

I. Definition

Sentiment analysis is a sub-field of natural language processing and employs machine learning, computational linguistics and data mining. The goal of sentiment analysis is to automatically detect the polarity of a text and try to systematically identify, extract, quantify, and study affective states and subjective information within text.

II. Project Overview

In research will build upon existed Arabic dataset from the web which were designed for detecting emotions for text, the goal is to optimize existing models toward designing a robust computational approach for analyzing and detecting emotions from Arabic dataset. Additionally, conducting research on the state of art techniques for Arabic sentimental analysis. The outcome from this research will contribute knowledge towards recognize, understand and detect accurate emotion from Arabic dataset called Arabic Sentiment Tweets Dataset (ASTD).

III. Problem Statement

This research project aims to explore the most accurate model to classify Arabic sentimental dataset. This will be look at through the evaluation for different proposed models to contribute in designing a robust model for evaluating emotions.

In our research project, we focused in ten different classification algorithms, used hyper tuning gridSearch method and resampling techniques to increase classification accuracy. We followed the following number of steps achieve our goal. First, Review the literature to evaluate various experiments and datasets that was used to extract sentiments from Arabic datasets and discuss the outcomes of experimentation by replicating selected experiments model with ASTD dataset [1]. Then, create an enhanced sentimental analysis model that is an upgraded version from the existing models that proved to gain better insights from sentimental dataset. This include a model to compare a set of reviewed experimented models and enhanced model developed after observation. This is done through tuning the parameters using gridSearch method in python and applied in both the standard architecture (replicated experiment) and the enhanced model (using different preprocessing steps and enhancing and testing the standard models before and after parameters tuning). Finally, we statistically evaluate, analyze and interpret the results obtained to assess the accuracy of sentimental analysis models using Arabic dataset.

IV. Metrics

We described the performance measurements by using numerical scalars and ROC curve to easily describe the accuracy, F-score and contingency table of the selected models.

- a. **Accuracy** is good measurement if test sets fully describe the diversity of data tuples to be classified.
- b. Precision and recall ignore true negatives, which focus more in measuring the performance. We will use the **F1 score** which can be interpreted as a weighted average of the precision and recall, where an F1 score reaches its best value at 1 and worst at 0.
- c. **Contingency table** will be constructed in the aim of measuring true positive, turn negative, false positive and false negative. Then, precision, recall and F-measure accurately analyzed to avoid selecting model that always predict negative.

V. Analysis

a. Data Exploration

The dataset has been explored first by converting it to CSV file and read it through DataFrame. The dataset consists of 10006 tweets, each categorized as OBJ, NEG, POS, NETURAL (Table 1). There are no empty data nor inconsistency between samples' classifications. While the prevalent category is OBJ with about 6675 tweets.

Table 1 Sample from ASDA Dataset

	Tweet	Classification
0	بعد استقالة رئيس المحكمة الدستورية ننتظر استق...	OBJ
1	أهني الدكتور أحمد جمال الدين، القيادي بحزب مصر...	POS
2	البرادعي يستقوى بأمريكا مره أخرى و يرسل عصام ال...	NEG
3	... الحرية والعدالة شاهد الآن: #ليلة_الاتحادية#	OBJ
4	...والدة لو اقولها بخاطري حشيشة تضحك بس من القول	NEUTRAL

To do more exploration, we investigate the existence of stop words in our dataset and found out that there is no need to remove any because there is no match between Arabic stop words in NLTK package and our dataset as such we skipped this step (code is commented in the notebook) and the fact that we were using the cleaned data from ASDA dataset. Also, to get better insight, frequency of each word calculated using CountVectorizer and 1-gram this allow us to construct a DataFrame consist of word-frequency for all the existent words in the dataset (Figure 1).

