

Configuration Manual

MSc Research Project DATA ANALYTICS (MSCDAD_JAN23C)

SHREEDHAR SHARMA Student ID: 22142835

School of Computing National College of Ireland

Supervisor: DR. ARGHIR NICOLAE MOLDOVAN

National College of Ireland

MSc Project Submission Sheet



School of Computing

Student Name: SHREEDHAR SHARMA

Student ID: 22142835

Programme: MSc DATA ANALYTICS

Module: MSc RESEARCH PROJECT

Lecturer: DR. ARGHIR NICOLAE MOLDOVAN

Submission Due

Date: 31.01.2024

Project Title: DISASTER TWEETS ENSEMBLE: ENSEMBLE FOR DISASTER TWEETS CLASSIFICATION

Word Count: 1055

Page Count: 14

Year: 2023-24

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: SHREEDHAR SHARMA

Date: 30.01.2024

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 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

SHREEDHAR SHARMA 22142835

1 Introduction

This document discusses the system requirements; hardware and software required for the implementation of the project. It also showcases the steps performed in the successful implementation of the research project," Disaster Tweets Ensemble: Ensemble for Disaster Tweets Classification".

2 System Configuration

The minimum system requirements for running the discussed project are specified as follows:

2.1 Hardware Requirement

- Processor: i5 or equivalent
- RAM: Minimum 8GB
- Storage: Minimum 5GB

2.2 Software Requirement

- Jupyter Notebook: Open source for interactive documents with code, widely used in data analysis and machine learning.
- Google Colab: cloud-based Python platform for collaborative machine learning tasks, eliminating the need for extensive local resources.

3 Environment Setup

In this step, we will install the necessary libraries using 'pip install' and import the necessary modules in Python for the successful execution of our project. In our project, we worked on three datasets Queensland, Nepal and Crisis and implemented four models Logistic Regression, Naïve Bayes, XGBoost and LSTM with three techniques CountVectorizer, Word2Vec and TF-IDF on all these datasets to check the best accuracy of these models on the dataset so that we can use the embedded approach to combine the models and enhance the accuracy. We are going ahead with the Queensland Dataset, considering the common steps in all three datasets (Queensland, Nepal and Crisis).

Fig 1. Installation of libraries and modules

4 Exploratory Data Analysis

4.1 Data Loading

In this step, we will be using the Pandas library to read the CSV (Comma-separated values). As we execute this command the data from queensland.csv can be accessed, analyzed and manipulated using the variable queensland. The same method we will be following for loading the CSV files of Nepal and Crisis in the variable Nepal and crisis for loading and working with respective datasets.

In [3]: # Record the start time
start_time = time.time()
queensland = pd.read_csv('queensland.csv')

Fig 2. Data Loading

4.2 Data Understanding

This step will used to view the first 5 rows of the dataset, this command is particularly used for getting the column names and the type of data stored in these columns. We are using the same command on the other datasets.

In [4]:	qu	ueensland.head()		
Out[4]:		tweet_id	text	label
	0	295541465013182465	I just though about the night I went clubbing	not_relevant
	1	295485717465935873	Looks like its going to be another long night	not_relevant
	2	297292225811972097	@LaniiBanani hahahaha I just told him id have	not_relevant
	3	297199686392094720	Off to meeting with so called Baaps of	not_relevant
	4	295527618449666049	Doubt I'll be getting much sleep tonight	not_relevant

Fig 3. Data Understanding

4.3 Finding Missing and Duplicate values

In this step, we will be analysing the data for the missing values so that we have properly balanced data for further computations. We are using the msno.bar(queensland) to visualise the missing values in the queensland dataset and then we will check for duplicate values and delete them.



Fig 5. Checking and deleting the duplicate values

4.4 Data Subsetting

Data subsetting is the process of choosing the specific data from the dataset. In the queensland dataset, we are choosing columns 'text' and 'label' for further analysis. We are dealing with the labelled data having tweets and the information whether they are relevant or not relevant to the disaster.





