

# Configuration Manual

MSc Research Project Artificial Intelligence

Naresh Kumar Satish Student ID: 23248441

School of Computing National College of Ireland

Supervisor: Anderson Simiscuka

#### National College of Ireland Project Submission Sheet School of Computing



Student Name:	Naresh Kumar Satish							
Student ID:	23248441							
Programme:	Artificial Intelligence							
Year:	2024-25							
Module:	MSc Research Project							
Supervisor:	Anderson Simiscuka							
Submission Due Date:	12/12/2024							
Project Title:	Configuration Manual							
Word Count:	340							
Page Count:	7							

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.

<u>ALL</u> 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:	Naresh Kumar Satish
Date:	25th January 2025

#### PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:

Attach a completed copy of this sheet to each project (including multiple copies).Attach a Moodle submission receipt of the online project submission, to<br/>each project (including multiple copies).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.

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

Naresh Kumar Satish 23248441

## 1 Introduction

An outline about the system requirements, configurations, tools used, step by step implementation, assumptions made, and the required dependencies are described in a systematic flow of the research work. The execution of the python scripts, environment utilized, and the packages used are well described in the below sections. The manual gives an overall visualization on the datasets used, data preprocessing stages, feature engineering techniques, model training and evaluation.

## 2 System Configurations and Setup

#### 2.1 Hardware Requirements

The processor and the memory utilized was 12th Generation with 10 cores, 12 Logical Processors, 1300Mhz, 8GB RAM.

#### 2.2 Software Environment

The research work was scripted using as the Python primary language in its latest version - 3.12.4, and JupyterLab is the primary IDE created from Anaconda distribution running on Windows 11 environment with 64-bit system type. Conda was used for the management of the libraries and packages and JupyterLab was accessed using Google Chrome browser.

The required packages and dependencies are shown in the below figure 1



Figure 1: Required Packages and Dependencies

## 3 Data Collection

The data was collected using four different APIs: Global Database of Events, Language, and Tone(GDELT), Mediastack, Reddit, and Kraken. The necessary data was retrieved and was downloaded in a ".csv" format file using Pandas package.

- 1. GDELT: https://github.com/gdelt retrieval.csv
- 2. Mediastack: https://github.com/mediastack retrieval.csv
- 3. Reddit: https://github.com/reddit retrieval.csv
- 4. Kraken: https://github.com/eth daily prices sept 2023 to oct 2024.csv