Figure 1 Word-Frequency of ASDA dataset

	frequency
دراما	17941
استنتاجات	3519
الغربة	8174
إفساد	2405
شيفنا	20447

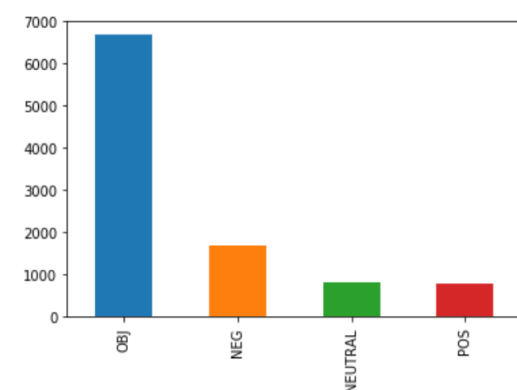
b. Exploratory Visualization

In order to explore the data, we used word cloud with word frequency to generate the most repeated words (Figure 2).

Figure 2 ASDA Dataset word frequency



The data distribution in in balanced classes is visualized using bar chart, OBJ is 6675, NEG is 1682, NETURAL is 831 and POS is 799.



With Down Sampled data (see VII. Data Preprocessing) we down sampled our data to 799 to match minority class POS. While with up sampled data we up sampled them to match the majority class OBJ to 6675.

c. Algorithms and Techniques

In our experiments, we used ten standard classifiers, eight of them were already specified in ASDA research paper [1] (read d. Benchmark section). We used token count and Frequency-inverse document frequency (Tfidf) to extract features from our corpus. We also used up streaming and down streaming techniques to balance out target classes. Additionally, we add two classifiers and they are Random Forest and Decision Trees as we wanted to evaluate their performance compared to the others. In the following subsections we will go through all of them.

1. Token Count

It takes each sentence (all the words) present in the data set in the review section and then splits each of the words present in the form of tokens. The occurrence of these tokens in the whole data set are counted in such a way that the count of the occurrence of each token in a positive, negative, neutral or objective are collected separately. Finally, the word frequency of the tokens is calculated. It is a medium to get statistical inference from the data to make a numerical analysis.

2. Frequency-inverse document frequency

Tfidf is another way to convert corpus data to statistical form to get better insight from the textual dataset. In such a way, it will calculate the frequency of appearance for each word features in each target class and compare it to the existence of the word in all the dataset. Then it calculates the logarithm to inverse the quantity and get the importance indicator value. The result is a measurement of how important the word in is classifying the class.

3. Up Balance

Up-balance will randomly duplicating tuples from the minority class POS, NEG and NUTERAL in order to reinforce the learning process. First, we separated samples from each class into different DataFrames. Next, we resampled every minority class with replacement, setting the number of samples to match that of the majority class. Finally, we combined the up-sampled minority classes' DataFrames with the original majority class DataFrame.

4. Down Balance

Down-balance will randomly remove some samples from the majority classes NEG, OBJ, and NUTRAL to prevent it from dominating the classification process. First, we separated our dataset samples from each class into different DataFrames. Next, we resampled the majority class without replacement, setting the number of samples to match that of the minority class. Finally, we combined the down-balanced majority class DataFrame with the original minority class DataFrame.

5. Decision Trees

Decision trees are the basic building block of random forest. To illustrate the concept, it starts by splitting the data from the root to leaf in order to find correct prediction for the data. Sometimes leafs appears in middle branches to stop the hierarchy of growth, or it will last till the last branch. Before each branch a good question to ask is what is the word that have most informative knowledge to split the data? This is calculated using gini index or entropy., etc. This a high-value question because it greatly reduced the scope of decision tree estimator. So, to arrive at an estimate, we used a series of questions, with each question narrowing our possible values until we were confident enough to make a single prediction. We repeat this decision process repeatedly with only the questions and answers changing.

6. Random Forest

There are too many features to consider, and chances are, each individual guess will be high or low. In decision trees, predictions have variance because they will be widely spread around the right answer. Random forest will allow predictions from hundreds or thousands of decision trees, some of which are high and some of which are low, taking ensample algorithm like random forest will allow the model to combine many decision trees into one. Thus, the probability of predictions made by decision trees being in-accurate will not have a major impact in prediction, the predictions will be closer to the mark on average from all the decision trees involved.

d. Benchmark

In this research, we will use ASDA designed and proposed in [1] as the benchmark model. ASDA stands for Arabic Sentimental Tweet Dataset which is used for Arabic social sentimental analysis. The dataset Consist of more than 10000 tweets gathered from twitter. Each tweet is Classified as: objective, subjective positive, subjective negative or subjective mixed.

