

Configuration Manual

MSc Research Project
Data Analytics

Marcelo Fischer
Student ID: 20118872

School of Computing
National College of Ireland

Supervisor: Rejwanul Haque

National College of Ireland
Project Submission Sheet
School of Computing



Student Name:	Marcelo Fischer
Student ID:	20118872
Programme:	Data Analytics
Year:	2021
Module:	MSc Research Project
Supervisor:	Rejwanul Haque
Submission Due Date:	16/08/2021
Project Title:	Configuration Manual
Word Count:	688
Page Count:	36

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:	15th September 2021

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

Marcelo Fischer
20118872

1 Introduction

This configuration manual lists all hardware and software requirements to reproduce this research. The steps taken from data acquisition to model implementation are shown in this document.

2 Hardware and Software Requirements

Table 1 shows the hardware specifications used in the research. Table 2 shows the programming language used, the libraries used and their respective versions.

Table 1: Hardware Specifications.

<i>RAM</i>	32GB
<i>Processor</i>	Intel(R) Core(TM) i7-8750H CPU @ 2.20GHz
<i>OS</i>	Windows 10 and Ubuntu 20.04

Table 2: Python Libraries and Versions.

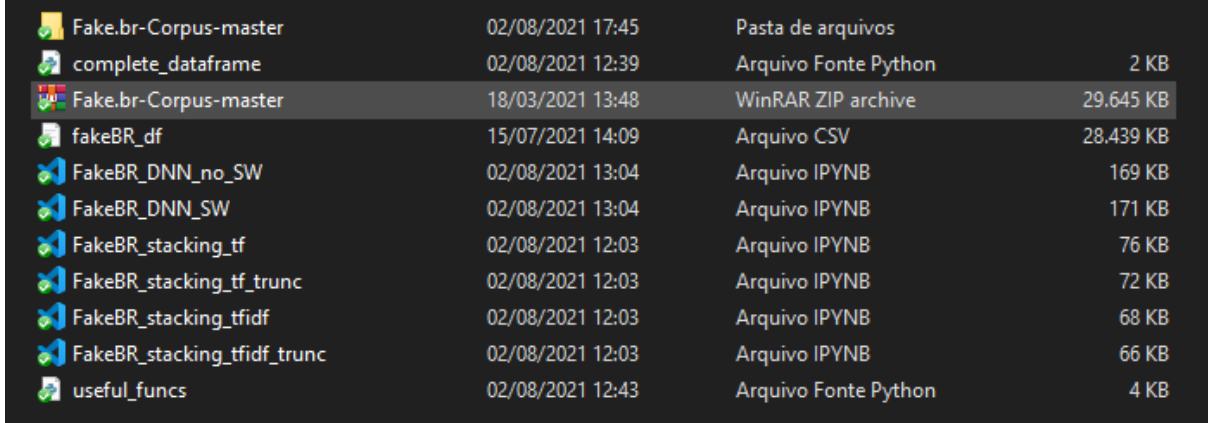
<i>Library</i>	<i>Version</i>
Python	3.8.5
Jupyter Lab	3.0.14
pandas	1.2.4
numpy	1.19.2
re	2.2.1
tensorflow	2.3.0
keras	2.4.3
scikit learn	0.24.2
nltk	3.6.2

3 Dataset

3.1 Folder Structure

The files need to be inside the *project* folder as shown in Figure 1 for the paths in the code to work correctly. Also, the paths must be changed if executing inside a Linux or

Windows machine. Each file shown in Figure 1 will be explained in this manual.



Fake.br-Corpus-master	02/08/2021 17:45	Pasta de arquivos	
complete_dataframe	02/08/2021 12:39	Arquivo Fonte Python	2 KB
Fake.br-Corpus-master	18/03/2021 13:48	WinRAR ZIP archive	29.645 KB
fakeBR_df	15/07/2021 14:09	Arquivo CSV	28.439 KB
FakeBR_DNN_no_SW	02/08/2021 13:04	Arquivo IPYNB	169 KB
FakeBR_DNN_SW	02/08/2021 13:04	Arquivo IPYNB	171 KB
FakeBR_stacking_tf	02/08/2021 12:03	Arquivo IPYNB	76 KB
FakeBR_stacking_tf_trunc	02/08/2021 12:03	Arquivo IPYNB	72 KB
FakeBR_stacking_tfidf	02/08/2021 12:03	Arquivo IPYNB	68 KB
FakeBR_stacking_tfidf_trunc	02/08/2021 12:03	Arquivo IPYNB	66 KB
useful_funcs	02/08/2021 12:43	Arquivo Fonte Python	4 KB

Figure 1: Structure of the *project* folder. All shown files must be inside this folder for the files and paths to work.

3.2 Dataset Creation

The dataset was download from a Github repository which can be found at <https://github.com/roneysco/Fake.br-Corpus>. The downloaded zip file contains three folders and a README file. The *full_texts* folder contains the full texts of the news, and also the metadata information about each label. The *size_normalized_texts* folder contains the truncated texts so that each fake-real pair has the same text length. The *preprocessed* folder contains a *.csv* file with three columns: index, label, and the pre-processed text. Pre-processed text means the removal of diacritic, accent, and Portuguese stopwords. Only the original full texts (the first folder) was used in this research.

Figure 2 shows the complete script to generate the master dataframe. For the code to work it is needed to have the *Fake.br-Corpus-master* folder inside the *project* folder. The necessary imports are shown at the top of the script.

4 Preprocessing

Before dealing with the actual data, some functions were defined to make the code cleaner and more organized. Figure 3 shows the necessary imports for the script. Figure 4 shows the function used to clean the texts. Figure 5 shows the function used to remove Portuguese stopwords from the texts. Figure 6 shows the function used to evaluate the models. Figure 7 shows the function used to save the models if wanted.

