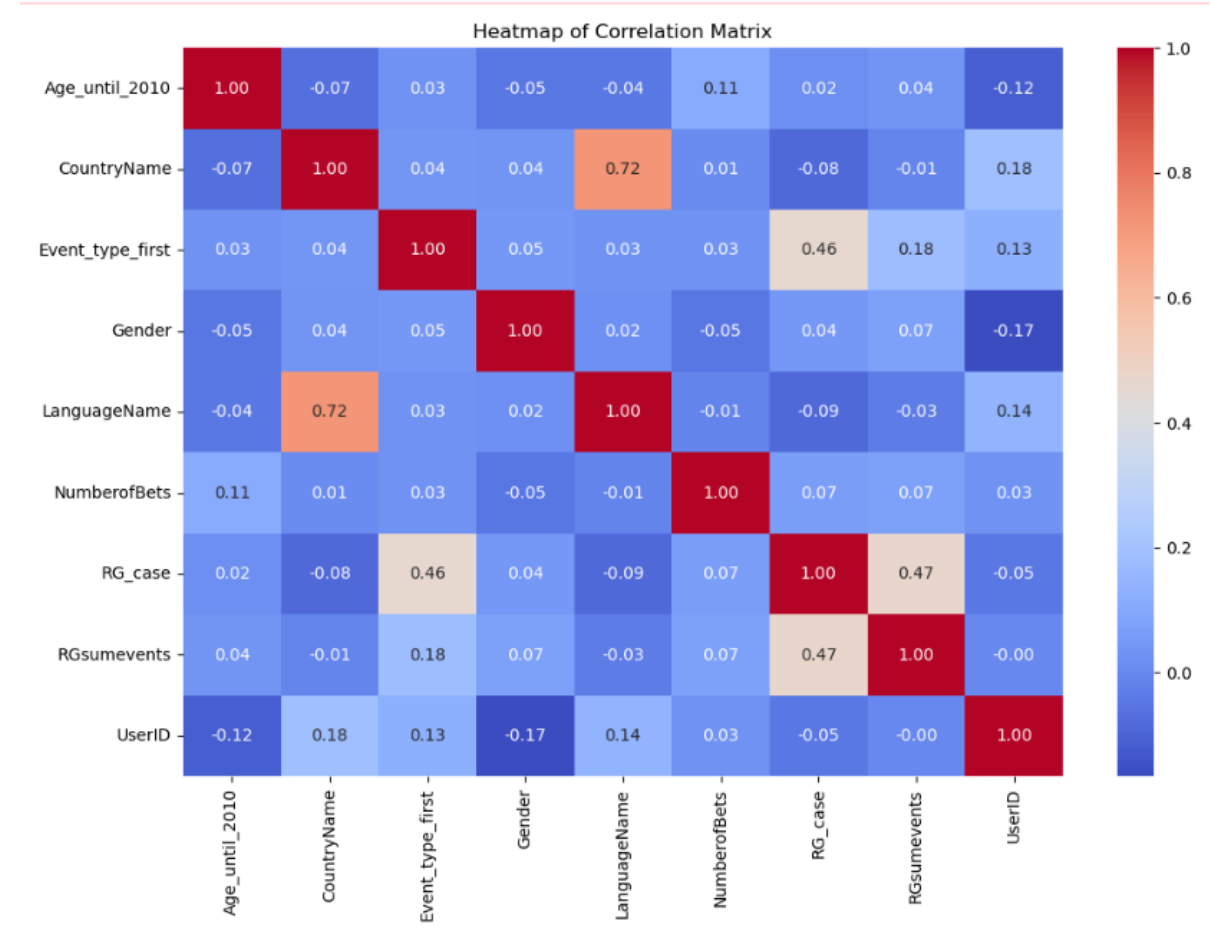


Figure 4: Heatmap for merged Dataset

6 Model Fitting

4 models have been used in total for the research project. Configurations for each model are given below.