The benchmark model was built and evaluated using corpus collected from twitter using eight standard supervised classification models with balanced classes (Figure 3 and Table 2) and unbalanced classes (Figure 3 and Table 3).

Figure 3 Balanced Data Distribution using Bar Chart

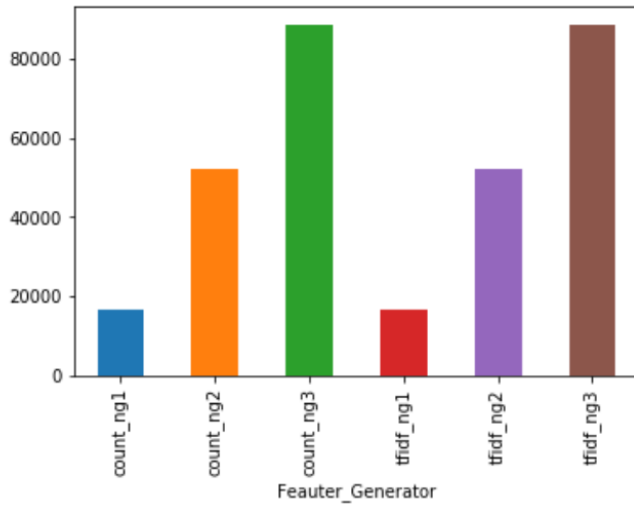


Table 2 Balanced Data feature count for each generator

DataSet Name	Feauter_Generator	Feauters_Count
4-balanced	count_ng1	16457
4-balanced	count_ng2	52036
4-balanced	count_ng3	88672
4-balanced	tfidf_ng1	16457
4-balanced	tfidf_ng2	52036
4-balanced	tfidf_ng3	88672

Figure 4 Un-Balanced Data Distribution using Bar Chart

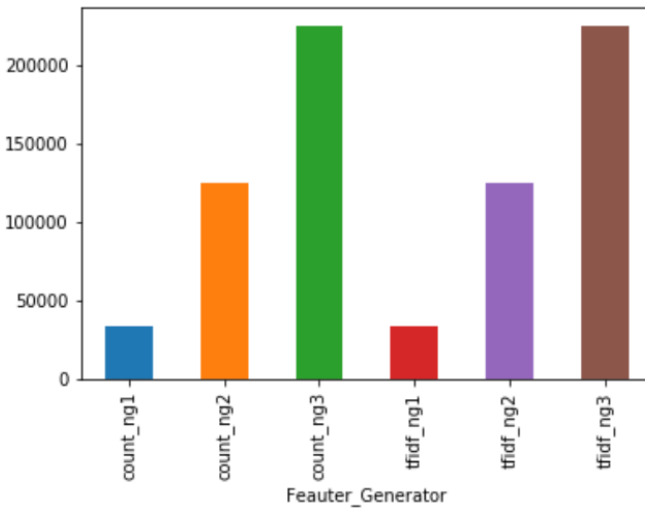


Table 3 Un-Balanced Data feature count for each generator

DataSet Name	Feauter_Generator	Feauters_Count
4-unbalanced	count_ng1	33365
4-unbalanced	count_ng2	124760
4-unbalanced	count_ng3	225111
4-unbalanced	tfidf_ng1	33365
4-unbalanced	tfidf_ng2	124760
4-unbalanced	tfidf_ng3	225111

The Classifiers used in their experiments are: Logistic Regression, Passive Aggressive, Linear Support Vector Machine, Perceptron, BNB, MNB, SGD, K-Nearest Neighbors. All of them used with the dataset taking 20% for testing without parameter tuning.

The dataset in their experiment is preprocessed. They cleaned text by removing special characters, non-Arabic text, white spaces and stop words. They used token count and Frequency-inverse document frequency (Tfidf) explained in IX. Data Preprocessing section to generate features.