5 Experiments

5.1 Term Frequency - Full Texts Experiments

Figure 8 shows the necessary imports for these experiments. All TF experiments did not remove stopwords from the texts. Figure 9 shows how to load the data. Figure 10 shows the pre-processing steps and the creation of the train and test sets.

```

1 from collections import defaultdict
2 from pathlib import Path
3 import pandas as pd
4
5 FAKE_FOLDER_PATH = "..\\project\\Fake.br-Corpus-master\\full_texts\\fake"
6 TRUE_FOLDER_PATH = "..\\project\\Fake.br-Corpus-master\\full_texts\\true"
7 SAVE_DF_FOLDER = "..\\project\\fakeBR_df.csv"
8
9 # Create the dataframe with the fake news
10 fakes = defaultdict(list)
11 for file in Path(FAKE_FOLDER_PATH).iterdir():
12     with open(file, "r", encoding='utf-8') as f:
13         fakes["file_name"].append(file.name)
14         fakes["text"].append(f.read())
15 fake_df = pd.DataFrame(fakes)
16
17 # Create the dataframe with the real news
18 trues = defaultdict(list)
19 for file in Path(TRUE_FOLDER_PATH).iterdir():
20     with open(file, "r", encoding='utf-8') as f:
21         trues["file_name"].append(file.name)
22         trues["text"].append(f.read())
23 true_df = pd.DataFrame(trues)
24
25 # Create the labels for both dataframes
26 fake_df['label'] = 1
27 true_df['label'] = 0
28
29 # Concatenate both dataframes to create the master dataframe
30 fakeBR_df = pd.concat([fake_df, true_df])
31
32 print(f"Shape of the fake_df dataframe: {fake_df.shape}")
33 print(f"Shape of the true_df dataframe: {true_df.shape}")
34
35 print(f"Saving the dataframe to {SAVE_DF_FOLDER}...")
36 fakeBR_df.to_csv(SAVE_DF_FOLDER, index=False)

```

Figure 2: Full Python script to generate the master dataframe.

```

1 from collections import Counter
2
3 # Tools to wrangle natural language
4 import re
5 from nltk.corpus import stopwords
6
7 # Evaluation
8 from sklearn.metrics import accuracy_score, precision_score, f1_score, recall_score
9
10 from joblib import dump
11

```

Figure 3: Necessary imports for the *useful_funcs* script.


```

79 def saveModels(model, model_file_path):
80     """Save a given machine learning model to a file in the model_file_path folder.
81
82     Args:
83         model (sklearn model): a trained machine learning model to be saved for later use.
84         file_path (str): path of the folder + name of the file to save the model.
85     """
86     print(
87         f"Saving the {model} model in {model_file_path} for later use..."
88     )
89     dump(model, model_file_path)
90     print(f"The {model} model has been saved.")

```

Figure 7: Function used to save the models.

```

1  import pandas as pd
2  import numpy as np
3  from useful_funcs import cleanText, modelEval, saveModels
4
5  # Machine Learning
6  from sklearn.model_selection import train_test_split
7  from sklearn.linear_model import LogisticRegression
8  from sklearn.ensemble import RandomForestClassifier, StackingClassifier
9  from sklearn.tree import DecisionTreeClassifier
10 from sklearn.neighbors import KNeighborsClassifier
11 from sklearn.naive_bayes import MultinomialNB
12 from sklearn.svm import SVC, LinearSVC
13 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
14 from sklearn.feature_extraction.text import CountVectorizer
15 import xgboost as xgb
16
17 # Graphics
18 import matplotlib.pyplot as plt

```

Figure 8: Necessary imports for the TF experiments.

```

1  path_linux = '../project/fakeBR_df.csv'
2  path_windows = "..\\project\\fakeBR_df.csv"
3
4  fakeBR_df = pd.read_csv(path_windows)
5  fakeBR_df.head()

```

Figure 9: Loading the dataframe inside the environment.

```

1 ##### DATA PREPROCESSING FOR THE MODELS #####
2 fakeBR_df['text'] = fakeBR_df['text'].apply(cleanText)
3
4 ##### CREATE THE X AND Y SETS #####
5 x = fakeBR_df['text']
6 y = fakeBR_df['label']
7
8 ##### TF #####
9 # Create the TF Vectorizer object
10 tf_vectorizer = CountVectorizer()
11
12 # Make the sparse matrix
13 cv = tf_vectorizer.fit_transform(x)
14
15 ##### TRAIN TEST SPLIT #####
16 x_train, x_test, y_train, y_test = train_test_split(cv, y, test_size=0.2, random_state=42)
17
18 print(f"Shape of the train data: {x_train.shape}")
19 print(f"Shape of the test data: {x_test.shape}")
20 print(f"Shape of the train labels: {y_train.shape}")
21 print(f"Shape of the test labels: {y_test.shape}")

```

Figure 10: Pre-processing steps and creation of train and test sets for the full texts.

Next, all machine learning models were trained. Their code snippets are shown below. For all models, the parameters were optimized with the use of grid search and the code snippets for these are also shown. Figure 11 depicts an example of the output of the *modelEval* function.

The KNN model was not optimized by the use of grid search but by the use of the elbow method. Figure 24 shows the method.

5.2 Term Frequency - Truncated Texts Experiments

All models and the code flow were exactly the same as shown for the Section 5.1. The only difference was in the *CountVectorizer* parameter *max_features* that was set to 200, which can be seen in Figure 37.

5.3 Term Frequency-Inverse Document Frequency Experiments

All models were trained and optimized exactly like shown in Section 5.1. The only differences are in the imports (Figure 38) and in the pre-processing steps (Figure 39). The only difference when considering truncated texts is that the *max_features* parameter of the *TfidfVectorizer* is set to 200.

5.4 Neural Networks

Figure 40 shows the necessary imports for these experiments. Figure 41 shows how to load the dataframe and define the path to the folder where the models will be saved. Figure 42 shows the pre-processing for these experiments. Two different scenarios were tried, with and without stopwords. The only difference in the code from one to the other is the third line in Figure 42 which is removed when the stopwords are not removed from

Logistic Regression

```
1 lr = LogisticRegression(random_state=42, max_iter=500)
2
3 lr.fit(x_train, y_train)
4
5 lr_predictions = lr.predict(x_test)
6
7 modelEval("Logistic Regression", y_test, lr_predictions)
```

The accuracy for the Logistic Regression model is: 0.9694
The f1_score for the Logistic Regression model is: 0.9694
The recall for the Logistic Regression model is: 0.9721
The precision for the Logistic Regression model is: 0.9668

Figure 11: Standard logistic regression model.

```
1 %%time
2 lr_params = {'C': [0.01, 0.1, 1, 10, 100]}
3 lr_search = GridSearchCV(LogisticRegression(random_state=42, max_iter=500),
4                           param_grid = lr_params,
5                           n_jobs = -1,
6                           cv = 5,
7                           verbose=1)
8
9 lr_search.fit(x_train, y_train)
```

Figure 12: Grid search for the logistic regression.

```
1 lr_search.best_params_

{'C': 0.1}

1 lr_best = LogisticRegression(C=0.1, max_iter=500, random_state=42)
2
3 lr_best.fit(x_train, y_train)
4
5 lr_best_predictions = lr_best.predict(x_test)
6
7 modelEval("best Logistic Regression", y_test, lr_best_predictions)
```

Figure 13: Optimized parameters for the logistic regression.

