

Empirical evaluation of shallow and deep learning classifiers for Arabic sentiment analysis

ALI BOU NASSIF, Computer Engineering Department, University of Sharjah, UAE

ABDOLLAH MASOUD DARYA, Electrical Engineering Department, University of Sharjah, UAE

ASHRAF ELNAGAR, Computer Science Department, University of Sharjah, UAE

This work presents a detailed comparison of the performance of deep learning models such as convolutional neural networks (CNN), long short-term memory (LSTM), gated recurrent units (GRU), their hybrids, and a selection of shallow learning classifiers for sentiment analysis of Arabic reviews. Additionally, the comparison includes state-of-the-art models such as the transformer architecture and the araBERT pre-trained model. The datasets used in this study are multi-dialect Arabic hotel and book review datasets, which are some of the largest publicly available datasets for Arabic reviews. Results showed deep learning outperforming shallow learning for binary and multi-label classification, in contrast with the results of similar work reported in the literature. This discrepancy in outcome was caused by dataset size as we found it to be proportional to the performance of deep learning models. The performance of deep and shallow learning techniques was analyzed in terms of accuracy and F1 score. The best performing shallow learning technique was Random Forest followed by Decision Tree, and AdaBoost. The deep learning models performed similarly using a default embedding layer, while the transformer model performed best when augmented with araBERT.

CCS Concepts: • **Information systems** → **Sentiment analysis**; *Data cleaning*; • **Computing methodologies** → **Supervised learning by classification**.

Additional Key Words and Phrases: Deep learning, shallow learning, learning curve, embedding, misclassification.

ACM Reference Format:

Ali Bou Nassif, Abdollah Masoud Darya, and Ashraf Elnagar. 2021. Empirical evaluation of shallow and deep learning classifiers for Arabic sentiment analysis. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.* 21, 1, Article 14 (November 2021), 24 pages. <https://doi.org/10.1145/3466171>

1 INTRODUCTION

With the rise of Web 2.0 and the ubiquity of online social networks, massive amounts of opinion-based data are generated by users that list their opinions and personal views online. Analyzing such data would produce valuable information regarding key trends, product evaluations, stock market predictions, and public opinion surveys. Sentiment analysis (SA) is a field that aims to extract sentiments relayed in pieces of text based on its contents. Unfortunately, automated systems of SA have faced challenges in accurately labeling human sentiment due to the complex nature of the semantics

Authors' addresses: Ali Bou Nassif, Computer Engineering Department, University of Sharjah, Sharjah, UAE, anassif@sharjah.ac.ae; Abdollah Masoud Darya, Electrical Engineering Department, University of Sharjah, Sharjah, UAE, abdollah.masoud@ieee.org; Ashraf Elnagar, Computer Science Department, University of Sharjah, Sharjah, UAE, ashraf@sharjah.ac.ae.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2021 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

1

used in human language. Moreover, the Arabic language adds another hurdle to automated SA due to its large number of dialects, its morphological richness and its ingrained ambiguity [6, 36].

Arabic-speaking internet users have grown by 9,348%¹ in the last twenty years, marking the largest growth value in any language. Furthermore, following English, Mandarin and Spanish, Arabic speakers are the fourth largest online linguistic population, tallied at 5.2% of the total online population². Despite being one of the fastest-growing languages in terms of online users, the field of Arabic SA is still not as mature as its English counterpart. However, the last few years have shown growing interests in the field of Arabic SA [1]. In addition to the task of SA [10], also referred to as opinion mining [15, 19], there is also great interest in the fields of emotion mining [16, 54] and text classification [32].

This work conducts a comprehensive comparison of some of the most widely used deep and shallow learning (SL) techniques using two of the largest publicly available datasets for multi-dialectal Arabic reviews: Hotel Arabic Reviews Dataset (HARD) [34] and the Book Reviews in Arabic Dataset (BRAD) [33, 35] (see section 3.1). It compares four of the main deep learning (DL) models, which requires further study in the field of Arabic natural language processing [2, 9, 51].

Section 2 presents the main contributions found in the literature of Arabic SA using DL classifiers. In section 3, the methodology and datasets used in this work are discussed, as well as the pre-processing steps followed. Section 4 discusses the results obtained, along with the learning curve. Finally, the conclusion summarizes the findings of this study and proposes potential avenues for future work.

1.1 Research Objectives

The results of this work prove that DL models are superior to SL models for SA of Arabic reviews, given that the dataset used is of adequate size. This is in contrast to the work done in [7], which found that support vector machines (SVM) performed significantly better than convolutional neural networks (CNN) and recurrent neural networks (RNN) for the task of Arabic sentiment polarity identification. We hypothesize that the findings of [7] were a result of using a dataset of an inadequate size that does not utilize the full potential of DL techniques, and therefore that they deliver an unfair comparison. Not only does this work consider binary sentiment classification, it also compares between the best performing SL and DL classifiers for multi-label classification (5-label), which strengthens the comparison established between SL and DL methods.