VI. Methodology

In order to achieve a solution for the problem statement, the main objectives we followed a clear plan. First, Review the literature to evaluate various experiments and datasets that was used to extract sentiments from Arabic datasets and discuss the outcomes of experimentation by replicating selected experiments model with ASTD dataset. Then, create an enhanced sentimental analysis model that is an upgraded version from the existing models that proved to gain better insights from sentimental dataset. This include a model to compare a set of reviewed experimented models and enhanced model developed after observation. This is done through tuning the parameters using Grid Search method in python and applied in both the standard architecture (replicated experiment) and the enhanced model (using different preprocessing steps and enhancing and testing the standard models before and after parameters tuning). Finally, we statistically evaluate, analyze and interpret the results obtained to assess the accuracy of sentimental analysis models using Arabic dataset. In the below sub-sections the technical detailed implemented in our experiment.

a. Data Preprocessing

In order to processes the data, we implement three different procedures, String token count (Bag of Word), Frequency–inverse document frequency (Tfidf) and data sampling.

First, we implement string token count using Sklearn library, CountVectorizer. It takes each sentence (all the words) present in the data set in the review section and then splits each of the words present in the form of tokens. The occurrence of these tokens in the whole data set are counted in such a way that the count of the occurrence of each token in a positive, negative, natural or objective feedback. Then they are collected separately. Finally, the word frequency of the tokens is calculated.

Then, after splitting the data we calculate Frequency–inverse document frequency (Tfidf). Dataset predicted without Tfidf called count_ng1, count_ng2 and count_ng3, while dataset predicted using Tfidf called tfidf_ng1, tfidf_ng2 and tfidf_ng3. The numbers refer to the features range. If n-gram equals to one, then each unique word will be one feature, while ng3 means that one word, two words and three words each as a bag will be considered three different features. Thus, the number of features increases with the increase of n-grams (more explanation in b. benchmark). One advantage of Tfidf is that it measures how important a word is to differentiate each category. It reduces the weightage of more common words like (stop words or common words) which occurs in all tweets.

Finally, to sample from our data and to follow ASDA research paper, the data was down sampled and up sampled to generate a balanced dataset, in addition to using the unbalanced dataset (statistics described in Exploratory Visualization).

In up sampling, we up-sampled the minority ratings these are POS, NEG and NEUTRAL. First, we separate samples from each class into different DataFrames. Next, we resampled every minority class with replacement, setting the number of samples to match that of the majority class. Finally, we combined the up-sampled minority classes' DataFrames with the original majority

class DataFrame. The process is like that of up-sampling, but we down sample all the majority classes to the minority class POS, without replacement as there is no need for it when data is reduced.

b. Implementation

We started by replicating the experiment in ASDA research paper. Then, we add two more classifiers to have a total of ten different classifiers. Additionally, we used the hyper tuning parameters to give our classifier different options and increase the chance of getting better accuracy score. The justification for using a lot of classification model is that we follow “no free lunch” theorem where we hyper tuned all the models to get as much options as we can. Figure 5 block shows our hyper parameters settings for each classifier.