Decision Tree ¶

```
1 dt = DecisionTreeClassifier(random_state=42)
2
3 dt.fit(x_train, y_train)
4
5 dt_predictions = dt.predict(x_test)
6
7 modelEval("Decision Tree", y_test, dt_predictions)
```

Figure 14: Standard decision tree model.

the texts. The rest of the images show the neural networks architectures, parameters and callbacks.

5.5 mBERT

The mBERT model was run using an online interface at <https://platform.peltarion.com/>. From Figure 56 to Figure 65 it is possible to see the necessary steps to reproduce the model.

```

1 %%time
2 # Measure the quality of each split
3 criterion = ['gini', 'entropy']
4 # The strategy used to choose the split at each node
5 splitter = ['best', 'random']
6 # Number of features to consider at every split
7 max_features = [None, 'auto', 'log2']
8 # Maximum number of levels in the tree
9 max_depth = [None, 2, 3, 4, 5, 6]
10 # Minimum number of samples required to split a node
11 min_samples_split = [2, 5, 7, 10]
12 # Minimum number of samples required at each leaf node
13 min_samples_leaf = [1, 2, 4]
14
15 # Create the random grid
16 dt_params = {'criterion': criterion,
17              'splitter': splitter,
18              'max_features': max_features,
19              'max_depth': max_depth,
20              'min_samples_split': min_samples_split,
21              'min_samples_leaf': min_samples_leaf}
22
23 dt_search = GridSearchCV(DecisionTreeClassifier(random_state=42),
24                          param_grid = dt_params,
25                          n_jobs = -1,
26                          cv = 5,
27                          verbose=1)
28
29 dt_search.fit(x_train, y_train)

```

Figure 15: Grid search for the decision tree.

```

1 dt_search.best_params_

{'criterion': 'entropy',
 'max_depth': 5,
 'max_features': None,
 'min_samples_leaf': 2,
 'min_samples_split': 2,
 'splitter': 'best'}

1 dt_best = DecisionTreeClassifier(splitter = 'best',
2                                   min_samples_split = 2,
3                                   min_samples_leaf = 2,
4                                   max_features = None,
5                                   max_depth = 5,
6                                   criterion = 'entropy',
7                                   random_state=42)
8
9 dt_best.fit(x_train, y_train)
10
11 dt_best_predictions = dt_best.predict(x_test)
12
13 modelEval("best Decision Tree", y_test, dt_best_predictions)

```

Figure 16: Optimized parameters for the decision tree.

Support Vector Machine

```

1 linear_svc = LinearSVC(random_state=42)
2
3 linear_svc.fit(x_train, y_train)
4
5 linear_svc_predictions = linear_svc.predict(x_test)
6
7 modelEval('LinearSVC', y_test, linear_svc_predictions)

```

Figure 17: Standard LinearSVC model.

```

1 %%time
2 # Create the regularization parameters for the SVC model.
3 C_values = [0.001, 0.003, 0.006, 0.01, 0.03, 0.06, 0.1, 0.3, 0.6]
4
5 # Create the random grid
6 linear_svc_params = {'C': C_values}
7
8 linear_svc_search = GridSearchCV(LinearSVC(random_state=42, max_iter=50000),
9                                 param_grid = linear_svc_params,
10                                n_jobs=-1,
11                                cv = 5,
12                                verbose=1)
13
14 linear_svc_search.fit(x_train, y_train)

```

Figure 18: Grid search for the LinearSVC.

```

1 linear_svc_search.best_params_

{'C': 0.006}

1 linear_svc_best = LinearSVC(C = 0.006, random_state=42)
2
3 linear_svc_best.fit(x_train, y_train)
4
5 linear_svc_best_predictions = linear_svc_best.predict(x_test)
6
7 modelEval("best LinearSVC", y_test, linear_svc_best_predictions)

```

Figure 19: Optimized parameters for the LinearSVC.

```

1 svc = SVC(random_state=42)
2
3 svc.fit(x_train, y_train)
4
5 svc_predictions = svc.predict(x_test)
6
7 modelEval('SVC', y_test, svc_predictions)

```

Figure 20: Standard SVC model.

```

1 %%time
2 # Create the random grid
3 svc_params = {'C': C_values,
4               'kernel': ['linear', 'rbf']}
5
6 svc_search = GridSearchCV(SVC(random_state=42),
7                           param_grid = svc_params,
8                           n_jobs = -1,
9                           cv = 5,
10                          verbose=1)
11
12 svc_search.fit(x_train, y_train)

```

Figure 21: Grid search for the SVC.

```

1 svc_search.best_params_

{'C': 0.01, 'kernel': 'linear'}

1 svc_best = SVC(C = 0.01, kernel = 'linear', random_state=42)
2
3 svc_best.fit(x_train, y_train)
4
5 svc_best_predictions = svc_best.predict(x_test)
6
7 modelEval("best SVC", y_test, svc_best_predictions)

```

Figure 22: Optimized parameters for the SVC.

K-Nearest Neighbours

```
1 error_rate = []
2
3 for i in range(1,40):
4     knn = KNeighborsClassifier(n_neighbors=i)
5     knn.fit(x_train, y_train)
6     pred_i = knn.predict(x_test)
7     error_rate.append(np.mean(pred_i != y_test))
```

Figure 23: Standard KNN model.

```
1 plt.figure(figsize=(10,6))
2 plt.plot(range(1,40),error_rate,color='blue', linestyle='dashed', marker='o',
3         markerfacecolor='red', markersize=10)
4 plt.title('Error Rate vs. K')
5 plt.xlabel('K')
6 plt.ylabel('Error Rate')
7 plt.plot()
```