The contributions of this work are summarized as follows:

- To the best of our knowledge, this is the first work that comprehensively compares the performance of SL and DL techniques for the purposes of Arabic SA. Not only is the evaluation done for the three main DL architectures (CNN, LSTM and GRU) but also all possible combinations of the aforementioned techniques. Furthermore, the state-of-the-art transformer architecture was studied and compared against the previously mentioned DL models.
- In this paper we prove that when a dataset of adequate size is used, DL classifiers outperform SL classifiers.
- We establish the point at which the best DL technique outperforms the best performing SL technique by plotting the learning curve for various dataset sizes. This was done for both binary and 5-label sentiment classification.
- We conduct a study highlighting the effect of different embedding layers on the performance of the transformer model. The comparison includes default embeddings derived from the analyzed dataset, custom embeddings acquired from Arabic websites, as well as pre-trained models such as araBERT.

¹<https://www.internetworldstats.com/stats7.htm>

²<https://www.internetworldstats.com/stats7.htm>

- This is the first study to comprehensively evaluate both the HARD and BRAD datasets, which are some of the largest publicly available Arabic review datasets.

2 LITERATURE REVIEW

Arabic SA is a highly challenging research area that involves complex tasks. The tasks that have been the most frequent subjects of study are subjectivity classification, sentiment classification, lexicon creation, aspects extraction, aspect sentiment classification and opinion spam detection [13, 17, 22, 23, 31, 42].

The work done in the field of Arabic SA is vast and encompasses many fields and niche applications which cannot be covered in this section, thus the reader is referred to the following surveys that present the most recent updates in the field [1, 2, 14, 28, 55]. However, with that being said, there are few contributions that will be highlighted in the following paragraph that are relevant to the scope of this work.

The authors in [6] investigated the results of using DL for sentence-level Arabic SA. In their work, they used Deep Belief Networks (DBN), Deep Autoencoders (DAE), Deep Neural Networks (DNN) and Recursive Auto Encoder (RAE). They used an annotated collection of 1,180 sentences selected from the Data Consortium Arabic Tree Bank (LDCATB) dataset. Their approach showed that RAE performed best out of the four models. The authors then worked on improving the results of the RAE model and introduced AROMA [5]: a recursive DL model for opinion mining in Arabic. An improvement in classification accuracy by upto 12% was achieved by performing word sentiment embedding and generating syntactic parse trees to be used as a reference to AROMA's recursion. The authors in [21] used Recursive Neural Tensor Networks (RNTN) for SA of Arabic text. They used a subset of the Qatar Arabic Language Bank (QALB) data set, consisting of 1,177 sentences. They claim that the proposed approach outperforms SVM, RAE, and LSTM. In [4], the authors use CNN and LSTM networks with word embedding features on the Arabic Sentiment Tweets Dataset (ASTD), which consists of 10,000 sentences. In their evaluation, they found that the LSTM performed best, with an accuracy of 87%. In [20], the authors compared an English state-of-the-art method to classify binary sentiment in Arabic tweets against a cluster-based sentiment classification approach and RAE. Their results show that the English state-of-the-art method outperformed the competition. A common trend seen in most related work is their reliance on datasets composed of fewer than 1000 sentences, while DL techniques demand the use of vast and comprehensive datasets.

The first use of hybrid DL techniques was done recently in the work presented in [8] where a combined CNN and LSTM architecture was used to analyze binary sentiment in several Twitter datasets. The largest dataset used in this study contains 2,500 sentences, which resulted in accuracy in the mid-90-s for tweet-level SA.

3 METHODOLOGY

In this section we start by introducing the selected datasets. We then discuss the pre-processing steps taken to prepare the data for the SA task, and we describe the SL and DL techniques used in our empirical evaluation.

3.1 Datasets

HARD. The first dataset used is the Hotel Arabic Reviews Dataset (HARD) [34], which contains 373,750 reviews³ (see Table 1 for a sample of the reviews). Two variants of the HARD dataset are available: balanced and unbalanced. The unbalanced dataset was chosen as the target of this study, because it is the larger of the two. The number of unique

³compiled from booking.com

Table 1. Two sample reviews from HARD with their literal English translation.

Rating	Review
4	استثنائي. سعر جدا مناسب وأتمنى ان لا يزيد فلمستقبل حتى نظل عملائكم ونظافة الغرف وخدمة البوفيه Exceptional. Price very reasonable and I wish it does not increase in the future so that we remain your customers and cleanliness of rooms and buffet service
2	كانت سيئه جدا بسبب سوء المعامله. الجو. كانت تجربه فاشله لاني تعرضت لسوء تعامل من موظفي الاستقبال It was very bad due to the bad treatment. The weather. It was a failed experiment because I was mistreated by the reception staff.

Table 2. HARD dataset variables.