Figure 5 Parameter Settings using gridSearch for all classifiers

MNB	{'alpha': [1.0, .1, .3, .6, .8, 0.0], 'fit_prior': [True, False]}
BNB	{'alpha': [1.0, .8, .6, .4, .2], 'binarize': [0.0, .25, .5], 'fit_prior': [True, False]}
SVM	{'C': [1.0, .7, .5, .3, .1], 'dual': [True], 'fit_intercept': [True, False], 'loss': ['hinge', 'squared_hinge'], 'random_state': [0], 'multi_class': ['ovr', 'crammer_singer'], 'intercept_scaling': [.3, .5, 1], 'fit_intercept': [True, False]}
Passive Aggressive	{'C': [1.0, .7, .5, .3, .1], 'loss': ['hinge'], '#n_iter': [15], 'n_jobs': [1], 'random_state': [0], 'shuffle': [True, False], 'verbose': [0], 'warm_start': [False, True]}
SGD	{'alpha': [0.0001, 0.0005], 'average': [False], 'class_weight': [None], 'epsilon': [0.1], 'eta0': [0.0], 'fit_intercept': [True], 'l1_ratio': [0.15, 0.001, 0.0001], 'learning_rate': ['optimal'], 'loss': ['hinge'], 'max_iter': [None], 'n_iter': [None], 'n_jobs': [1], 'penalty': ['l2'], 'power_t': [0.5], 'random_state': [None], 'shuffle': [True], 'tol': [None], 'verbose': [0], 'warm_start': [False]}
Logistic Regression	{'C': [1.00000000e+00, 2.78255940e+00, 7.74263683e+00, 2.15443469e+01, 5.99484250e+01, 1.66810054e+02, 4.64158883e+02, 1.29154967e+03, 3.59381366e+03, 1.00000000e+04]}
Linear Perceptron	{'alpha': [1, .8, .5, .3, .1, .01, .001, .0001], 'fit_intercept': [True, False], '#n_iter': [5, 10, 15], 'random_state': [0], 'penalty': ['l1', 'l2'], 'shuffle': [True, False], 'warm_start': [True, False], 'eta0': [1.0, .5]}
KNN	{'algorithm': ['auto'], 'leaf_size': [30], 'metric': ['minkowski'], 'n_jobs': [1], 'n_neighbors': [2, 4, 8], 'p': [2], 'weights': ['uniform']}
Random Forest	{'bootstrap': [True], 'class_weight': [None], 'criterion': ['gini'], 'max_depth': [None], 'max_features': ['auto'], 'max_leaf_nodes': [None], 'min_impurity_decrease': [0.0], 'min_impurity_split': [None], 'min_samples_leaf': [1], 'min_samples_split': [8, 12], 'min_weight_fraction_leaf': [0.0], 'n_estimators': [10], 'n_jobs': [1], 'oob_score': [False], 'random_state': [0], 'verbose': [0], 'warm_start': [False]}
Decition Trees	{'criterion': ['gini'], '#', 'entropy', 'max_depth': [2, 4], 'presort': [False], 'random_state': [0], 'splitter': ['best']}

After applying preprocessing steps mentioned earlier, we looped through all the dataset structures, split, train and test our data before and after tuning each classifiers our statistical results proved to be better than the benchmark model in all structures.

One weakness noticed is that we calculated f1 score using external function for the visuals only, it turned our that the function EvaluateResult does not measure the F1 score correctly because there is a big differene between our experiment results and the visuals F1-Score, due to time and deadlines we did not fix this issue, however, we make sure that all the results in Table 4 are correct and based on Sklearn Accuracy and F-Score functions with Average='macro'.

Table 4 Accuracy and F-Score for Up Balanced, Down Balanced and Un Balanced datasets using ten different tuned classifiers