5 Natural language Understanding (NLU)

The following steps are performed on the datasets to do the case standardization, removal of punctuation and other unnecessary values, additional text analysis to find the composition of stop words and then removal of stopwords and finally lemmatization (reducing the word to base form) and stemming (trimming word to root form). These are the vital steps to prepare the data and to make it uniform and understandable for processing and analysis.

```
In [12]: queensland['text'] = queensland['text'].astype(str)
queensland['text'] = queensland['text'].apply(lambda x:x.lower())
```



```
In [13]: import re
def pre_process_text(x):
    x = x.strip()
    x = re.sub(r' +', ' ', x)
    x = re.sub(r"[-()\"#/@;:{}'+=~|.!?,'0-9]", "", x)
    return(x)
In [14]: queensland['text'] = queensland['text'].apply(pre_process_text)
In [15]: def clean_text(x):
    for i in x.split(" "):
        return "".join(re.sub(r"[^a-zA-Z ]","",x))
    queensland['text'] = queensland['text'].apply(clean_text)
In [16]: queensland['label'] = queensland['label'].map({'relevant':0,'not_relevant':1})
    queensland
```

Fig 8. Removal of Punctuation and other values from Queensland Dataset

```
In [17]: # Additional Text Analysis
queensland['numeric'] = queensland['text'].apply(lambda x: len([x for x in x.split() if x.isdigit()]))
queensland['alphanumeric'] = queensland['text'].apply(lambda x: len([x for x in x.split() if x.isalnum()]))
queensland['alphabetics'] = queensland['text'].apply(lambda x: len([x for x in x.split() if x.isalpha()]))
queensland["Stopwords"] = queensland['text'].apply(lambda x: len([w for w in str(x).lower().split() if w in eng_stopwords]))
queensland["Stopwords"] = queensland['text'].apply(lambda x: len([w for w in str(x).lower().split() if w in eng_stopwords]))
                                                                                                                                                  text label numeric alphanumeric alphabetics Stopwords
 Out[17]:
                                       0 i just though about the night i went clubbing ...
                                                                                                                                                                    1
                                                                                                                                                                                              0
                                                                                                                                                                                                                                  14
                                                                                                                                                                                                                                                                14
                                                                                                                                                                                                                                                                                               8
```

1	looks like its going to be another long night	1	0	14	14	4
2	laniibanani hahahaha i just told him id have t	1	0	30	30	18
3	off to meeting with so called baaps of mineral	1	0	10	10	5
4	doubt ill be getting much sleep tonight	1	0	7	7	1
10022	is flood seriously that trending rt channelnew	0	0	13	13	2
10028	grafton queensland flood peaks at meters torr	0	0	15	15	3
10029	helicopters deployed to rescue flood victims i	0	0	16	16	3
10031	skeletonunicorn the waters upto the door my ol	0	0	20	20	11
10032	queenslands flood crisis deepens as death toll	0	0	16	16	3

⁹⁰⁴⁴ rows × 6 columns

Fig 9. Additional text analysis performed on the Queensland Dataset

```
In [18]: plt.figure(figsize=(12, 6))
sns.kdeplot(queensland[queensland['label'] == 0]['Stopwords'], label='Relevant Tweet')
sns.kdeplot(queensland[queensland['label'] == 1]['Stopwords'], label='Not Relevant Tweet', alpha=0.5, shade=True, color='#ea4335')
plt.title('Stopwords Distribution')
plt.tigend(title='Tweet')
                     plt.show()
```



Fig 10. Visualization of Stpwords on the Queensland Dataset



Fig 11. Removal of Stopwords from the Queensland Dataset

```
In [20]: def lemmatize_stemming(text):
             stemmer=SnowballStemmer('english')
              return stemmer.stem(WordNetLemmatizer().lemmatize(text, pos='v'))
In [21]: import nltk
         nltk.download('omw-1.4')
         queensland['text'] = queensland['text'].apply(lambda x:lemmatize_stemming(x))
```



6 Data Preparation

The process of organizing and structuring the data suitable for feeding to the machine learning model. The train-test split is useful to evaluate the model performance as the model was trained on the training set of data and then evaluated on the test set of data to understand how well the model behaves to new and unseen data. On the other hand data representation converts the data into the format that the machine learning algorithm can learn from.

```
Train-Test Split
In [24]: x_queens = queensland['text'
             y_queens = queensland['label']
In [25]: train_x_queens, test_x_queens, train_y_queens, test_y_queens = train_test_split(x_queens,y_queens,random_state=101,test_size=0.3)
                                                        Fig 13. Train-Test split of the Queensland Dataset
            Representations Queensland Dataset
 In [26]: count_vect = CountVectorizer(analyzer='word', token_pattern=r'\w{1,}', ngram_range=(1, 2))
            xtrain_count_queens = count_vect.fit_transform(train_x_queens)
           tfidf_transformer = TfidfTransformer()
X_traIn_tfidf_queens = tfidf_transformer.fit_transform(xtrain_count_queens)
xtest_count_queens = count_vect.transform(test_x_queens)
            X test tfidf queens = tfidf transformer.transform(xtest count queens)
 In [27]: import gensim
            queensland_model = gensim.models.Word2Vec(sentences=queensland['text'], vector_size=100, window=5, min_count=5, workers=4)
queensland_word_vectors = queensland_model.wv
            def queensland_word2vec(text):
               queensland_word.
vecs = []
for word in text.split():
    if word in queensland_word_vectors:
        vecs.append(queensland_word_vectors[word])
    if len(vecs) > 0:
        return np.mean(vecs, axis=0)
    }
}

                else:
return np.zeros(100)
            gueensland['text vecs'] = gueensland['text'].apply(gueensland word2vec)
 In [28]: queensland features = np.array(list(map(np.array, queensland['text vecs'])))
 In [29]: queensland_train_x, queensland_test_x, queensland_train_y, queensland_test_y = train_test_split(queensland_features, queensland['label'], test_size=0.3, random_state=42)
```