	А	В	С	D	E	F		А	В	C	D	E	-	F	1	В	С	D
1	Date	Open Price	High Price	Low Price	Close Pric	Volume	1	Title	SeenDate	URL	Domain	Date	1	Title		Published I	JRL	Source
2	******	1628.67	1644.5	1627.69	1636.89	3243.676	2	Cheating	S#######	https://t	ombmmagaz	i ########	2	Canâ€"t withdr	awal	*******	nttps://ww	Reddit
3	******	1636.89	1646.97	1625.11	1635.79	4577.329	з	JP Morgar	##06-10-	2024 10:00	ns insidebitc	(#######	3	When using Gro	th16 on Ethere	( <i>#######</i> # 1	nttps://ww	Reddit
4	******	1635.8	1643.37	1616.21	1629.17	5476.967	4	US Spot B	: ******	# https://t	ectechrepor		4	What is the Ethe	ereum Virtual N toking? That io	1 <i>488488888</i>	nttps://lea	Reddit
5	******	1629.17	1646.81	1608.13	1633.46	16519.12	5	Evaluating	3 #######	# https://g	gis gisuser.co	*****	5	Staking of hot's	taking? That is Ethoroum Now		https://ww	Reddit
6	*****	1633.46	1668.43	1609.06	1632.29	18376.35	6	Kava : Top	*****	# https://g	gis gisuser.co	*****	7	Re-staking FTH	Ethereuminew	*********	https://we	Reddit
7	******	1632.3	1669.36	1623.02	1647.55	13857.45	7	Difference	e#######	# https://g	gis gisuser.co	> ########	8	Dissection of an	ERC-20 Stable	********	nttps://w	Reddit
8	*******	1647.56	1656.91	1616.22	1636.2	14380.29	8	Elastos ( I		# https://	ww.tickerrepo		9	Please Help Us	Understand Sr	*******	nttps://ww	Reddit
9	*******	1030.2	1636.95	1629.73	1635.39	2551.947	9	Ethereum		# nttps://t	DIZ DIZTOC.CO		10	Ethereum Q3ã€	*2024 Highligh		nttps://ww	Reddit
10	*******	1035.30	1035.50	1500.05	1010.04	5111.704	10	NowiCon		# https://l	ora bravenew		11	I got hacked and	d stolen \$45K o	f ######## 1	nttps://ww	Reddit
12	********	1651 51	1624.44	1530.05	1502.05	20764.37	12	Magic Ede		# https://i	ns insidebito	*********	12	Need Sepolia E	"H (testnet)!	*******	nttps://ww	Reddit
14	12  ######## 1351.51 1 15/2.44 1549.53 1593.05 14544.45													Rec	ldit			
		EINI	-inanciai Data						GD	ELIData						Di	ita	
1	A		В	С	C		E		F	G	Н		1	J	K	L		Μ
1	Title	Put	lished.	Source	URL													
2	Ether	eum 202	24-10-0	america	nt https	://www.a	ameri	canban	kingnev	s.com	/2024/10	/05/ethe	ereum-	classic-etc	-market-ca	ap-reach	es-2-7	7-billio
3	Ether	eum 202	24-10-0	america	nt https	://www.a	amerio	canban	kingnev	s.com	/2024/10	/05/ethe	reum-	eth-trading	-9-7-lower	over-last	t-week	.html
4	iShar	es Etl 202	4-10-0	america	nt https	://www.a	ameri	canban	kingnev	s.com	/2024/10	/05/isha	res-et	hereum-trus	st-etf-nasc	lagetha-t	rading	-down
5	How	To Ac 202	4-10-0	pressrel	ahttps	://press	releas	enetwo	rk.com	/site/20	024/10/05	5/how-to	-acce	nt-crypto-pa	avments-o	n-vour-w	ebsite	1
6	iShar	es Eti 202	4-10-0	etfdailyn	ev https	://www.e	etfdail	vnews.	com/20	24/10/	04/ishare	s-ether	eum-tr	ust-etf-nase	lagetha-sl	ares-do	wn-0-8	-whats
7	Ether	eum 202	24-10-0	america	nh https	·//\\\\\\\	ameri	anhan	kingnew	s com	/2024/10	/03/ethe	reum.	classic-hits	-1-day-vol	ume-of-1	68-48	millio
8	Ether	eum 202	24-10-0	america	nh https	·//\\\\\\\	ameri	ranhan	kingnew	is com	/2024/10	/03/ethe	reum.	classic-etc	trading_9.	6-lower-	over-la	st-7-d
0	Ethor	oum 202	4-10-0	america	ak httpo	.//	mori	anban	kingnou	13.COM	2024/10	/02/oth	reum-	oth market	ann raad	0-100001-0		on htn
9	Ether	eum 202	4-10-0	america	nunttps	.//www.a	americ	Janban	kinghev	s.com	2024/10	rusrethe	reum-	etn-market	-cap-react	185-279	10-0100	un.ntff
10	ETH:	Grays 202	24-10-0	Seeking	Al https	://seekir	ngalph	a.com/	article/	472480	08-graysc	ale-ethe	ereum-	mini-trust-r	naking-sei	nse-of-un	ider-pe	rforma
1	Ether	eum 202	24-10-0	america	nt https	://www.a	amerio	canban	kingnev	s.com	/2024/10	/03/ethe	ereum-	name-servi	ce-hits-24	-hour-tra	ding-v	olume-
12	2 Visa I	ntrod 202	24-10-0	Financia	l F https	://financ	ialpos	st.com/	pmn/bu	siness	-wire-nev	vs-relea	ses-pn	nn/visa-intro	oduces-the	e-visa-tol	cenized	l-asse
									Mediast	ack Data								

Figure 2: Datasets

## 4 Data Processing and Transformation

The data was preprocessed by undergoing a series of steps, by removing the URL, HTML tags, numeric characters, tokenization, stemming, lemmatization, removing the duplicate values, removing the abbreviation and was merged with the financial data of Ethereum. Figure 3 shows the script for data preprocessing.

```
[] df1 = df1.drop_duplicates(subset='Title')
df1 = df1.dropna(subset=['Title'])

df1['Title'] = df1['Title'].str.lower()
df1['Title'] = df1['Title'].str.replace(f'[{string.punctuation}]', '')
df1['Title'] = df1['Title'].str.replace(r'\d+', '')
df1['Title'] = df1['Title'].str.replace(r'\s+', ' ', regex=True)
df1['Title'] = df1['Title'].apply(word_tokenize)
stop_words = set(stopwords.words('english'))
df1['Title'] = df1['Title'].apply(lambda x: [word for word in x if word not in stop_words])
lemmatizer = WordNetLemmatizer()
df1['Title'] = df1['Title'].apply(lambda x: [lemmatizer.lemmatize(word) for word in x])
ethereum_news = pd.merge(ethereum_news, price_data, on='Date', how='inner')
```