Column #	Variable	More Information
1	Rating	Possible values are 1-5
2	Review ID	Unique ID assigned to each review
3	Hotel ID	ID assigned to each hotel
4	User ID	Unique ID assigned to each user
5	Room ID	Unique ID assigned to each room
6	# of Nights	Number of nights
7	Review	Reviewers opinion written in MSA/DA

Table 3. Two sample reviews from BRAD with their literal English translation.

Rating	Review
2	الروايه مفيهاش احداث كثير . ممله لحد كبير. بس الاسلوب في الكتابه كان حلو جدا. بس كنت متوقعه الروايه احلي من كده The novel does not have many events. Boring to a large extent. But the style of writing is very good. But I was expecting the novel better than this.
5	احمد مراد بأسلوبه المبدع يتألق في سرد ورواية التاريخ و يلقى الضوء على شخصيات في ظل هذه الحقبه التاريخيه والأحداث Ahmad Murad (the author) with his creative style shines in narrating history and casting a light on characters under this historical era and events.

tokens in the HARD dataset is in the order of 90K. Note: the variables contained in the HARD dataset are listed in Table 2. What makes this dataset unique is that it combines reviews in Modern Standard Arabic (MSA) and Dialectal Arabic (DA). It is also the largest publicly available dataset of Arabic hotel reviews.

BRAD. The Book Reviews in Arabic Dataset (BRAD) [35] consists of a collection containing 508,538 book reviews⁴ in total (see Table 3 for a sample of the reviews). Similar to HARD, BRAD contains reviews in MSA and DA. The BRAD dataset is made up of around 220K unique tokens. Note: the variables contained in the BRAD dataset are listed in Table 4. These datasets are considered to be two of the largest available for the applications of Arabic SA and machine learning.

⁴taken from www.goodreads.com

dataset of 91,000 reviews. As BRAD was a larger dataset, 60,000 samples were taken randomly from 63,104 negative and 249,737 positive reviews to form a balanced dataset of 120,000 reviews to be used for binary classification. For 5-label classification, the HARD dataset was balanced around the smallest label, which was the lowest rating of 1. Since rating 1 contained around 12,500 reviews, random samples of the same size (12,500) were taken from each label, adding up to a full dataset of 62,500 reviews. Note that this process was implemented after all the previous data cleaning and pre-processing steps. The training size was selected to be 90% of the balanced dataset, while 10% was selected for testing using the holdout technique. Out of the 90% training set, 10% was selected for validation for DL, and 9-fold cross-validation for SL. Furthermore, the performance metrics of the DL techniques were obtained as an average of 10 runs, where the dataset is randomly resampled at the beginning of each run.

Step.9 Finally, tokenization was implemented, by splitting the text into separate tokens (words). Next, the integer vector was defined for each token and each sentence was padded to a length of 100 tokens.

Note: no other preprocessing techniques were used as they could bias the results towards a specific model and therefore provide an unfair advantage, i.e., feature selection and extraction [45, 62]. Furthermore, several techniques generally seen in preliminary trial and error tests to either reduce performance or leave it unchanged at the cost of added computation—such as stemming [9] or removing stopwords—were not implemented. The resultant datasets were then tested and compared using SL and DL techniques.

3.3 Shallow Learning

Decision Tree (DT): Decision tree [59] is a classification algorithm that divides a set of data into smaller sets based on tests that are defined at each node in the tree. The tree consists of a root node, internal nodes and terminal nodes (leaves). Each node in the tree has one parent node and two or more child nodes. Based on this, each data set is classified by assigning it a location according to the framework defined by the tree.

Random Forest (RF): A random forest [24] consists of a combination of tree predictors where each tree is dependent on the variables of a random vector sampled independently, with all the trees in the forest having the same distribution. In random forests, the best split is used among a randomly chosen subset of predictors at each node. This adds robustness against overfitting and improved performance at classification [49].

Support Vector Machines (SVM): Support Vector Machines (SVM) [24] is the most widely used machine learning network for two-group classification problems. Conceptually, input vectors are mapped non-linearly to a feature space with high dimensionality. A linear decision surface plane or line is constructed in the feature space to separate the two classes also known as a hyperplane. In this work, SVM was implemented using two kernels: linear (lin.) and radial basis function (RBF).

k-nearest neighbors (KNN): K-Nearest Neighbours (KNN) is a simple classification technique [11]. For instance, to classify a data record n_{test} , K nearest neighbors are saved to form the neighborhood for n_{test} . The classification is based on the most frequent nearest neighbors, either with or without distance weighting. But to apply KNN an appropriate value of K needs to be considered which makes the technique biased.

Multilayer Perceptron (MLP): Multi-Layer Perceptron consists of one input layer of input neurons followed by one or more hidden layers and a single output layer [58]. Each layer is composed of nodes or neurons that are fully connected to nodes in the subsequent layers up to the output layers. The activation function for the initial input layer is linear with no thresholds. Subsequently, however, the hidden node layers have non-linearity in their activation function

in addition to the threshold. The output layer activity is dependent on linear function and threshold. Each hidden unit node is related to the addition of the weighted sum of every input node in the initial layer and the associated threshold. The output layer is associated with the hidden layers in the same manner.