Fig 14. Representation of the Queensland Dataset

7 Implemented Dataset Models

The following machine learning models (Logistic Regression, XGBoost, Naïve Bayes and LSTM) with vectors (CountVectorizer, Word2Vec and TF-IDF) were implemented on all three datasets (Nepal, Crisis and Queensland) to find out the best accuracy so that the same can be used in developing the ensemble for the project which can perform better. The following is the implementation code for the Queensland dataset models. Further, we will be comparing the accuracy and ROC-AUC Score of all the models on all the datasets.

1.1 Logistic Regression With CountVectorizer

```
In [30]: logreg = LogisticRegression()
logreg.fit(xtrain_count_queens, train_y_queens)
Out[30]: LogisticRegression()
In [31]: pred_lg_cnt_queens = logreg.predict(xtest_count_queens)
queens_logit_cnt = accuracy_score(test_y_queens, pred_lg_cnt_queens).round(3)
print("Accuracy:",accuracy_score(test_y_queens, pred_lg_cnt_queens).round(3))
queens_logit_cnt_auc = roc_auc_score(test_y_queens, pred_lg_cnt_queens).round(3)
print("ROC-AUC Score: ",queens_logit_cnt_auc)
print(classification_report(test_y_queens, pred_lg_cnt_queens,digits=3))
```

Fig 15. Logistic Regression with CountVectorizer on Queensland Dataset

```
1.2 Logistic Regression With TF-IDF
```

```
In [33]: logreg = LogisticRegression()
logreg.fit(X_train_tfidf_queens, train_y_queens)
Out[33]: LogisticRegression()
In [34]: pred_lg_tfidf_queens = logreg.predict(X_test_tfidf_queens)
queens_logit_tfidf = accuracy_score(test_y_queens, pred_lg_tfidf_queens).round(3)
print("Accuracy:",accuracy_score(test_y_queens, pred_lg_tfidf_queens).round(3))
queens_logit_tfidf_auc = roc_auc_score(test_y_queens, pred_lg_tfidf_queens).round(3)
print("ROC-AUC Score: ",queens_logit_tfidf_auc)
print(classification_report(test_y_queens, pred_lg_tfidf_queens, digits=3))
```

Fig 16. Logistic Regression with TF-IDF on Queensland Dataset

1.3 Logistic Regression with Word2Vec

```
In [36]: queens_lr_model = LogisticRegression(max_iter=1000)
queens_lr_model.fit(queensland_train_x, queensland_train_y)
queens_lr_pred_y = queens_lr_model.predict(queensland_test_x)
acc_queens_word2vec_logit = accuracy_score(queensland_test_y, queens_lr_pred_y).round(3)
print("Accuracy:", accuracy_score(queensland_test_y, queens_lr_pred_y).round(3))
queens_logit_word2vec_auc = roc_auc_score(queensland_test_y, queens_lr_pred_y).round(3)
print("ROC-AUC Score: ",queens_logit_word2vec_auc)
print(classification_report(queensland_test_y, queens_lr_pred_y,digits=3))
```

Fig 17. Logistic Regression with Word2Vec on Queensland Dataset

2.1 Naive Bayes With CountVectorizer



Fig 18. Naïve Bayes with CountVectorizer on Queensland Dataset

2.2 Naive Bayes With TF-IDF

```
In [41]: nb = MultinomialNB()
# Fit the data
nb.fit(X_train_tfidf_queens, train_y_queens)
Out[41]: MultinomialNB()
In [42]: pred_nb_tfidf_queens = nb.predict(X_test_tfidf_queens)
queens_nb_tfidf = accuracy_score(test_y_queens, pred_nb_tfidf_queens).round(3)
print("Accuracy:",accuracy_score(test_y_queens, pred_nb_tfidf_queens).round(3))
queens_nb_tfidf_auc = roc_auc_score(test_y_queens, pred_nb_tfidf_queens).round(3)
print("ROC-AUC Score: ",queens_logit_cnt_auc)
print(classification_report(test_y_queens, pred_nb_tfidf_queens,digits=3))
```