Features	Tf-Idf	Up Balanced ASDA dataset			Down Balanced ASDA dataset			Un Balanced ASDA dataset		
		1g	1g+2g	1g+2g+3g	1g	1g+2g	1g+2g+3g	1g	1g+2g	1g+2g+3g
MNB	No	0.932/0.931	0.951/0.950	0.950/0.949	0.528/0.526	0.516/0.510	0.511/0.506	0.697/0.613	0.697/0.613	0.695/0.599
	Yes	0.937/0.936	0.951/0.950	0.951/0.949	0.531/0.527	0.522/0.516	0.527/0.522	0.687/0.573	0.686/0.570	0.686/0.569
BNB	No	0.914/0.913	0.938/0.937	0.944/0.943	0.567/0.564	0.567/0.567	0.558/0.556	0.690/0.586	0.683/0.562	0.681/0.558
	Yes	0.914/0.912	0.939/0.937	0.944/0.943	0.567/0.563	0.567/0.567	0.558/0.559	0.689/0.586	0.683/0.562	0.681/0.558
SVM	No	0.950/0.948	0.970/0.970	0.975/0.974	0.545/0.542	0.548/0.549	0.547/0.547	0.700/0.630	0.701/0.632	0.699/0.629
	Yes	0.938/0.935	0.951/0.950	0.950/0.949	0.542/0.540	0.514/0.511	0.513/0.509	0.696/0.623	0.694/0.641	0.695/0.636
Passive Aggressive	No	0.950/0.95	0.965/0.965	0.970/0.970	0.514/0.514	0.523/0.525	0.516/0.512	0.671/0.636	0.695/0.637	0.700/0.632
	Yes	0.949/0.948	0.959/0.958	0.957/0.957	0.512/0.510	0.525/0.523	0.514/0.512	0.691/0.648	0.689/0.646	0.694/0.647
SGD	No	0.947/0.946	0.966/0.967	0.971/0.971	0.542/0.543	0.531/0.529	0.534/0.537	0.691/0.622	0.697/0.624	0.695/0.620
	Yes	0.908/0.905	0.924/0.922	0.926/0.924	0.528/0.525	0.523/0.519	0.517/0.516	0.698/0.632	0.699/0.637	0.700/0.639
Logistic Regression	No	0.948/0.946	0.968/0.968	0.973/0.973	0.545/0.543	0.541/0.539	0.536/0.536	0.695/0.633	0.699/0.628	0.698/0.620
	Yes	0.943/0.941	0.950/0.949	0.950/0.949	0.533/0.530	0.517/0.514	0.506/0.502	0.689/0.584	0.691/0.637	0.694/0.644
Linear Perceptron	No	0.940/0.938	0.957/0.956	0.958/0.957	0.511/0.509	0.539/0.540	0.516/0.512	0.644/0.620	0.682/0.632	0.676/0.614
	Yes	0.939/0.937	0.945/0.943	0.941/0.939	0.525/0.524	0.519/0.516	0.516/0.512	0.646/0.620	0.664/0.630	0.661/0.632
KNN	No	0.892/0.883	0.877/0.867	0.883/0.876	0.369/0.345	0.373/0.326	0.394/0.328	0.673/0.557	0.674/0.561	0.675/0.558
	Yes	0.912/0.907	0.904/0.897	0.898/0.890	0.420/0.419	0.427/0.428	0.425/0.426	0.685/0.603	0.690/0.618	0.692/0.620
Random Forest	No	0.967/0.967	0.972/0.972	0.971/0.971	0.517/0.516	0.509/0.520	0.525/0.530	0.689/0.596	0.692/0.597	0.688/0.585
	Yes	0.965/0.965	0.967/0.966	0.966/0.969	0.519/0.518	0.508/0.504	0.528/0.520	0.690/0.601	0.689/0.602	0.689/0.585
Decition Trees	No	0.936/0.933	0.930/0.926	0.930/0.926	0.439/0.437	0.464/0.461	0.455/0.454	0.683/0.576	0.683/0.576	0.683/0.577
	Yes	0.924/0.919	0.916/0.910	0.915/0.909	0.422/0.415	0.418/0.416	0.459/0.459	0.687/0.583	0.687/0.585	0.684/0.577

c. Refinement

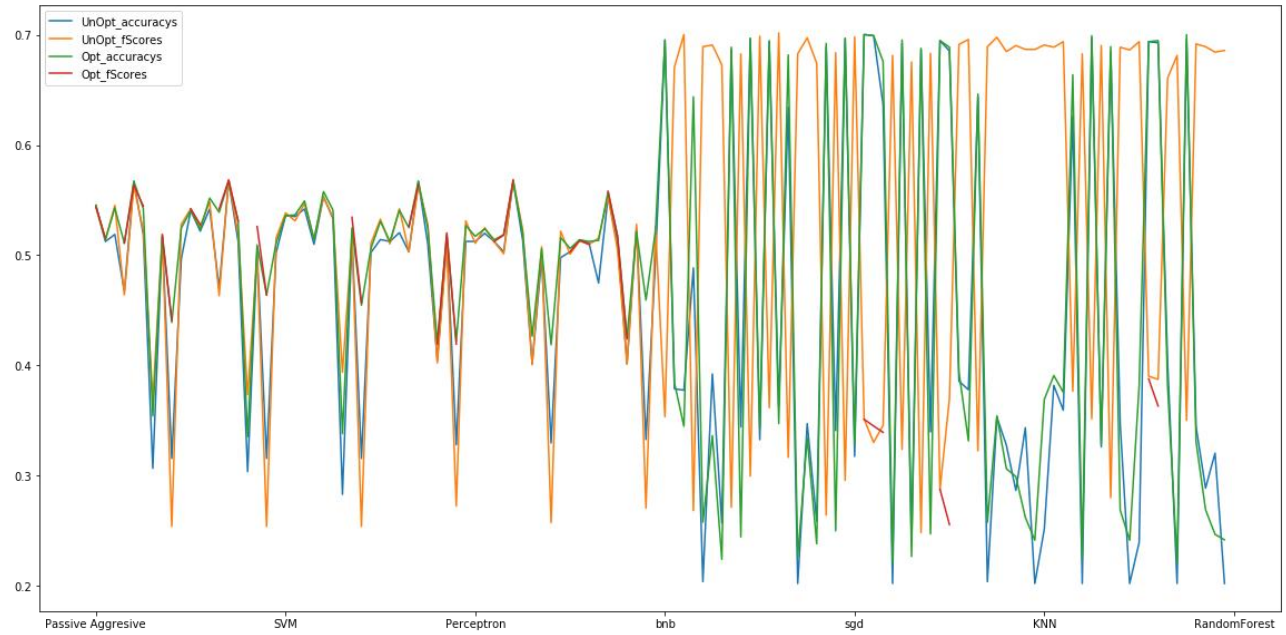
Our method in up and down resampling dataset out performed the balanced data in classification accuracy reported in ASDA research paper. The most obvious differences between our model and ASDA model is that they used standard classifiers without cross validation nor parameter tuning, in addition to dis-similarity in resampling methods. For our experiment, we choose to up sample the data and down sample it to make it balanced (check Data Preprocessing section). In ASDA research paper, they did not explicitly state their methodology for balancing the data, however, it seems to us that it was manually manipulated.