[[

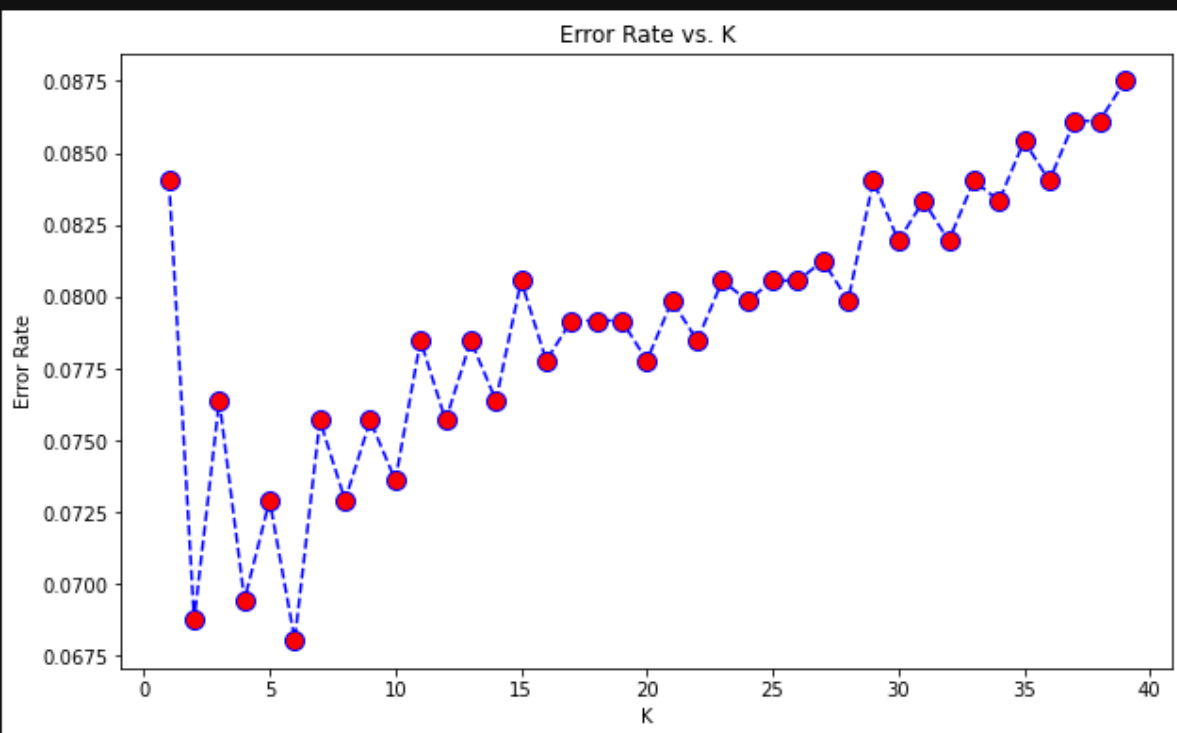


Figure 24: Elbow method for the KNN.

```
1 knn = KNeighborsClassifier(n_neighbors=6)
2 knn.fit(x_train, y_train)
3
4 knn_predictions = knn.predict(x_test)
5
6 modelEval("KNN", y_test, knn_predictions)
```

Figure 25: Optimized parameters for the KNN.

Random Forest

```
1 # Basic random forest model
2 rf1 = RandomForestClassifier(n_jobs=4, random_state=42)
3 rf1.fit(x_train, y_train)
4
5 rf1_predictions = rf1.predict(x_test)
6
7 modelEval('Random Forest', y_test, rf1_predictions)
```

Figure 26: Standard random forest model.


```

1 %%time
2 # Number of trees in random forest
3 n_estimators = list(range(10, 160, 20))
4 # Number of features to consider at every split
5 max_features = ['auto', 'sqrt']
6 # Maximum number of levels in tree
7 max_depth = [2**x for x in range(6)]
8 max_depth.append(None)
9 # Method of selecting samples for training each tree
10 bootstrap = [True, False]
11 # Create the random grid
12 rf_params = {'n_estimators': n_estimators,
13              'max_features': max_features,
14              'max_depth': max_depth,
15              'bootstrap': bootstrap}
16
17 rf_search = GridSearchCV(RandomForestClassifier(random_state=42),
18                           param_grid = rf_params,
19                           n_jobs = -1,
20                           cv = 5,
21                           verbose=1)
22
23 rf_search.fit(x_train, y_train)

```

Figure 27: Grid search for the random forest.

```

1 rf_search.best_params_

{'bootstrap': False,
 'max_depth': None,
 'max_features': 'auto',
 'n_estimators': 150}

1 rf_best = RandomForestClassifier(n_estimators = 150,
2                                 max_features = "auto",
3                                 max_depth = None,
4                                 bootstrap = False,
5                                 random_state = 42)
6 rf_best.fit(x_train, y_train)
7
8 rf_best_predictions = rf_best.predict(x_test)
9
10 modelEval('best Random Forest', y_test, rf_best_predictions)

```

Figure 28: Optimized parameters for the random forest.

Multinomial Naive Bayes

```

1 nb = MultinomialNB()
2
3 nb.fit(x_train, y_train)
4
5 nb_predictions = nb.predict(x_test)
6
7 modelEval('Multinomial Naive Bayes', y_test, nb_predictions)

```

Figure 29: Standard Naive Bayes model.

```

1 %%time
2 nb_params = {'alpha': [0.0001, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1]}
3 nb_search = GridSearchCV(MultinomialNB(),
4                           param_grid = nb_params,
5                           n_jobs = -1,
6                           cv = 5,
7                           verbose=1)
8
9 nb_search.fit(x_train, y_train)

```

Figure 30: Grid search for the Naive Bayes.

```

1 nb_search.best_params_

{'alpha': 0.1}

1 nb_best = MultinomialNB(alpha=0.1)
2
3 nb_best.fit(x_train, y_train)
4
5 nb_best_predictions = nb_best.predict(x_test)
6
7 modelEval('best Multinomial Naive Bayes', y_test, nb_best_predictions)

```

Figure 31: Optimized parameters for the Naive Bayes.

Stacking

```

1 %%time
2 estimator_lst = [
3     ("lr", lr_best),
4     ("dt", dt_best),
5     ("rf", rf_best),
6     ("linear_svc", linear_svc_best),
7     ("svc", svc_best),
8     ("knn", knn),
9     ("nb", nb_best)
10 ]
11
12 stack_model = StackingClassifier(estimators=estimator_lst,
13                                 final_estimator=LinearSVC(random_state=42, max_iter=500000),
14                                 n_jobs = -1)
15
16 stack_model.fit(x_train, y_train)

```

Figure 32: Stacking model.

```

1 stack_model_predictions = stack_model.predict(x_test)
2
3 modelEval("Stacking Model", y_test, stack_model_predictions)

```

Figure 33: Predictions for the stacking model.