Fig 19. Naïve Bayes with TF-IDF on Queensland Dataset

2.3 Naive Bayes with Word2Vec

```
In [44]: # Train and Evaluate Naive Bayes with Word2Vec features for Queensland Dataset
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import classification_report, accuracy_score
queens_nb_model = GaussianNB()
queens_nb_model.fit(queensland_train_x, queensland_train_y)
queens_nb_pred_y = queens_nb_model.predict(queensland_test_x)
acc_queens_nb_word2vec = accuracy_score(queensland_test_y, queens_nb_pred_y).round(3)
print("Accuracy:", accuracy_score(queensland_test_y, queens_nb_pred_y).round(3))
queens_nb_word2vec_auc = roc_auc_score(queensland_test_y, queens_nb_pred_y).round(3)
print("ROC-AUC Score: ",queens_nb_word2vec_auc)
print(classification_report(queensland_test_y, queens_nb_pred_y,digits=3))
```

Fig 20. Naïve Bayes with Word2Vec on Queensland Dataset

3.1 XgBoost With Count-Vectorizer

<pre>xgb_estimator = XGBClassifier(n_estimators=200,</pre>
[16:55:05] W&RNING: C:\huildkite-agent\huildk\huildkite-windows-cou-autoscaling-groun-i-07503ffd91cd9da33-1\vghoost-vghoost-ci-windows\scc\learner.cc:767:
Parameters: { "verbose" } are not used.
<pre>XGBClassifier(base_score=None, booster=None, callbacks=None,</pre>
<pre>pred_xgb_cnt_queens = xgb_estimator.predict(xtest_count_queens) queens_xgb_cnt = accuracy_score(test_y_queens, pred_xgb_cnt_queens).round(3) print("Accuracy:",accuracy_score(test_y_queens, pred_xgb_cnt_queens).round(3)) queens_xgb_cnt_auc = roc_auc_score(test_y_queens, pred_xgb_cnt_queens).round(3) print("ROC-AUC_Score: ",queens_logit_cnt_auc)</pre>
<pre>print(classification_report(test_y_queens, pred_xgb_cnt_queens,digits=3))</pre>

Fig 21. XGBoost with CountVectorizer on Queensland Dataset

3.2 XgBoost With TF-IDF

n [49]: xgb_est	<pre>imator_queens_tfidf = XGBClassifier(n_estimators=200,</pre>
	n_jobs=-1,
xgb_est	<pre>uerous=ri imator_queens_tfidf.fit(X_train_tfidf_queens, train_y_queens)</pre>
[16:55: Paramet	96] WARNING: C:\buildkite-agent\builds\buildkite-windows-cpu-autoscaling-group-i-07593ffd91cd9da33-1\xgboost\xgboost-ci-windows\src\learner.cc:767 ers: { "verbose" } are not used.
[49]: XGBClas	sifier(base_score=None, booster=None, callbacks=None,
	colsample_bylevel=None, colsample_bynode=None,
	enable categorical=False, eval meric=None, Feature types=None.
	gamma-None, gpu_id=None, grow_policy=None, importance_type=None,
	interaction_constraints=None, learning_rate=None, max_bin=None,
	max_cat_threshold=Wone, max_cat_to_onehot=Wone, max_data_ten=Wone max_datb=Wone, max_baysc=Wone
	may using step-intent, may using ana, may leaves-work, and the second states and the sec
	n estimators=200, n jobs=-1, num parallel tree=None,
	predictor=None, random_state=42,)
[50]: pred_xg	b_tfidf_queens = xgb_estimator_queens_tfidf.predict(X_test_tfidf_queens)
queens_	<pre>xgb_tfidf = accuracy_score(test_y_queens, pred_xgb_tfidf_queens).round(3)</pre>
queens	Accuracy: ,accuracy_score(test_yqueens, pred_xgp_triar_queens).round(3) xeb tfild auc = noc auc score(test y queens, pred xgb tfild queens).round(3)
print("	ROC-AUC Score: ",queens logit_cnt_auc)
print(c	lassification_report(test_y_queens, pred_xgb_tfidf_queens,digits=3))
	Fig 22. XGBoost with TF-IDF on Queensland Dataset
	2.2 VaPoost With word?voc
	5.5 Agboost With word2vec
In [52]:	# Train and Evaluate XGBoost with Word2Vec features for Queensland Dataset
200 200020	from xgboost import XGBClassifier
	<pre>queens xgb model = XGBClassifier(use label encoder=False, eval metric='logloss')</pre>
	queens yeb model fit(queensland train y queensland train y)
	ducens_xgb_moder.fre(ducensiana_crain_x, ducensiana_crain_y)
	<pre>queens_xgb_prea_y = queens_xgb_model.predict(queensiand_test_x)</pre>
	print("Accuracy:", accuracy score(queensland test y, queens xgb pred y).round(3))
	are guarder with word? Was - accuracy scans(guard) and tast y guards with prod y) pound(2)