Figure 3: Preprocessing and Transformation of the Dataset

## 5 Feature Engineering and Exploratory Data Analysis

The data was analyzed for its insights, and there were necessary feature engineering steps implemented for creation of new features which added more value to the dataset. The below figure 4 shows some of the main parts like the sentiment analyzers, hybrid sentiments implementations, and average sentiments in the feature engineering technique,

```
analyzer = SentimentIntensityAnalyzer()
D
     ethereum_news['VADER Sentiment'] = ethereum_news['Tokens'].apply(lambda x: analyzer.polarity_scores(x)['compound'])
[ ] ethereum_news['Price Change (%)'] = ((ethereum_news['Close Price (USD)'] - ethereum_news['Low Price (USD)']) / ethereum_news['Low Price (USD)'])*100
    def bert_sentiment(text):
C
         result = general_sentiment_classifier(text)
         sentiment result = result[0]
         score = sentiment_result['score
         label = sentiment result['label'
         # Map the star labels to numeric scores
         star_map = {
    '1 star': -2,
             '2 stars': -1,
             '3 stars': 0,
'4 stars': 1,
             '5 stars': 2
         numeric_score = star_map.get(label, 0)
         return numeric_score, score
 ethereum_news['TextBlob Sentiment'] = ethereum_news['Tokens'].apply(lambda x: TextBlob(x).sentiment.polarity)
 extreme_sentiment = ethereum_news[(ethereum_news['Average Sentiment'] > 0.8) | (ethereum_news['Average Sentiment'] < 0.2)]</pre>
      # correlation for extreme sentiment cases only
      correlation_extreme = extreme_sentiment[['Price Change (%)', 'Average Sentiment']].corr()
     print(correlation_extreme)
🜔 ethereum_news['Volume Norm'] = (ethereum_news['Volume'] - ethereum_news['Volume'].min()) / (ethereum_news['Volume'].max() - ethereum_news['Volume'].m
     # Hybrid Sentiments
     ethereum_news['Sentiment and Volume'] = ethereum_news['Average Sentiment'] * ethereum_news['Volume Norm']
     # correlation between this new feature and price change
correlation_sentiment_volume = ethereum_news[['Price Change (%)', 'Sentiment and Volume']].corr()
    print(correlation_sentiment_volume)
```

Figure 4: Feature Engineering

After performing the feature engineering on the dataset, the dataset was analyzed for its insights by visualizing some of the seasonal trends and downfall of the cryptocurrency within the short period. Some of the important exploratory data analysis which revealed the maximum insights from the data was correlation, tope 5 sources of the datasets, lagged effects of the price, rolling volatility of the price, and stationarity tests of the time series data, the below figure 5 and 6 highlights the steps mentioned

```
print("Correlation Matrix:")
    print(correlation_matrix)
    plt.figure(figsize=(10, 4))
    sns.heatmap(correlation_matrix, annot=True, cmap='Spectral', fmt=".2f", vmin=-1, vmax=1)
    plt.title("Correlation Matrix Heatmap")
    plt.show()
top_n = 5
    top_sources = difference.most_common(top_n)
    sources, counts = zip(*top_sources)
    colors = ['crimson', 'limegreen'] * (len(sources) // 2 + 1)
    plt.figure(figsize=(5, 5))
    plt.barh(sources, counts, color=colors[:len(sources)])
    plt.xlabel('Count')
    plt.ylabel('Source')
   plt.title(f'Top {top_n} News Sources for Ethereum News')
    plt.show()
```

Figure 5: Exploratory Data Analysis



Figure 6: Exploratory Data Analysis

## 6 Model Training

There were three models which was trained for a comparative analysis of ML and DL models, the implementation of Random Forest Regressor 7, XG Boost Regressor 8, and LSTM 9 are shown below:



Figure 7: Random Forest Regressor



Figure 8: XG Boost Regressor with Grid Search



Figure 9: LSTM

The models were tested on various hyperparameter tuning like adjusting the maximum depth of the trees, nodes, learning rate, epochs, optimizers, dropout rates, number of layers and the results produced are the best hyperparameters used in the model.