VII. Results

a. Model Evaluation and Validation

Our model examined the accuracy using cross validation and can predict the data with accuracy reached 70% higher than the accuracy reported in ASDA research paper using un-balanced data and LinearSVC with two folds.

Figure 6 Down-sampled and Un-Balanced ASDA dataset



Also, the model reaches all benchmark model with the down streamed data with accuracy as high as %70 with SGD and Passive Aggressive (Table 5).

Table 5 Down Balanced ASDA model performance before and after parameter tuning

	classifiers	data	feauter_generator	Best_Settings	UnOpt_accuarcys	UnOpt_accuarcys	Opt_accuarcys	Opt_fScores
81	Passive Aggresive	unBalanced	count_ng3	{'warm_start': False, 'loss': 'hinge', 'C': 1.0, 'n_jobs': 1, 'shuffle': True, 'verbose': 0, 'fit_intercept': True, 'average': False, 'max_iter': None, 'random_state': 0, 'tol': None, 'n_iter': None, 'class_weight': None}	0.7	0.351	0.7	0.351
115	sgd	unBalanced	tfidf_ng3	{'warm_start': False, 'loss': 'hinge', 'n_jobs': 1, 'eta0': 0.696, 'epsilon': 0.1, 'average': False, 'max_iter': None, 'penalty': 'l2', 'power_t': 0.5, 'random_state': 0, 'tol': None, 't1_ratio': 0.15, 'n_iter': None, 'alpha': 0.0001, 'learning_rate': 'optimal', 'class_weight': None}	0.696	0.35	0.7	0.366

Upstream Sampling out perform all the above with accuracy as high as %97 with RandomForest using 3-ngram range of features (Table 6).

Table 6 Up Balanced ASDA dataset performance

	classifiers	feature_generator	UnOpt_accuracy	UnOpt_fScores	Opt_accuracy	Opt_fScores
7	RandomForest	count_ng1	0.960209	0.959525	0.97422	0.973823
12	SVM	count_ng2	0.969923	0.969272	0.970297	0.969717
17	RandomForest	count_ng2	0.96488	0.964357	0.975528	0.975166
20	Logistic Regression	count_ng3	0.971044	0.970465	0.973286	0.972798
21	Passive Aggressive	count_ng3	0.970297	0.96985	0.970484	0.970042
22	SVM	count_ng3	0.97366	0.973239	0.974594	0.974148
27	RandomForest	count_ng3	0.964506	0.964174	0.978517	0.978283