acc_queens_xgb_word2vec = accuracy_score(queensland_test_y, queens_xgb_pred_y).round(3)
queens_xgb_word2vec_auc = roc_auc_score(queensland_test_y, queens_xgb_pred_y).round(3)
print("ROC-AUC_Score: ",queens_xgb_word2vec_auc)

print(classification_report(queensland_test_y, queens_xgb_pred_y,digits=3))

Fig 23. XGBoost with Word2Vec on Queensland Dataset

4.1 LSTM with Word2vec

In [54]:	<pre>queens_word2vec_lstm = Sequential() queens_word2vec_lstm.add(LSTM(64, input_shape=(queensland_train_x.shape[1], 1))) queens_word2vec_lstm.add(Dense(1, activation='sigmoid')) queens_word2vec_lstm.add(Dense(1, activation='sigmoid')) queens_word2vec_lstm.compile(Loss='binary_crossentropy', optimizer='adam', metrics=['accuracy']) queens_word2vec_lstm_history = queens_word2vec_lstm_history = queens_word2vec_lstm_history = queens_word2vec_lstm_history = queens_word2vec_lstm_history = queens_word2vec_lstm_history = ('queensland_train_x, queensland_train_y, epochs=5, batch_size=64, validation_data=(queensland_test_x, queensland_test_y)) plt.plot(queens_word2vec_lstm_history: \label='Train Accuracy') plt.title('queensland Word2vec LSTM Accuracy'], label='Validation Accuracy') plt.title('queensland Word2vec LSTM Accuracy') plt.ylabel('Accuracy') plt.ylabel('Accuracy') plt.label('Accuracy') plt.show()</pre>
In [55]:	<pre>pred_lstm_word2vec = queens_word2vec_lstm.predict(queensland_test_x) print("Accuracy:", accuracy_score(queensland_test_y, pred_lstm_word2vec.round(),normalize=True).round(3)) word2vec_lstm_acc = accuracy_score(queensland_test_y, pred_lstm_word2vec.round(),normalize=True).round(3)</pre>
In [56]	<pre>lstm_word2vec_auc = roc_auc_score(queensland_test_y, pred_lstm_word2vec).round(3) print("ROC-AUC Score: ",lstm_word2vec_auc) print(classification_report(queensland_test_y, pred_lstm_word2vec.round(),digits=3))</pre>

Fig 24. LSTM with Word2Vec on Queensland Dataset

4.2 LSTM with TF-IDF

```
In [57]: X train tfidf reshaped queens = np.reshape(X train tfidf queens.toarray(), (X train tfidf queens.shape[0], 1, X train tfidf queens.shape[1]))
              X_test_tfidf_reshaped_queens = np.reshape(X_test_tfidf_queens.toarray(), (X_test_tfidf_queens.shape[0], 1, X_test_tfidf_queens.shape[1]))
              queensland_tfidf_lstm = Sequential()
              queensland_tfidf_lstm.add(LSTM(64, input_shape=(1, X_train_tfidf_queens.shape[1])))
              queensland_tfidf_lstm.add(Dropout(0.2))
              queensland_tfidf_lstm.add(Dense(1, activation='sigmoid'))
              queensland_tfidf_lstm.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
              queensland_tfidf_lstm_history = queensland_tfidf_lstm.fit(X_train_tfidf_reshaped_queens,queensland_train_y , epochs=5, batch_size=64, validation_data=(X_test_tfidf_reshaped_queens, queensland_test y))
              plt.plot(queensland_tfidf_lstm_history.history['accuracy'], label='Train Accuracy')
              plt.plot(queensland_tfidf_lstm_history.history['val_accuracy'], label='Validation Accuracy')
              plt.title('Queensland TF-IDF LSTM Accuracy')
              plt.xlabel('Epochs')
              plt.ylabel('Accuracy')
              plt.legend()
              plt.show()
 In [58]: pred_lstm_tfidf = queensland_tfidf_lstm.predict(X_test_tfidf_reshaped_queens)
               print("Accuracy:", accuracy_score(queensland_test_y, pred_lstm_tfidf.round(),normalize=True).round(3))
               tfidf_lstm_acc = accuracy_score(queensland_test_y, pred_lstm_tfidf.round(),normalize=True).round(3)
In [59]:
               lstm tfidf_auc = roc_auc_score(queensland_test y, pred_lstm tfidf).round(3)
               print("ROC-AUC Score: ",lstm tfidf auc)
               print(classification report(queensland test y, pred lstm tfidf.round(),digits=3))
```