XGBoost

```

1 xgboost = xgb.XGBClassifier(random_state = 42,
2                             use_label_encoder = False,
3                             n_jobs = -1)
4
5 xgboost.fit(x_train, y_train)
6
7 xgboost_predictions = xgboost.predict(x_test)
8
9 modelEval('XGBoost', y_test, xgboost_predictions)

```

Figure 34: XGBoost model.

```

1 %%time
2 # Define the grid of hyperparameters to search
3 xgb_params = {
4     'n_estimators': [100, 300, 500, 700],
5     'max_depth': [2, 3, 4, 5, 6],
6     'learning_rate': [0.01, 0.05, 0.1, 0.3],
7     'min_child_weight': [1, 2, 3, 4, 5],
8     'booster': ['gbtree', 'gblinear', 'dart']
9 }
10
11 xgb_search = RandomizedSearchCV(estimator = xgb.XGBClassifier(random_state = 42, use_label_encoder = False),
12                                param_distributions = xgb_params,
13                                cv = 5,
14                                n_iter = 50,
15                                n_jobs = -1,
16                                verbose = 3,
17                                random_state = 42)
18
19 xgb_search.fit(x_train, y_train)

```

Figure 35: Grid search for the XGBoost.

```

1 xgb_search.best_params_

{'n_estimators': 700,
 'min_child_weight': 2,
 'max_depth': 2,
 'learning_rate': 0.05,
 'booster': 'gbtree'}

1 xgboost_best = xgb.XGBClassifier(n_estimators = 700,
2                                 min_child_weight = 2,
3                                 max_depth = 2,
4                                 learning_rate = 0.05,
5                                 booster = 'gbtree',
6                                 random_state = 42,
7                                 use_label_encoder = False,
8                                 n_jobs = 6)
9
10 xgboost_best.fit(x_train, y_train)
11
12 xgboost_best_predictions = xgboost_best.predict(x_test)
13
14 modelEval("xgboost_best", y_test, xgboost_best_predictions)

```

Figure 36: Optimized parameters for the XGBoost.

```

1 ##### DATA PREPROCESSING FOR THE MODELS #####
2 fakeBR_df['text'] = fakeBR_df['text'].apply(cleanText)
3
4 ##### CREATE THE X AND Y SETS #####
5 x = fakeBR_df['text']
6 y = fakeBR_df['label']
7
8 ##### TF #####
9 # Create the TF Vectorizer object
10 tf_vectorizer = CountVectorizer(max_features=200)
11
12 # Make the sparse matrix
13 cv = tf_vectorizer.fit_transform(x)
14
15 ##### TRAIN TEST SPLIT #####
16 x_train, x_test, y_train, y_test = train_test_split(cv, y, test_size=0.2, random_state=42)
17
18 print(f"Shape of the train data: {x_train.shape}")
19 print(f"Shape of the test data: {x_test.shape}")
20 print(f"Shape of the train labels: {y_train.shape}")
21 print(f"Shape of the test labels: {y_test.shape}")

```

Figure 37: Pre-processing steps and creation of train and test sets for the truncated texts.

```

1 import pandas as pd
2 import numpy as np
3 from useful_funcs import cleanText, modelEval, saveModels
4
5 # Machine Learning
6 from sklearn.model_selection import train_test_split
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.ensemble import RandomForestClassifier, StackingClassifier
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.neighbors import KNeighborsClassifier
11 from sklearn.naive_bayes import MultinomialNB
12 from sklearn.svm import SVC, LinearSVC
13 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
14 from sklearn.feature_extraction.text import TfidfVectorizer
15 import xgboost as xgb
16
17 # Graphics
18 import matplotlib.pyplot as plt

```

Figure 38: Pre-processing steps and creation of train and test sets for the truncated texts.

```

1 ##### DATA PREPROCESSING FOR THE MODELS #####
2 fakeBR_df['text'] = fakeBR_df['text'].apply(cleanText)
3
4 ##### CREATE THE X AND Y SETS #####
5 x = fakeBR_df['text']
6 y = fakeBR_df['label']
7
8 ##### TF-IDF #####
9 # Create the TF-IDF Vectorizer object
10 tfidf_vectorizer = TfidfVectorizer()
11
12 # Make the sparse matrix
13 tfidf = tfidf_vectorizer.fit_transform(x)
14
15 ##### TRAIN TEST SPLIT #####
16 x_train, x_test, y_train, y_test = train_test_split(tfidf, y, test_size=0.2, random_state=42)
17
18 print(f"Shape of the train data: {x_train.shape}")
19 print(f"Shape of the test data: {x_test.shape}")
20 print(f"Shape of the train labels: {y_train.shape}")
21 print(f"Shape of the test labels: {y_test.shape}")

```

Figure 39: Pre-processing steps and creation of train and test sets for the full texts and the TF-IDF technique.

```

1 import pandas as pd
2 import numpy as np
3 from useful_funcs import cleanText, modelEval, saveModels, removeStopwords
4 import os
5
6 from sklearn.model_selection import train_test_split
7
8 import tensorflow as tf
9 from keras.models import Sequential, load_model
10 from keras.layers import Embedding, Dense, Dropout, GlobalMaxPooling1D, Conv1D, GRU, LSTM
11 from keras.preprocessing.text import Tokenizer
12 from keras.preprocessing.sequence import pad_sequences
13 from keras import callbacks
14 from tensorflow.keras.wrappers.scikit_learn import KerasClassifier
15
16 # Graphics
17 import matplotlib.pyplot as plt
18
19 # Needed for Linux
20 os.environ['TF_FORCE_GPU_ALLOW_GROWTH'] = 'true'

```

Figure 40: Necessary imports for the NN experiments.

```

1 MODELS_FOLDER_LINUX = "../project/models/"
2 MODELS_FOLDER_WINDOWS = "..\\project\\models\\"

```

```

1 path_linux = '../project/fakeBR_df.csv'
2 path_windows = "..\\project\\fakeBR_df.csv"
3
4 fakeBR_df = pd.read_csv(path_windows)
5 fakeBR_df.head()