6.1 K - Means Clustering

K - means clustering in the Fig. 5 require the data to be pivoted and features added as layers to 3d dataframe (Kobylin and Lyashenko; 2020). All the low correlated columns are passed through the [user data 3d] column. This values can be configured as per need.

K means Evaluation

```
In [24]: 1 inertias = []
2 for k_value in range(1, 6):
3     model = KMeans(n_clusters=k_value, random_state=0)
4     model.fit(user_data_array)
5     inertias.append(model.inertia_)
6
7 plt.plot(range(1, 6), inertias, marker='o')
8 plt.xlabel('Number of Clusters (K)')
9 plt.ylabel('Inertia')
10 plt.title('Elbow Method for Optimal K')
```

C:\Users\abhiv\AppData\Local\anaconda3\Lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of 'n_init' explicitly to suppress the warning
warnings.warn(
C:\Users\abhiv\AppData\Local\anaconda3\Lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of 'n_init' explicitly to suppress the warning
warnings.warn(
C:\Users\abhiv\AppData\Local\anaconda3\Lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of 'n_init' explicitly to suppress the warning
warnings.warn(
C:\Users\abhiv\AppData\Local\anaconda3\Lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of 'n_init' explicitly to suppress the warning
warnings.warn(
C:\Users\abhiv\AppData\Local\anaconda3\Lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of 'n_init' will change from 10 to 'auto' in 1.4. Set the value of 'n_init' explicitly to suppress the warning
warnings.warn(
Out[24]: Text(0.5, 1.0, 'Elbow Method for Optimal K')

Figure 5: K - Means Clustering Elbow method code

The optimal K value is set as 3 using Elbow method as shows in the Fig. 6. This can be configured as per requirement if better K value is found for different game type apart from "2".

```
7
8 # Pivot the DataFrame to create a 3D array with entries as rows, features as columns, and dates as depth
9 user_data_3d = filtered_df.pivot(index='UserID', columns='Aggregate_Date', values=['Turnover', 'Hold', 'NumberOfBets', 'Age_u
10
11 # Fill missing values with zeros (if any)
12 user_data_3d = user_data_3d.fillna(0)
13
14 # Convert empty strings ( ' ') to float values of zero (0)
15 user_data_3d = user_data_3d.replace(' ', 0).astype(float)
16 user_data_3d = user_data_3d.astype(float)
17
18 # Convert the pivot table to a NumPy array
19 user_data_array = user_data_3d.to_numpy()
20
21 #k fold value
22 k = 3
23
24 # Perform K-means clustering
25 model = KMeans(n_clusters=k, random_state=0)
26 y_pred = model.fit_predict(user_data_array)
27
28 # Apply PCA to reduce dimensionality to 2D
29 pca = PCA(n_components=2)
30 user_data_pca = pca.fit_transform(user_data_array)
31
32 # Visualize the clustered entries using PCA components
33 plt.figure(figsize=(8, 6))
34 scatter = plt.scatter(user_data_pca[:, 0], user_data_pca[:, 1], c=y_pred, cmap='viridis')
35 plt.xlabel('PCA Component 1')
36 plt.ylabel('PCA Component 2')
37
```

Figure 6: K - Means Clustering Model

Cluster Labels can be configured based on Cluster Analysis done on the mean hold and turnover as in Fig. 7.

```

In [28]: 1 # Create a scatter plot
2 plt.figure(figsize=(10, 6))
3 sns.scatterplot(data=agguser_allfeature, x='Mean_Hold', y='Mean_NumberofBets', hue='Cluster', palette='viridis', alpha=0.7)
4 plt.xlabel('Mean_Hold')
5 plt.ylabel('Mean_NumberofBets')
6 plt.title('Cluster Visualization based on Mean Hold and Mean Number of Bets')
7 #plt.legend(title='Cluster', loc='upper right')
8 legend_labels = ['Moderate Problem Gamblers', 'Early Players', 'Problem Gamblers']# [f'Cluster {info["cluster_num"]}' for info
9 legend = plt.legend(handles=scatter.legend_elements()[0], title='Cluster', labels=legend_labels)
10
11 plt.show()