Fig 25. LSTM with TF-IDF on Queensland Dataset

4.3 LSTM with CountVectorizer

Fig 26. LSTM with CountVectorizer on Queensland Dataset

Comparison

In [64]:	Que acc	<pre>sens_df = {'Classifiers':['Logi 'NaiveBayes_ 'XgBoos 'Accuracy' c_queens_xgb_word2vec,queens_ls }</pre>	stic_Regr CountVect t_Word2Ve : [queens queens stm_cnt,tf	ession_CountVectorizer','Logistic_Regression_TF-IDF','Logistic_Regression_Word2Vec', orizer','NaiveBayes_TF-IDF','NiveBayes_Word2Vec','XGBoost_CountVectorizer','XGBoost_TF-IDF', c','LSTM_CountVectorizer','LSTM_TF-IDF','LSTM_Word2Vec'], _logit_cnt,queens_logit_tfidf,acc_queens_word2vec_logit, _nb_cnt,queens_nb_tfidf,acc_queens_nb_word2vec,queens_xgb_cnt,queens_xgb_tfidf, idf_lstm_acc,word2vec_lstm_acc],
In [65]:	Que Que	eens_df = pd.DataFrame(Queens_c eens_df	lf).round(3)
Out[65]:		Classifiers	Accuracy	
	0	Logistic_Regression_CountVectorizer	0.957	
	1	Logistic_Regression_TF-IDF	0.953	
	2	Logistic_Regression_Word2Vec	0.508	
	3	NaiveBayes_CountVectorizer	0.941	
	4	NaiveBayes_TF-IDF	0.934	
	5	NiveBayes_Word2Vec	0.510	
	6	XGBoost_CountVectorizer	0.958	
	7	XGBoost_TF-IDF	0.955	
	8	XgBoost_Word2Vec	0.509	
	9	LSTM_CountVectorizer	0.497	
	10	LSTM_TF-IDF	0.503	
	11	LSTM_Word2Vec	0.516	

Fig 27. Comparison of Accuracies of all models on the Queensland Dataset

Queensland Dataset





AUC Comparison

In [67]:	queer queer	ns_df_auc = {'Classifiers':['NaiveBayes 'XgBoo 'AUC': [q ns_xgb_word2vec_auc,lstm_cnt	'Logis _Count st_Wor ueens_ _auc,1	stic_Regression_CountVectorizer','Logistic_Regression_TF-IDF','Logistic_Regression_Word2Vec', :Vectorizer','NaiveBayes_TF-IDF','NiveBayes_Word2Vec','XGBoost_CountVectorizer','XGBoost_TF-IDF', 'd2Vec','LSTM_CountVectorizer','LSTM_TF-IDF','LSTM_Word2Vec'], logit_cnt_auc,queens_logit_tfidf_auc,queens_logit_word2vec_auc,queens_nb_cnt_auc,queens_nb_tfidf_auc,queens_nb_word2vec_auc,queens_xgb_cnt_auc,queens_xgb_tfidf_auc, lstm_tfidf_auc,lstm_word2vec_auc]}
In [68]:	queer queer	ns_df_auc = pd.DataFrame(que ns_df_auc	ens_df	f_auc)
Out[68]:		Classifiers	AUC	
	<mark>0</mark> L	ogistic_Regression_CountVectorizer	0.957	
	1	Logistic_Regression_TF-IDF	0.952	
	2	Logistic_Regression_Word2Vec	0.521	
	3	NaiveBayes_CountVectorizer	0.942	
	4	NaiveBayes_TF-IDF	0.934	
	5	NiveBayes_Word2Vec	0.523	
	6	XGBoost_CountVectorizer	0.958	
	7	XGBoost_TF-IDF	0.955	
	8	XgBoost_Word2Vec	0.523	
	9	LSTM_CountVectorizer	0.494	
	10	LSTM_TF-IDF	0.493	
	11	ISTM Word2Vec	0.525	





Fig 30. Graphical representation of the AUC-ROC of all models on the Queensland Dataset

Comparison



Classifiers

ZVec

"Wectorizer



IVes TF-IDF

8.