```

Figure 41: The folder to save the models and how to load the dataframe.

```

1 ##### DATA PREPROCESSING FOR THE MODELS #####
2 fakeBR_df['text'] = fakeBR_df['text'].apply(cleanText)
3 fakeBR_df['text'] = fakeBR_df['text'].apply(removeStopwords)
4
5 ##### CREATE THE X AND Y SETS #####
6 x_text = fakeBR_df['text']
7 y = fakeBR_df['label']
8
9 ##### TRAIN TEST SPLIT #####
10 dnn_text_train, dnn_text_test, y_train, y_test = train_test_split(x_text, y, test_size=0.2, random_state=42)
11
12 maxlen = 300
13 embedding_dim = 100
14
15 # Tokenize words
16 tokenizer = Tokenizer()
17 tokenizer.fit_on_texts(dnn_text_train)
18 X_train = tokenizer.texts_to_sequences(dnn_text_train)
19 X_test = tokenizer.texts_to_sequences(dnn_text_test)
20
21 # Adding 1 because of reserved 0 index
22 vocab_size = len(tokenizer.word_index) + 1
23
24 # Pad sequences with zeros
25 X_train = pad_sequences(X_train, padding='post', maxlen=maxlen)
26 X_test = pad_sequences(X_test, padding='post', maxlen=maxlen)

```

Figure 42: Pre-processing steps for the NN experiments.

Convolutional Neural Network

```
1 cnn = Sequential()
2 cnn.add(Embedding(vocab_size, embedding_dim,
3                   input_length=maxlen,
4                   trainable=True))
5 cnn.add(Dropout(0.3))
6 cnn.add(Conv1D(128, 4, activation='relu'))
7 cnn.add(GlobalMaxPooling1D())
8 cnn.add(Dropout(0.3))
9 cnn.add(Dense(128, activation='relu'))
10 cnn.add(Dense(1, activation='sigmoid'))
11
12 cnn.summary()
```

Model: "sequential_8"

Layer (type)	Output Shape	Param #
embedding_8 (Embedding)	(None, 300, 100)	6810100
dropout_10 (Dropout)	(None, 300, 100)	0
conv1d_2 (Conv1D)	(None, 297, 128)	51328
global_max_pooling1d_2 (Glob	(None, 128)	0
dropout_11 (Dropout)	(None, 128)	0
dense_13 (Dense)	(None, 128)	16512
dense_14 (Dense)	(None, 1)	129
Total params: 6,878,069		
Trainable params: 6,878,069		
Non-trainable params: 0		

Figure 43: CNN architecture.

```

1 cnn.compile(optimizer='adam',
2             loss='binary_crossentropy',
3             metrics=['accuracy'])
4
5 es_cnn = callbacks.EarlyStopping(monitor='val_accuracy', patience=4, min_delta=0.0001, verbose=1)
6 mc_cnn = callbacks.ModelCheckpoint(f'{MODELS_FOLDER_WINDOWS}cnn_best.h5', monitor='val_accuracy', save_best_only=True, verbose=1)
7 cb_list_cnn = [es_cnn, mc_cnn]

```

```

1 history_cnn = cnn.fit(X_train, y_train, epochs=30, batch_size = 8, validation_split=0.3, callbacks=cb_list_cnn)

```

Figure 44: CNN callbacks and fit.

```

1 cnn_predictions = (cnn.predict(X_test) > 0.5).astype("int32")
2 cnn_predictions

```

```

array([[1],
       [1],
       [0],
       ...,
       [0],
       [1],
       [1]])

```

```

1 modelEval("cnn", y_test, cnn_predictions)

```

Figure 45: CNN predictions.

```

1 best_cnn = load_model(f'{MODELS_FOLDER_WINDOWS}cnn_best.h5')
2 best_cnn_predictions = (best_cnn.predict(X_test) > 0.5).astype("int32")
3 modelEval("best CNN", y_test, best_cnn_predictions)

```

Figure 46: Best CNN model being loaded.

```

1 acc = history_cnn.history['accuracy']
2 val_acc = history_cnn.history['val_accuracy']
3 loss = history_cnn.history['loss']
4 val_loss = history_cnn.history['val_loss']
5 epochs = range(len(acc))
6
7 fig, ax = plt.subplots(1, 2, figsize=(16,8))
8 ax[0].plot(epochs, acc, 'r', label='Training accuracy')
9 ax[0].plot(epochs, val_acc, 'b', label='Validation accuracy')
10 ax[0].set_title('Training and validation accuracy')
11 ax[0].legend(loc=0)
12
13 ax[1].plot(epochs, loss, 'r', label='Training loss')
14 ax[1].plot(epochs, val_loss, 'b', label='Validation loss')
15 ax[1].set_title('Training and validation loss')
16 ax[1].legend(loc=0)
17
18 # plt.savefig("cnn_acc.png", dpi=300, bbox_inches='tight')
19
20 plt.show()

```

Figure 47: Code snippet to plot the loss and accuracy of the CNN.

Gated Recurrent Unit

```
1 gru = Sequential()
2 gru.add(Embedding(vocab_size,
3                  embedding_dim,
4                  input_length=maxlen,
5                  trainable=False))
6 gru.add(GRU(128))
7 gru.add(Dropout(0.25))
8 gru.add(Dense(1, activation='sigmoid'))
9
10 gru.summary()
```

Model: "sequential_2"

Layer (type)	Output Shape	Param #
embedding_2 (Embedding)	(None, 300, 100)	6810100
gru_1 (GRU)	(None, 128)	88320
dropout_3 (Dropout)	(None, 128)	0
dense_3 (Dense)	(None, 1)	129

Total params: 6,898,549
Trainable params: 88,449
Non-trainable params: 6,810,100

Figure 48: GRU architecture.

```
1 gru.compile(optimizer='adam',
2             loss='binary_crossentropy',
3             metrics=['accuracy'])
4
5 es_gru = callbacks.EarlyStopping(monitor='val_accuracy', patience=4, min_delta=0.0001, verbose=1)
6 mc_gru = callbacks.ModelCheckpoint(f'{MODELS_FOLDER_WINDOWS}gru_best.h5', monitor='val_accuracy', save_best_only=True, verbose=1)
7 cb_list_gru = [es_gru, mc_gru]
8
9 history_gru = gru.fit(X_train, y_train,
10                    epochs = 20,
11                    batch_size = 8,
12                    validation_split=0.3,
13                    callbacks=cb_list_gru)
```

Figure 49: GRU callbacks and fit.

```

1 gru_predictions = (gru.predict(X_test) > 0.5).astype("int32")
2 gru_predictions

array([[1],
       [1],
       [0],
       ...,
       [0],
       [1],
       [1]])

1 modelEval("GRU", y_test, gru_predictions)

```

Figure 50: GRU predictions.

```

1 best_gru = load_model(f"{MODELS_FOLDER_WINDOWS}gru_best.h5")
2 best_gru_predictions = (best_gru.predict(X_test) > 0.5).astype("int32")
3 modelEval("best GRU", y_test, best_gru_predictions)

```

Figure 51: Best GRU model being loaded.