b. Justification

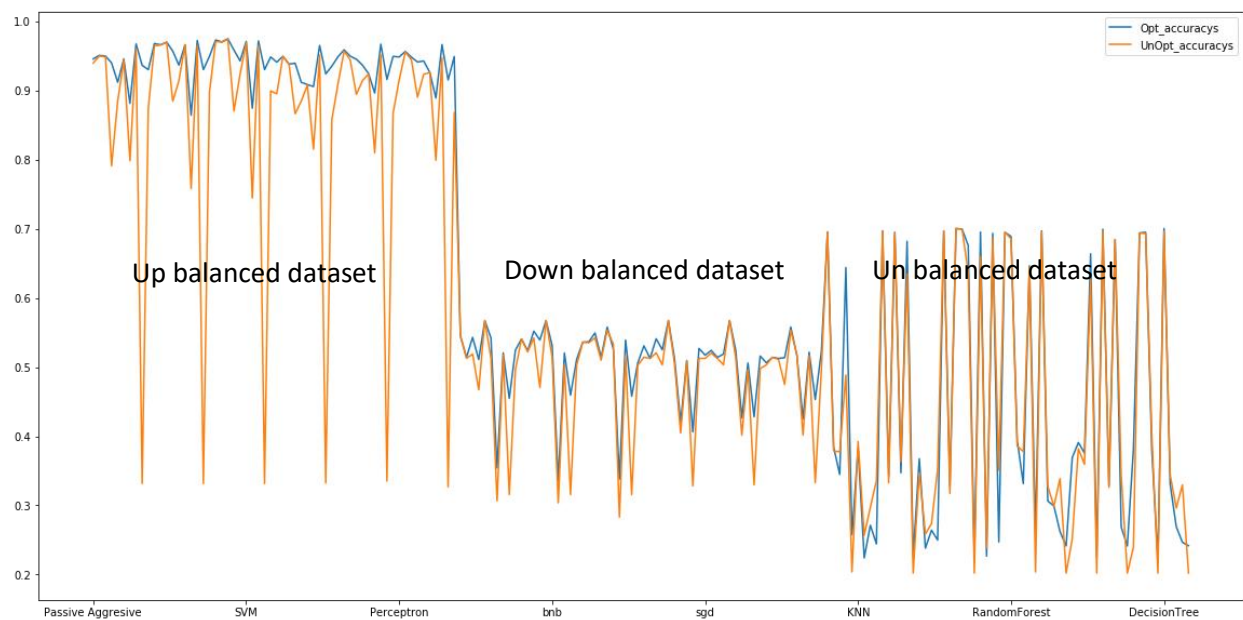
We clearly can see the up-balanced sampled dataset out perform all our models using the ten classifiers, despite of using Grid Search to tune the hyperparameters. However, it highly likely that this model will suffer from overfitting. This due to the fact we increased the number of copies to reach the majority class OBJ for all the other target classes.

VIII. Conclusion

a. Free-Form Visualization

In our experiment we succeeded in improving the performance of the benchmark model which reaches as high as %69.4. Our model using un balanced ASDA dataset reached %70, with significant drop of recall and precision with F-Score merely reached %38. Using down balanced dataset our model successfully competed with the benchmark with accuracy and F-score higher than %49.3 reported in ASDA research paper and reached around %56.7 for both measurements. The highest classification accuracy and F-Score gained with up balanced data with accuracy and F-score above %95. The shows the performance of our ten classifiers using up balanced, down balanced and unbalanced ASDA dataset respective.

Figure 7 Performance of our ten classifiers with up-balanced, down-balanced and un-balanced ASDA dataset.



b. Reflection

In our experiment, we started by replicating ASDA research paper to get the same results as they reported. Once we finalized their work, we started our own by implementing the steps mentioned in the Data Preprocessing section. We find out that Up Sampling increases the accuracy while the model will be vulnerable for overfitting for unseen data.

Time is a big challenge because I am using my personal laptop and each algorithm, especially the up sampled models, takes prohibitively very long time to process. I found this project very challenging and interesting at the same time, I have learned how to replicate a research paper and get the result, analyze them and report them. Writing is also one of the most challenging tasks.

c. Improvements

We aim to improve our model by measuring the time each classifier takes. We better reduce the number of features using feature reduction techniques and may use Sparks to increase performance.

Additionally, we need to increase the dataset and test it using unseen corpus to get better insight about its performance, especially with the classifiers trained using the up balanced dataset.

IX. References

[1] Nabil, M., Aly, M., & Atiya, A. (2015). Astd: Arabic sentiment tweets dataset. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing* (pp. 2515-2519).