eBayles Countly ctorizer

Regri

'ession_TF-JDF

'ession Wordzvec

AUC Comparison

Accuracy

In [66]:	nep	bal_df_auc = {'Classifiers':[' 'NaiveBayes 'XgBoo 'AUC': [n ls	Logist _Count st_Wor epal_1 ne ne tm_cnt	<pre>ic_Regression_CountVectorizer', 'Logistic_Regression_TF-IDF', 'Logistic_Regression_Word2Vec', Vectorizer', 'NaiveBayes_TF-IDF', 'NiveBayes_Word2Vec', 'XGBoost_CountVectorizer', 'XGBoost_TF-IDF', d2Vec', 'LSTM_CountVectorizer', 'LSTM_Word2Vec'], ogit_cnt_auc,nepal_logit_tfidf_auc,nepal_logit_word2vec_auc, pal_nb_cnt_auc,nepal_nb_tfidf_auc,nepal_nb_word2vec_auc, pal_xgb_cnt_auc,nepal_xgb_tfidf_auc,nepal_xgb_word2vec_auc, _auc,lstm_tfidf_auc,lstm_word2vec_auc,]}</pre>
In [67]:	nep nep	bal_df_auc = pd.DataFrame(nepa bal_df_auc	l_df_a	uc)
Out[67]:		Classifiers	AUC	
	0	Logistic_Regression_CountVectorizer	0.769	
	1	Logistic_Regression_TF-IDF	0.772	
	2	Logistic_Regression_Word2Vec	0.503	
	3	NaiveBayes_CountVectorizer	0.770	
	4	NaiveBayes_TF-IDF	0.760	
	5	NiveBayes_Word2Vec	0.516	
	6	XGBoost_CountVectorizer	0.757	
	7	XGBoost_TF-IDF	0.743	
	8	XgBoost_Word2Vec	0.502	
	9	LSTM_CountVectorizer	0.490	
	10	LSTM_TF-IDF	0.846	
	11	LSTM_Word2Vec	0.519	

Fig 33. Comparison of AUC-ROC of all models on the Nepal Dataset







Comaprison

In [72]:	<pre>crisis_df = {'Classifiers':['Logistic_Regression_CountVectorizer','Logistic_Regression_TF-IDF','Logistic_Regression_Word2Vec', 'NaiveBayes_CountVectorizer','NaiveBayes_TF-IDF','NiveBayes_Word2Vec','XGBoost_CountVectorizer','XGBoost_TF-IDF', 'XgBoost Word2Vec','LSTM CountVectorizer','LSTM TF-IDF','LSTM Word2Vec'],</pre>
	'Accuracy': [crisis_logit_cnt,crisis_logit_tfidf,acc_crisis_word2vec_logit,
	<pre>crisis_nb_cnt,crisis_nb_tfidf,acc_crisis_nb_word2vec,</pre>
	crisis_xgb_cnt,crisis_xgb_tfidf,acc_crisis_xgb_word2vec,
	}
In [73]:	<pre>crisis_df = pd.DataFrame(crisis_df) crisis_df</pre>

ClassifiersAccuracy10logistic_Regression_CountVectorizer0.8611Logistic_Regression_TF-IDF0.8522Logistic_Regression_Word2Vec0.6053NaiveBayes_CountVectorizer0.8284NaiveBayes_TF-IDF0.7445NiveBayes_Word2Vec0.6006XGBoost_CountVectorizer0.8487XGBoost_TF-IDF0.8498XgBoost_Word2Vec0.8019LSTM_CountVectorizer0.81810LSTM_Word2Vec0.819
0Logistic_Regression_CountVectorizer0.8611Logistic_Regression_TF-IDF0.8522Logistic_Regression_Word2Vec0.6053NaiveBayes_CountVectorizer0.8284NaiveBayes_TF-IDF0.7745NiveBayes_Word2Vec0.6006XGBoost_CountVectorizer0.8397XGBoost_CountVectorizer0.8398XgBoost_Word2Vec0.6059LSTM_CountVectorizer0.81510LSTM_Word2Vec0.59811LSTM_Word2Vec0.604
1Logistic_Regression_TF-IDF0.8522Logistic_Regression_Word2Vec0.6053NaiveBayes_CountVectorizer0.8284NaiveBayes_TF-IDF0.7745NiveBayes_Word2Vec0.6006XGBoost_CountVectorizer0.8397XGBoost_TF-IDF0.8398XgBoost_Word2Vec0.6059LSTM_CountVectorizer0.81510LSTM_Word2Vec0.604
2Logistic_Regression_Word2Vec0.6053NaiveBayes_CountVectorizer0.8284NaiveBayes_TF-IDF0.7745NiveBayes_Word2Vec0.6006XGBoost_CountVectorizer0.8397XGBoost_TF-IDF0.8398XgBoost_Word2Vec0.6059LSTM_CountVectorizer0.81510LSTM_TF-IDF0.59811LSTM_Word2Vec0.604
3 NaiveBayes_CountVectorizer 0.828 4 NaiveBayes_TF-IDF 0.774 5 NiveBayes_Word2Vec 0.600 6 XGBoost_CountVectorizer 0.839 7 XGBoost_TF-IDF 0.839 8 XgBoost_Word2Vec 0.605 9 LSTM_CountVectorizer 0.815 10 LSTM_Word2Vec 0.604
4 NaiveBayes_TF-IDF 0.774 5 NiveBayes_Word2Vec 0.600 6 XGBoost_CountVectorizer 0.846 7 XGBoost_TF-IDF 0.839 8 XgBoost_Word2Vec 0.605 9 LSTM_CountVectorizer 0.815 10 LSTM_Word2Vec 0.604
5 NiveBayes_Word2Vec 0.600 6 XGBoost_CountVectorizer 0.846 7 XGBoost_TF-IDF 0.839 8 XgBoost_Word2Vec 0.605 9 LSTM_CountVectorizer 0.815 10 LSTM_TF-IDF 0.598 11 LSTM_Word2Vec 0.604
6 XGBoost_CountVectorizer 0.846 7 XGBoost_TF-IDF 0.839 8 XgBoost_Word2Vec 0.605 9 LSTM_CountVectorizer 0.815 10 LSTM_Word2Vec 0.604 11 LSTM_Word2Vec 0.604
7 XGBoost_TF-IDF 0.839 8 XgBoost_Word2Vec 0.605 9 LSTM_CountVectorizer 0.815 10 LSTM_TF-IDF 0.598 11 LSTM_Word2Vec 0.604
8 XgBoost_Word2Vec 0.605 9 LSTM_CountVectorizer 0.815 10 LSTM_TF-IDF 0.598 11 LSTM_Word2Vec 0.604
9 LSTM_CountVectorizer 0.815 10 LSTM_TF-IDF 0.598 11 LSTM_Word2Vec 0.604
10 LSTM_TF-IDF 0.598 11 LSTM_Word2Vec 0.604
11 LSTM_Word2Vec 0.604