Long-short Term Memory

```
1 lstm = Sequential()
2 lstm.add(Embedding(vocab_size, embedding_dim,
3                   input_length=maxlen,
4                   trainable=True))
5 lstm.add(LSTM(256, return_sequences=True))
6 lstm.add(Dropout(0.3))
7 lstm.add(LSTM(64))
8 lstm.add(Dense(32, activation='relu'))
9 lstm.add(Dense(1, activation='sigmoid'))
10
11 lstm.summary()
```

Model: "sequential_6"

Layer (type)	Output Shape	Param #
=====		
embedding_6 (Embedding)	(None, 300, 100)	6810100
=====		
lstm_5 (LSTM)	(None, 300, 256)	365568
=====		
dropout_7 (Dropout)	(None, 300, 256)	0
=====		
lstm_6 (LSTM)	(None, 64)	82176
=====		
dense_9 (Dense)	(None, 32)	2080
=====		
dense_10 (Dense)	(None, 1)	33
=====		
Total params: 7,259,957		
Trainable params: 7,259,957		
Non-trainable params: 0		
=====		

Figure 52: LSTM architecture.

```

1 lstm.compile(optimizer='adam',
2               loss='binary_crossentropy',
3               metrics=['accuracy'])
4
5 es_lstm = callbacks.EarlyStopping(monitor='val_accuracy', patience=4, min_delta=0.0001, verbose=1)
6 mc_lstm = callbacks.ModelCheckpoint(f'{MODELS_FOLDER_WINDOWS}lstm_best.h5', monitor='val_accuracy', save_best_only=True, verbose=1)
7 cb_list_lstm = [es_lstm, mc_lstm]

```

```

1 history_lstm = lstm.fit(X_train, y_train,
2                          epochs=20,
3                          batch_size = 8,
4                          validation_split=0.3,
5                          callbacks=cb_list_lstm)

```

Figure 53: LSTM callbacks and fit.

```

1 lstm_predictions = (lstm.predict(X_test) > 0.5).astype("int32")
2 lstm_predictions

```

```

array([[1],
       [1],
       [0],
       ...,
       [0],
       [1],
       [1]])

```

```

1 modelEval("LSTM", y_test, lstm_predictions)

```

Figure 54: LSTM predictions.

```

1 best_lstm = load_model(f'{MODELS_FOLDER_WINDOWS}lstm_best.h5')
2 best_lstm_predictions = (best_lstm.predict(X_test) > 0.5).astype("int32")
3 modelEval("LSTM", y_test, best_lstm_predictions)

```

Figure 55: Best LSTM model being loaded.

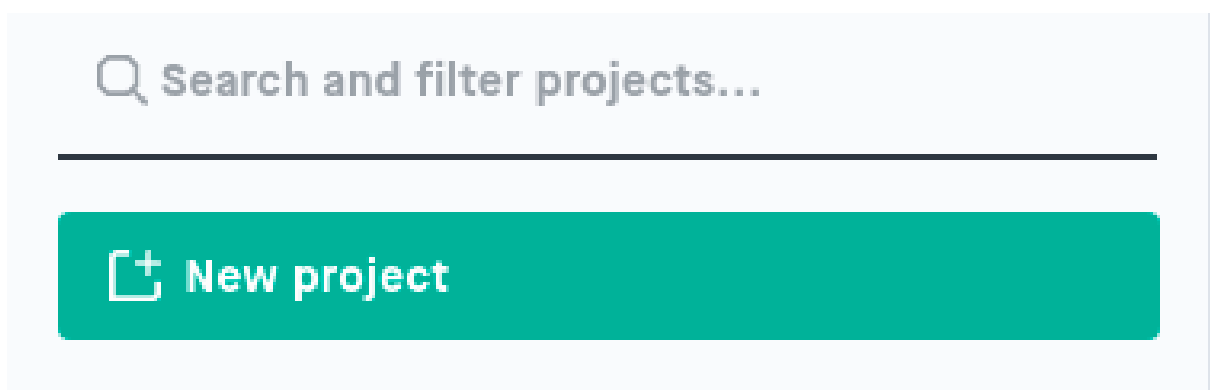


Figure 56: Create a project. It is located in the top left of the page.

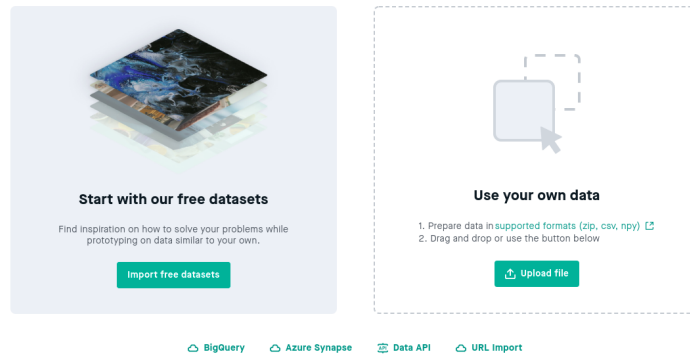


Figure 57: Upload the dataset as a csv file.