```

Figure 7: K - Means Clustering Label code

After Cluster Analysis. The moderate cluster user id need to be populated in the variable shown in Fig. 8.

```

In [34]: 1 # features to consider : turnover, hold , number of bets for time series prediction

In [35]: 1 moderate_pg_players= [868583, 1175809, 1411743, 1457496, 1486136, 1662632, 1679490, 1776178, 1790848, 1921204, 2070894, 2150
2
3 count = len(moderate_pg_players)
4 print("Count of moderate addicted players:", count)
5

```

Count of moderate addicted players: 40

Figure 8: K - Means Moderate players config.

6.2 ARIMA/ SARIMA

Before Fitting into ARIMA/ SARIMA models the moderate cluster id needs to be segregated into stationary and non stationary (Franses; 1991). For this ADF test is being used. The threshold can be configured as in Fig. 9. A confidence interval of 0.05 is being used as default to reject the null hypothesis.

the [is_stationary_df] is a df which stores the results of ADF test and is used to segregate into ARIMA/ SARIMA or LSTM model data input.

Both ARIMA and SARIMA models in Fig. 10 and 11 are implemented over a loop for each user. Thus both models are coded in the same loop for efficiency. SARIMA model is implemented using SARIMAX. It is SARIMA model with exogenous factors. The ARIMA/SARIMAX parameters can be configured based on seasonality and repetitive trends identified in the manual examination of the usage plot (Kumar Dubey et al.; 2021).

6.3 LSTM

LSTM model is used at the last and is implemented for non stationray datapoints (Kumar Dubey et al.; 2021). The Hyperparameter tuning section in the Fig. 12 is commented out intentionally to save time. It can be un-commented and be used only if there is change in the dataset. The LSTM model has already undergone hyper parameter tuning and the results are commented for reference as can be seen in the Fig. 13. The best parameters obtained are directly used in the LSTM implementation over a loop to forecast for each user as seen in Fig. 14.

Filtering of data based on hypothesis testing for Stationarity

```
In [37]: 1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from statsmodels.tsa.stattools import adfuller
5
6 # Create a dictionary to store user IDs, their corresponding 'is_stationary' values, and 'hypothesis_test_result'
7 is_stationary_data = {'UserID': moderate_pg_players, 'is_stationary': [], 'hypothesis_test_result': []}
8
9 # Iterate through the List of user IDs
10 for user_id in moderate_pg_players:
11     # Extract the user's turnover data
12     user_turnover = user_data_3d['Turnover'][user_data_3d.index == user_id].values.ravel()
13
14     # Perform the ADF test for stationarity
15     result = adfuller(user_turnover)
16
17     # Check if the time series is stationary based on the p-value
18     if result[1] <= 0.05:
19         is_stationary = 1 # Stationary
20         hypothesis_test_result = 'Reject Null Hypothesis' # Stationarity is significant
21     else:
22         is_stationary = 0 # Not stationary
23         hypothesis_test_result = 'Fail to Reject Null Hypothesis' # Stationarity is not significant
24
25     # Append the 'is_stationary' and 'hypothesis_test_result' values to the List
26     is_stationary_data['is_stationary'].append(is_stationary)
27     is_stationary_data['hypothesis_test_result'].append(hypothesis_test_result)
28
29 # Create a new DataFrame from the dictionary
30 is_stationary_df = pd.DataFrame(is_stationary_data)
31
32 # Print the new DataFrame
33 print(is_stationary_df)
34
```