Gaussian Naïve Bayes (GNB): The Naïve Bayes (NB) classifier applies the Bayes theorem with the assumption that each pair of features are independent. The NB classifier finds the probability that a given instance belongs to a certain class. Gaussian Naïve Bayes, on the other hand, executes the classification by assuming the likelihood of the features to be Gaussian.

AdaBoost (AB): The Adaptive Boosting (AdaBoost) algorithm [37] is one of the first proposed boosting algorithms and works by combining several relatively weak and inaccurate models to create an accurate prediction model. AdaBoost can be used to substantially reduce learning algorithms errors and is mainly aimed at classification applications [38].

3.4 Deep Learning

CNN: Convolutional Neural Networks (CNN) [47] are feedforward neural networks originally conceived for the field of computer vision and have shown to be effective for natural language processing (NLP) applications [46]. They utilize a layer with convolving filters applied to local features. They feature convolution in place of general matrix multiplication, which is present in standard neural networks. This decreases the number of weights, thereby reducing the complexity of the network, causing it to be one of the best-performing DL techniques in terms of execution time. Furthermore, another advantage of CNN is that it requires minimal preprocessing. Along with its low complexity, this paved the way for its use in NLP, speech and handwriting recognition, and image classification, amongst many others [50].

LSTM: Long short term memory (LSTM) [43] network is a type of recurrent neural networks (RNN) that is effective at learning problems related to sequential data. It tackles these problems by capturing long-term temporal dependencies. LSTM does not suffer from the optimization issues affecting the basic form of RNN due to its complex nature and the repetition of its modules [60]. The basic idea behind the LSTM architecture is a memory cell that maintains its state over time and nonlinear gating units that control information flow in and out of the cell [41]. It also features three main gates: the input, forget and output gates. The input block is connected to the output block and all the gates.

GRU: The gated recurrent unit (GRU) framework was proposed by [25] in 2014. Similar to LSTM, GRU contains gating units that control the flow of information. While in LSTM networks, the gate controls the amount of memory that is utilized by other units in the network, in GRU all contents are exposed without any restriction. It has been reported, however, that GRU outperforms LSTM for nearly all tasks except language modeling [44]. Furthermore, the gap between the performance of LSTM and GRU networks can be minimized by initializing LSTM's forget gate bias to one. GRU was previously used in several Arabic NLP tasks such as [3].

Transformer: The transformer (TRANS) model was first proposed in [63]. The transformer consists of an encoder and a decoder. The input sequence is taken by the encoder and mapped into a higher dimensional space. The decoder then produces an output sequence from the mapped input. It has been reported to train significantly faster than recurrent and convolutional architectures for translation tasks [63]. Transformers (feed-forward architecture) allow for efficient training on huge datasets, with the objective of simply predicting words based on their context. Building such models is very expensive; however, a variety of models are published and ready to use in downstream tasks such as SA. It has been reported that fine-tuning these models on some smaller, supervised datasets can improve classification results. This process is called transfer learning. Therefore, instead of building a transformer model from scratch, we will simply take an existing one (in this case: araBERT [18]) and utilize its parameters to initialize the sentiment classifier and achieve the result (see Figure 1).

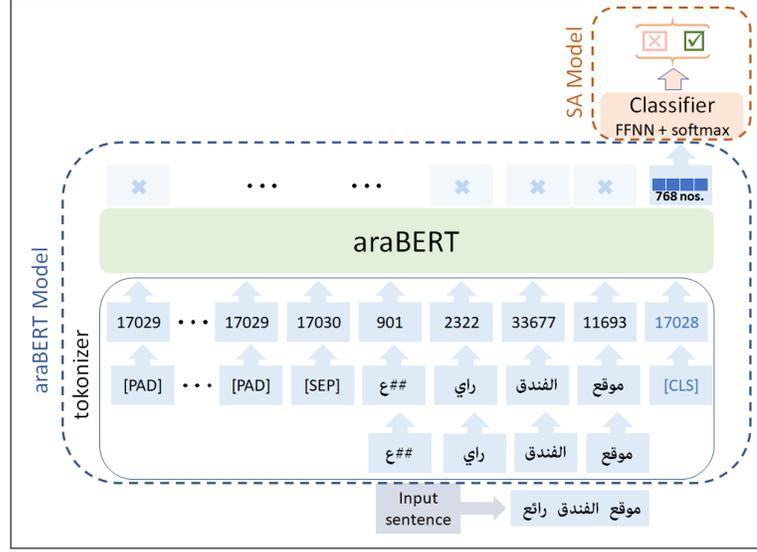


Fig. 1. The workflow is summarized in the pipeline diagram above showing how the sentiment classifier (the fine-tuning task) is added on top of the araBERT transformer model.