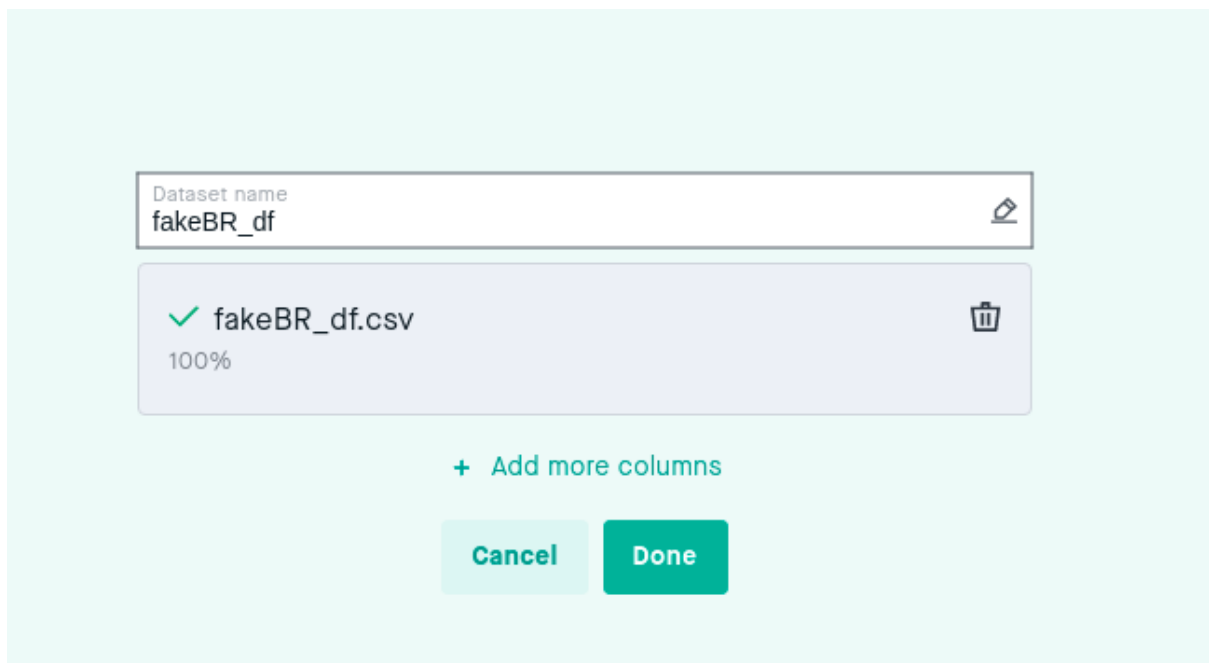


Figure 58: Click in "done".

[← All my datasets](#)

fakeBR_df

That's your dataset

Take a look in **Table view** or **Feature view** to get further details.

We've set the parameters for each feature, but they are easy to change; Click the  icon to change a feature parameter.

If you want more editing possibilities, click **Show advanced settings**.

 **Use in new experiment**

 **Show advanced settings**

Figure 59: Click in "Use in new experiment".

Experiment wizard

×

Experiment name
Experiment 1

Dataset

Inputs / target

Problem type

Select the data you want to use.

Dataset fakeBR_df	▼
Version #1	▼
Split Default split	▼

Create custom experiment

Next

Figure 60: Choose these settings for the first tab.

Experiment wizard

Experiment name
Experiment 1

Dataset	Inputs / target	Problem type
---------	-----------------	--------------

Select the input(s) and target features you want to use.
For similarity search you may not need a specific target feature if you are not tuning the model. In that case, you can pick any valid one in the list and then add an output block in the model to extract the embeddings you need.

[-]

Search feature

[Q]

☐ file_name

☒ text

☐ label

Search feature

[Q]

☐ file_name

☐ text

☒ label

Previous

Create custom experiment

Next

Figure 61: Choose these settings for the second tab.

Experiment wizard

Experiment name
Experiment 1

Dataset

Inputs / target

Problem type

Select problem type.

Problem type

Single-label text classification

Single-label text classification is when a deep learning model predicts one class for each example.

Previous

Create custom experiment

Create

Figure 62: Choose these settings for the third tab and click in "create".

Build

Settings

Dataset

Dataset

fakeBR_df

Version

#1

Split

Default split

Run settings

Batch size

32

Epochs

2

Optimizer

Adam

Learning rate

0.00002

Learning rate schedule

Triangular

Warm-up epochs

0.5

Decrement per epoch

0.000008

Early stopping

Run full number of epochs

Figure 63: The parameters used are shown in this figure.

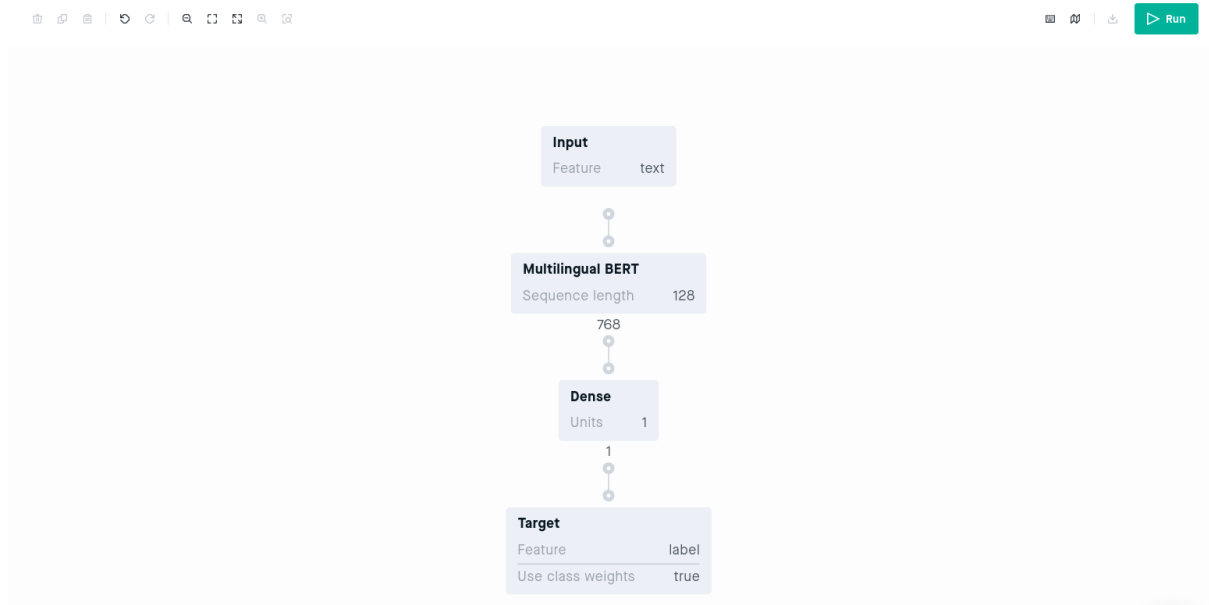


Figure 64: Click in "run" on the top right corner to run the model.



Figure 65: Go to the evaluation tab and check the results after the model ran.