Figure 9: Stationarity Test

```
28 # Scale the data
29 scaler = MinMaxScaler(feature_range=(0, 1))
30 x_single_user_turnover_ravel_scaled = scaler.fit_transform(x_single_user_turnover_ravel.reshape(-1, 1))
31
32 # Split the data into training and test sets
33 split_ratio = 0.97 # 97% for training, 3% for testing
34 split_index = int(len(x_single_user_turnover_ravel) * split_ratio)
35
36 train_data = x_single_user_turnover_ravel_scaled[:split_index]
37 test_data = x_single_user_turnover_ravel_scaled[split_index:]
38
39 # Create a time index for your data
40 time_index_train = pd.date_range(start=start_date, periods=len(train_data), freq=frequency)
41
42 # Create a time index for your data
43 time_index_test = pd.date_range(start=time_index_train[-1], periods=len(test_data), freq=frequency)
44
45 #ARIMA
46
47 # Define the ARIMA model
48 model_arima = ARIMA(train_data, order=(1, 1, 1))
49
50 # Fit the ARIMA model to the training data
51 FITmodel_arima = model_arima.fit()
52
53 # Forecast the test series using ARIMA
54 FITmodel_arima_forecast = FITmodel_arima.predict(start=split_index, end=len(x_single_user_turnover_ravel) - 1)
55
56 # Inverse scale the ARIMA forecasted values
57 FITmodel_arima_forecast = scaler.inverse_transform(FITmodel_arima_forecast.reshape(-1, 1)).reshape(-1)
```

Figure 10: ARIMA Model

```

59 #SARIMA
60
61 # Define the SARIMA model with seasonal difference and order
62 model_sarima_monthly = SARIMAX(train_data, order=(1, 1, 1), seasonal_order=(1, 1, 1, 14))
63
64 # Fit the model to the training data
65 FITmodel_sarima_monthly = model_sarima_monthly.fit()
66
67 # Forecast the test series
68 FITmodel_sarima_monthly_forecast = FITmodel_sarima_monthly.forecast(steps=len(test_data))
69
70 # Inverse scale the forecasted values
71 FITmodel_sarima_monthly_forecast = scaler.inverse_transform(FITmodel_sarima_monthly_forecast.reshape(-1, 1)).reshape(-1)
72
73 # Inverse scale the training data
74 train_data_inverse = scaler.inverse_transform(train_data.reshape(-1, 1)).reshape(-1)
75
76 # Inverse scale the test data
77 test_data_inverse = scaler.inverse_transform(test_data.reshape(-1, 1)).reshape(-1)

```

Figure 11: SARIMA Model

```

57
58 # Define hyperparameters for tuning
59 units_values = [64, 128, 256]
60 learning_rate_values = [0.01, 0.001, 0.0001]
61
62 best_rmse = float('inf')
63 best_params = None
64
65 # Perform grid search
66 for units in units_values:
67     for learning_rate in learning_rate_values:
68         model = create_lstm_model(units, learning_rate)
69
70         # Train the model
71         model.fit(x_train, y_train, batch_size=1, epochs=3, verbose=0)
72
73         # Prepare the testing data
74         test_data = scaled_data[training_data_monthly_len - 7 * test_weeks:, :]
75         x_test = []
76
77         for i in range(7 * test_weeks, len(test_data)):
78             x_test.append(test_data[i - (7 * test_weeks):i, 0])
79
80         x_test = np.array(x_test)
81         x_test = np.reshape(x_test, (x_test.shape[0], x_test.shape[1], 1))
82
83         # Get predictions
84         predictions = model.predict(x_test)
85         predictions = scaler.inverse_transform(predictions)