Fig 35. Comparison of Accuracy of all models on the Crisis Dataset

AUC Comparison

In [75]:	<pre>crisis_df_auc = {'Classifiers':['Logistic_Regression_CountVectorizer','Logistic_Regression_TF-IDF','Logistic_Regression_Word2Vec',</pre>			
In [76]:				
Out[76]:		Classifiers	AUC	
	0	Logistic_Regression_CountVectorizer	0.857	
	1	Logistic_Regression_TF-IDF	0.837	
	2	Logistic_Regression_Word2Vec	0.510	
	3	NaiveBayes_CountVectorizer	0.802	
	4	NaiveBayes_TF-IDF	0.725	
	5	NiveBayes_Word2Vec	0.511	
	6	XGBoost_CountVectorizer	0.845	
	7	XGBoost_TF-IDF	0.840	
	8	XgBoost_Word2Vec	0.509	
	9	LSTM_CountVectorizer	0.513	
	10	LSTM_TF-IDF	0.493	
	11	LSTM Word2Vec	0.889	

Fig 36. Comparison of AUC-ROC of all models on the Crisis Dataset

Embedded Model

```
In [70]: from sklearn.ensemble import VotingClassifier
model1 = LogisticRegression()
model2 = MultinomialNB()
model3 = XGBClassifier()
model = VotingClassifier(estimators=[('lr', model1), ('nb', model2),('xgb',model3)], voting='hard')
model.fit(xtrain_count_queens, train_y_queens)
model.score(xtest_count_queens,test_y_queens).round(3)
```

Out[70]: 0.959

Fig 37. Embedded model with accuracy on Queensland Dataset

Embedded Model

```
In [69]: from sklearn.ensemble import VotingClassifier
model1 = LogisticRegression()
model2 = MultinomialNB()
model3 = XGBClassifier()
model = VotingClassifier(estimators=[('lr', model1), ('nb', model2),('xgb',model3)], voting='hard')
model.fit(xtrain_count_nepal, train_y_nepal)
model.score(xtest_count_nepal,test_y_nepal).round(3)
```

```
Out[69]: 0.779
```

Fig 38. Embedded model with accuracy on Nepal Dataset

Embedded Model

```
In [78]: from sklearn.ensemble import VotingClassifier
model1 = LogisticRegression()
model2 = MultinomialNB()
model3 = XGBClassifier()
model = VotingClassifier(estimators=[('lr', model1), ('nb', model2), ('xgb', model3)], voting='hard')
model.fit(xtrain_count_crisis, train_y_crisis)
model_score = model.score(xtest_count_crisis, test_y_crisis).round(3)
model_score
```

```
Out[78]: 0.862
```

Fig 39. Embedded model with accuracy on Crisis Dataset