This work evaluates the three main DL architectures discussed above. Those are CNN, LSTM and GRU, and all of their possible combinations, as well as the TRANS model. Note that bidirectional LSTM (Bi-LSTM) has been used in place of LSTM, as it performed better in preliminary trial and error tests. Training, testing and pre-processing were executed on Google Colab (<https://colab.research.google.com/>), with a Tesla K80 GPU, an Intel Xeon processor and 12 GB of RAM.

4 RESULTS AND DISCUSSION

In this study, the performance of the techniques discussed above was tested using some of the most commonly used performance metrics. These are accuracy, precision, recall, F1, area under the curve of the receiver operating characteristics, training time and testing time [26]. Accuracy is simply the number of correct predictions divided by the total number of predictions. Accuracy can be represented as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

where the sum of true positives and true negatives is divided by the sum of true positives, true negatives, false positives and false negatives. True positives (TP) are positive examples labeled correctly as positives. True negatives (TN) are negative examples labeled correctly as negative. False positives (FP) are negative examples labeled incorrectly as positive. False negatives (FN) are positive examples labeled incorrectly as negative [26]. Precision is defined as:

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

and recall can be presented as:

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

Table 5. Binary classification results from training and testing the HARD (H) and BRAD (B) datasets using SL. Note: all results are rounded to three significant digits. **Bold** results represent the best performance from each category.

Classifier	Accuracy		Precision		Recall		F1		AUC	
	H	B	H	B	H	B	H	B	H	B
DT	0.871	0.741	0.868	0.732	0.876	0.747	0.872	0.740	0.870	0.740
RF	0.897	0.809	0.891	0.887	0.905	0.700	0.898	0.783	0.900	0.810
SVM (RBF)	0.670	0.701	0.607	0.675	0.858	0.431	0.711	0.526	0.670	0.700
KNN	0.677	0.666	0.680	0.680	0.672	0.605	0.676	0.640	0.680	0.660
MLP	0.526	0.513	0.611	0.667	0.149	0.018	0.240	0.036	0.530	0.500
GNB	0.502	0.588	0.501	0.673	0.980	0.315	0.663	0.429	0.500	0.580
AB	0.842	0.690	0.823	0.676	0.872	0.712	0.846	0.693	0.840	0.690
SVM (Linear)	0.477	0.497	0.481	0.493	0.573	0.779	0.523	0.603	0.480	0.500

The F1 measure can be deduced from the following:

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

The area under the curve (AUC) of the receiver operating characteristics (ROC) presents an important metric for binary classification, as it portrays the ability of the system to distinguish between good and bad reviews [26]. The best-performing system would have AUC values closer to one, while the worst-performing systems would have AUC values closer to zero. The parameters of all models used are presented in section 6. For multi-label classification, other metrics are used such as the macro/micro/weighted average for precision, recall and F1 [52]. We note, however, that since we are using a balanced dataset, the metrics obtained from the macro/micro/weighted average would yield similar results. Therefore, for simplification, we are referring to the precision, recall and F1 for multi-label classification as average precision, average recall and average F1 respectively.

4.1 Shallow Learning

Binary Classification. The results in Table 5 and Figure 2 clearly show the superiority of three SL techniques from the rest. Random Forest performed the best, followed by Decision Tree and AdaBoost. It can also be seen that in most cases, the SL techniques performed best when being trained and tested on the HARD dataset as compared to BRAD, even though BRAD is a larger dataset than HARD (120K vs 90K reviews). This can be attributed to the size of the vocabulary being utilized in each dataset. Where BRAD contains unique vocabulary in the order of 220K, HARD contains only 90K, less than half that amount. This agrees with the findings of [48], where it was found that datasets with shorter reviews and less unique vocabulary performed best. We note that the RBF kernel SVM outperformed the linear kernel SVM. This means that the data analyzed was not linearly separable. We also highlight that other studies that utilize significantly smaller datasets found SVM to be the best performing SL classifier [21, 61]. Using our larger datasets we prove that ensemble learning methods such as RF, DT and AB perform best.

5-label Classification. Since the SL classifiers performed better when trained and tested on the HARD dataset for binary sentiment classification, the next step was to evaluate the classifiers for 5-label sentiment classification using the same dataset. The results in Table 6 and Figure 3 show the performance of SL classifiers on the task of 5-label sentiment classification. Similar to the binary classification case, Random Forest performed best, followed by Decision Tree and