### 7 Evaluation

The model's were evaluated on various metrics like R-squared, Mean squared error, Root mean squared error, time complexity and visualizations of model's actual vs predicted data points.

```
print(f"Training Time: {training_time:.2f} seconds")
 print(f"Prediction Time: {prediction_time:.2f} seconds")
print(f"RMSE: {rmse}")
 print(f"MSE: {mse}")
 print(f"R<sup>2</sup> Score: {r2}")
last_embedding = "TFIDF"
 y_test_last = y_test
 y_pred_last = model.predict(X_test)
 plt.figure(figsize=(5, 5))
 plt.scatter(y_test_last, y_test_last, alpha=0.6, color='blue', label='Actual Values')
 plt.scatter(y_test_last, y_pred_last, alpha=0.6, color='orange', label='Predicted Values')
 plt.plot([y_test_last.min(), y_test_last.max()],
          [y_test_last.min(), y_test_last.max()],
          'r--', lw=2, label='Perfect Prediction Line')
 plt.title(f'Predicted vs Actual Price Change (%) using {last_embedding} and Sentiments')
 plt.xlabel('Actual Price Change (%)')
 plt.ylabel('Predicted Price Change (%)')
 plt.legend()
 plt.grid()
 plt.show()
```



XG Boost model showed high and accurate results outperforming other two models due to various reasons.

## 8 Explainable AI

Partial Dependence Plot, Individual Conditional Expectation and Shap was implemented only on the best working model among the three, and the results were visualized for better understanding and evaluation. Figure 11 represents the XAI implemented on XG Boost

```
numeric_features = list(range(X_train.shape[1]))
 D
      core_features = [numeric_features[column_names.index(f)] for f in ['Average_Sentiment', 'Sentiment_and_Volume'] if f in column_names]
      features_to_analyze = core_features
      fig, ax = plt.subplots(figsize=(12, 6))
      PartialDependenceDisplay.from_estimator(
          best_model, X_train, features=features_to_analyze,
kind="average", grid_resolution=50, ax=ax
      plt.title('Partial Dependence Plot for Average Sentiment and Sentiment-and-Volume')
      plt.tight_layout()
     plt.show()
numeric_features = list(range(X_train.shape[1]))
     core_features = [numeric_features[column_names.index(f)] for f in ['Average_Sentiment', 'Sentiment_and_Volume'] if f in column_names]
     features_to_analyze = core_features
     fig, ax = plt.subplots(figsize=(12, 6))
     PartialDependenceDisplay.from_estimator(
         best_model, X_train, features=features_to_analyze,
         kind="individual", grid_resolution=50, ax=ax
     plt.title('ICE Plot for Average Sentiment and Sentiment-and-Volume')
     plt.tight_layout()
     plt.show()
[ ] explainer = shap.TreeExplainer(best_model)
     shap_values = explainer.shap_values(X_test)
     sentiment_volume_idx = X_test.columns.get_loc('Sentiment_and_Volume')
     sentiment_volume_shap = shap_values[:, sentiment_volume_idx]
high_sentiment_volume = X_test['Sentiment_and_Volume'] > X_test['Sentiment_and_Volume'].quantile(0.75)
     high_sentiment_shap = sentiment_volume_shap[high_sentiment_volume]
     positive_contributions = (high_sentiment_shap > 0).sum()
     total_contributions = len(high_sentiment_shap)
     print(f"Positive SHAP contributions: {positive_contributions} out of {total_contributions}")
\rightarrow Positive SHAP contributions: 39 out of 40
explainer = shap.TreeExplainer(best_model)
     shap values interaction = explainer.shap interaction values(X test)
     shap.summary_plot(shap_values_interaction, X_test, feature_names=X_test.columns)
explainer = shap.TreeExplainer(best_model)
     shap_values = explainer.shap_values(X_test)
     X_test_subset = X_test.iloc[:, :6]
     shap_values_subset = shap_values[:, :6]
     shap.summary_plot(shap_values_subset, X_test_subset, feature_names=X_test_subset.columns)
```

Figure 11: XAI implementation