```

Figure 12: Hyperparameter Tuning

```

In [52]: 1 # Mean: 35.68018188057653
2 # Standard Deviation: 122.09556811748713
3 # Variance: 14907.32775393194
4 # 1/1 [=====] - 1s 853ms/step
5 # Units: 64, Learning Rate: 0.01, Test RMSE: 375.2553253207926
6 # 1/1 [=====] - 1s 840ms/step
7 # Units: 64, Learning Rate: 0.001, Test RMSE: 231.66883181925292
8 # 1/1 [=====] - 1s 907ms/step
9 # Units: 64, Learning Rate: 0.0001, Test RMSE: 299.90232400037746
10 # 1/1 [=====] - 1s 902ms/step
11 # Units: 128, Learning Rate: 0.01, Test RMSE: 398.97323688410444
12 # 1/1 [=====] - 1s 879ms/step
13 # Units: 128, Learning Rate: 0.001, Test RMSE: 241.82620950011136
14 # 1/1 [=====] - 1s 854ms/step
15 # Units: 128, Learning Rate: 0.0001, Test RMSE: 279.9525872191042
16 # 1/1 [=====] - 1s 886ms/step
17 # Units: 256, Learning Rate: 0.01, Test RMSE: 398.22153453188594
18 # 1/1 [=====] - 1s 876ms/step
19 # Units: 256, Learning Rate: 0.001, Test RMSE: 289.69737162552474
20 # 1/1 [=====] - 1s 879ms/step
21 # Units: 256, Learning Rate: 0.0001, Test RMSE: 266.59826253364014
22 # Best Hyperparameters: {'units': 64, 'learning_rate': 0.001}, Best RMSE: 231.66883181925292

```

Figure 13: Results of Hyperparameter Tuning saved in comments to reduce computation time

```

85
86 # Define the best hyperparameters
87 best_units = 64
88 best_learning_rate = 0.001
89
90 # Build and compile the LSTM model with the best hyperparameters
91 best_model = Sequential()
92 best_model.add(LSTM(units=best_units, return_sequences=True, input_shape=(x_train.shape[1], 1)))
93 best_model.add(LSTM(units=best_units, return_sequences=False))
94 best_model.add(Dense(units=25))
95 best_model.add(Dense(units=1))
96
97 optimizer = Adam(learning_rate=best_learning_rate)
98 best_model.compile(optimizer=optimizer, loss='mean_squared_error')
99
100 # Train the model
101 best_model.fit(x_train, y_train, batch_size=1, epochs=1) # Adjust epochs as needed
102
103 # Prepare the testing data
104 test_data = scaled_data[training_data_monthly_len - (7*test_weeks):, :]
105 x_test = []
106
107 for i in range(7*test_weeks, len(test_data)):
108     x_test.append(test_data[i-(7*test_weeks):i, 0])
109
110 x_test = np.array(x_test)
111 x_test = np.reshape(x_test, (x_test.shape[0], x_test.shape[1], 1))
112
113 # Get predictions
114 predictions = best_model.predict(x_test)
115 predictions = scaler.inverse_transform(predictions)
116

```

Figure 14: LSTM model

References

- BUYRUKOĞLU, S. and AKBAŞ, A. (2022). Machine learning based early prediction of type 2 diabetes: A new hybrid feature selection approach using correlation matrix with heatmap and sfs, *Balkan Journal of Electrical and Computer Engineering* **10**: 110–117.
- Division on Addiction, Cambridge Health Alliance, a. H. M. S. t. h. (2021). Behavioral characteristics of internet gamblers who trigger corporate responsible gambling interventions, Medford, MA: Division on Addiction, The Transparency Project [database distributor].
- Franses, P. H. (1991). Seasonality, non-stationarity and the forecasting of monthly time series, *International Journal of Forecasting* **7**(2): 199–208.
URL: <https://www.sciencedirect.com/science/article/pii/016920709190054Y>
- Gray, H. M., LaPlante, D. A. and Shaffer, H. J. (2012). Behavioral characteristics of Internet gamblers who trigger corporate responsible gambling interventions., *Psychology of Addictive Behaviors* **26**(3): 527–535.
- Kobylin, O. and Lyashenko, V. (2020). Time series clustering based on the k-means algorithm, *Journal La Multiapp* **1**: 1–7.
- Kumar Dubey, A., Kumar, A., García-Díaz, V., Kumar Sharma, A. and Kanhaiya, K. (2021). Study and analysis of sarima and lstm in forecasting time series data, *Sustainable Energy Technologies and Assessments* **47**: 101474.
URL: <https://www.sciencedirect.com/science/article/pii/S2213138821004847>