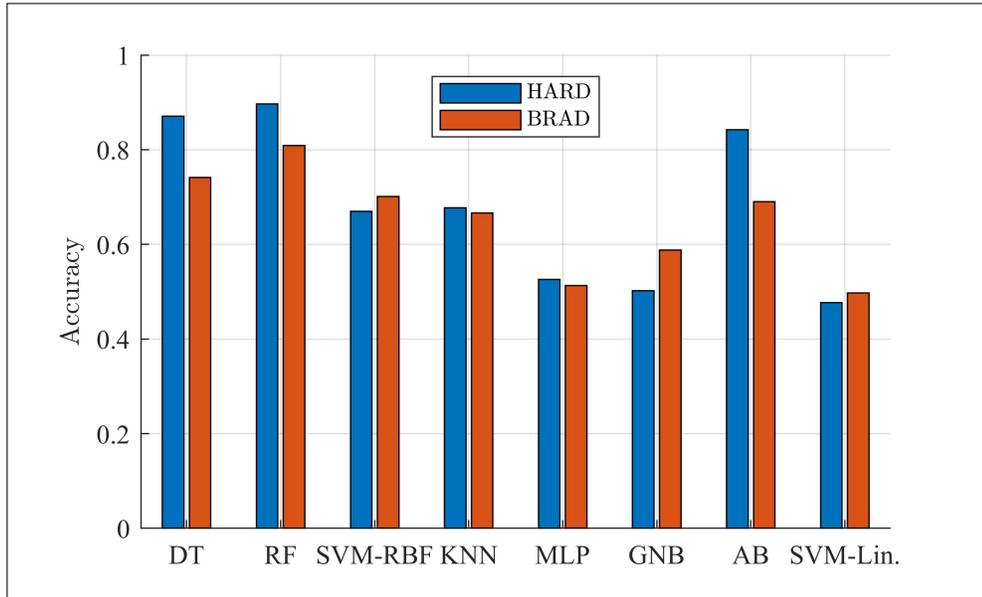


Fig. 2. Comparison between the binary classification accuracy obtained for SL techniques for both datasets.

Table 6. 5-label classification results from training and testing the HARD dataset using SL. Note: all results are rounded to three significant digits. **Bold** results represent the best performance from each category.

Classifier	Accuracy	Avg. Precision	Avg. Recall	Avg. F1	Avg. AUC
DT	0.664	0.660	0.660	0.660	0.790
RF	0.742	0.740	0.740	0.740	0.839
SVM (RBF)	0.390	0.500	0.390	0.370	0.617
KNN	0.392	0.390	0.390	0.390	0.619
MLP	0.220	0.200	0.210	0.140	0.509
GNB	0.196	0.240	0.200	0.080	0.500
AB	0.630	0.660	0.630	0.630	0.768
SVM (Linear)	0.223	0.230	0.220	0.200	0.513

AdaBoost. This further shows the superiority of these ensemble methods. The other classifiers performed considerably worse, with the GNB classifier performing worse than the baseline accuracy of 0.2.

4.2 Deep Learning

Binary Classification. The previously mentioned DL techniques were tested in a similar manner to the SL classifiers (see Table 7). Accuracy is seen to be around the mid-90-s for the HARD dataset and around 82% for the BRAD dataset. The results are fairly close for all techniques (see Figure 4). The best-performing classifier in terms of accuracy was the Bi-LSTM+CNN hybrid model for both the HARD and BRAD datasets.

To check if the most accurate model is statistically different from the other tested models, we conducted the non-parametric Wilcoxon test [39] between the classifier with the highest accuracy value (the Bi-LSTM+CNN hybrid model) and other tested classifiers. Based on the p -values reported in Table 7, we noticed that the Bi-LSTM+CNN classifier is

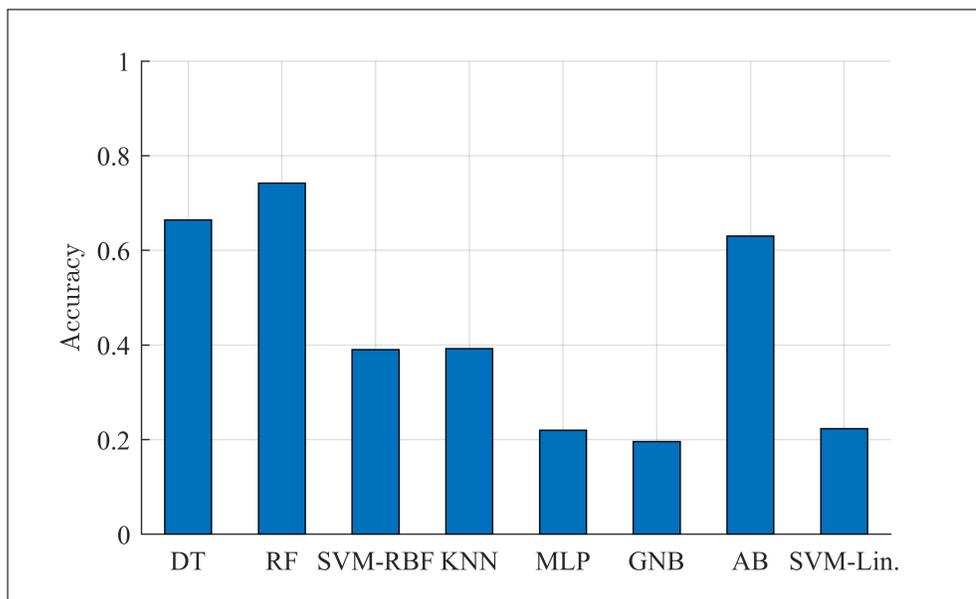


Fig. 3. Comparison between the 5-label classification accuracy obtained for shallow learning techniques for the HARD dataset.

Table 7. Binary classification results from training and testing the HARD (H) and BRAD (B) datasets using DL. Note: all results are rounded to three significant digits. **Bold** results represent the best performance from each category.

Classifier	Accuracy		Precision		Recall		F1		AUC		p -value
	H	B	H	B	H	B	H	B	H	B	H
CNN	0.939	0.824	0.933	0.820	0.947	0.824	0.939	0.822	0.939	0.824	0.0883
GRU	0.939	0.831	0.924	0.827	0.956	0.831	0.940	0.829	0.939	0.831	0.0342
Bi-LSTM	0.937	0.835	0.931	0.825	0.945	0.844	0.938	0.834	0.939	0.835	0.0008
CNN+GRU	0.936	0.815	0.934	0.804	0.939	0.826	0.937	0.815	0.937	0.814	0.0211
CNN+Bi-LSTM	0.940	0.812	0.933	0.804	0.949	0.816	0.941	0.810	0.940	0.813	0.2120
GRU+CNN	0.937	0.823	0.928	0.819	0.948	0.823	0.938	0.820	0.937	0.822	0.0045
GRU+Bi-LSTM	0.941	0.825	0.932	0.815	0.952	0.836	0.942	0.825	0.941	0.826	0.7623
Bi-LSTM+CNN	0.942	0.835	0.931	0.826	0.955	0.842	0.943	0.834	0.941	0.834	-
Bi-LSTM+GRU	0.923	0.830	0.904	0.821	0.949	0.838	0.926	0.829	0.924	0.830	0.0003
TRANS	0.922	0.813	0.918	0.801	0.928	0.825	0.923	0.813	0.923	0.813	0.0002

statistically different, i.e., having a p -value that is lower than 0.05, from 6 models based on a 95% confidence interval. We also observed that the Bi-LSTM+CNN classifier is statistically different from the CNN classifier at a 90% confidence interval. Thus we conclude that the Bi-LSTM+CNN hybrid classifier is statistically different from the majority of the tested DL classifiers.

5-label Classification. As the DL classifiers performed better when trained and tested on the HARD dataset for binary sentiment classification, the next step was to evaluate the classifiers for 5-label sentiment classification using the same dataset. The results in Table 8 and Figure 5 show the performance of DL classifiers on the task of 5-label sentiment classification. We find that the DL classifiers have similar performance as they did in the binary classification task, with

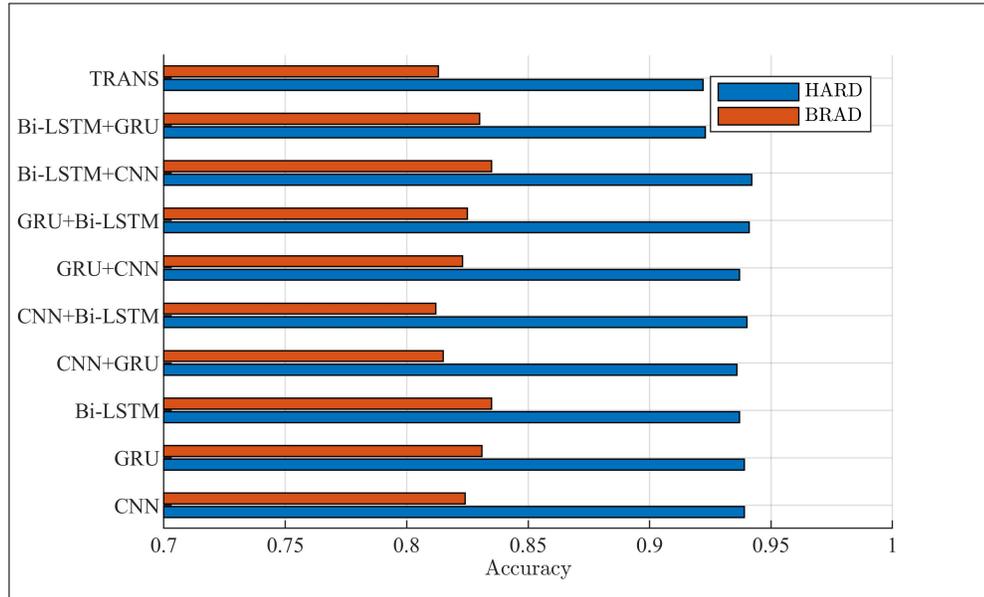


Fig. 4. Comparison between the binary classification accuracy obtained for deep learning techniques for both datasets.

Table 8. 5-label classification results from training and testing the HARD dataset using deep learning. Note: all results are rounded to three significant digits. **Bold** results represent the best performance from each category.

Classifier	Accuracy	Avg. Precision	Avg. Recall	Avg. F1	Avg. AUC
CNN	0.674	0.773	0.674	0.720	0.812
GRU	0.651	0.754	0.653	0.698	0.798
Bi-LSTM	0.631	0.725	0.632	0.674	0.785
CNN+GRU	0.649	0.754	0.648	0.697	0.797
CNN+Bi-LSTM	0.621	0.715	0.621	0.664	0.779
GRU+CNN	0.634	0.770	0.633	0.697	0.793
GRU+Bi-LSTM	0.653	0.713	0.655	0.681	0.793
Bi-LSTM+CNN	0.634	0.727	0.635	0.678	0.787
Bi-LSTM+GRU	0.667	0.709	0.667	0.687	0.799
TRANS	0.603	0.683	0.602	0.640	0.766

classification accuracy being measured at around 0.65 and with CNN having the highest accuracy at 0.674. However, we note that while the DL classifiers outperform most of their SL counterparts for the task of 5-label classification, the RF classifier achieved higher accuracy at 0.742. This point will be discussed in greater detail in the following section.

4.3 Learning Curve

Binary Classification. An experiment was conducted to understand the effect of dataset size on the accuracy of the best-performing DL and SL models, i.e., Bi-LSTM+CNN and RF. This was done by training and testing the best-performing techniques using the HARD dataset at varying dataset sizes (see Table 9). The changes in accuracy can then be captured and compared. To ensure the true behavior of the techniques is observed, 100 random samples were extracted from the

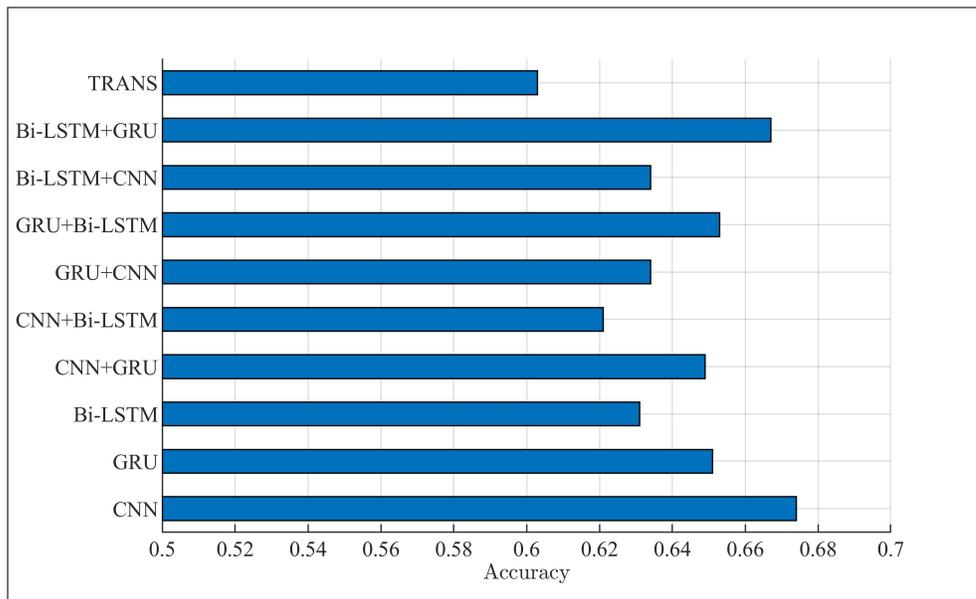


Fig. 5. Comparison of classification accuracy in 5-label sentiment analysis obtained for deep learning techniques for the HARD dataset.

Table 9. The performance of Bi-LSTM+CNN and RF versus dataset size variation for binary classification of HARD. Accuracy is averaged over 100 runs, where each run randomly resamples the dataset.

Dataset Size	Number of Reviews	Bi-LSTM+CNN		RF	
		Mean Accuracy	Standard Deviation	Mean Accuracy	Standard Deviation
1%	910	0.895	0.0327	0.921	0.0248
5%	4,550	0.915	0.0140	0.919	0.0121
10%	9,100	0.926	0.0095	0.913	0.0104
50%	45,500	0.938	0.0038	0.902	0.0046
100%	91,000	0.942	0.0026	0.899	0.0035

dataset and split by (80-10-10) for training, testing and validation (note that cross-validation was used for RF). The models were trained using the split data and the mean accuracy and standard deviation values for the 100 runs were recorded in Table 9.

From Figure 6, it can be seen that the mean accuracy of the models used is affected by dataset size variation. As the utilized dataset size increased, so did the mean accuracy of the DL models. On the other hand, the SL classifier saw a gradual decrease in mean accuracy with dataset size increase. Furthermore, the standard deviation value for both classifiers is seen to decrease as the number of reviews increases. An important feature to note in this case is the percentage of the dataset needed for the mean accuracy of Bi-LSTM+CNN (a DL technique) to overtake the mean accuracy of RF (a SL technique), which is around 6% of the dataset, corresponding to 5,460 reviews. Note that this value

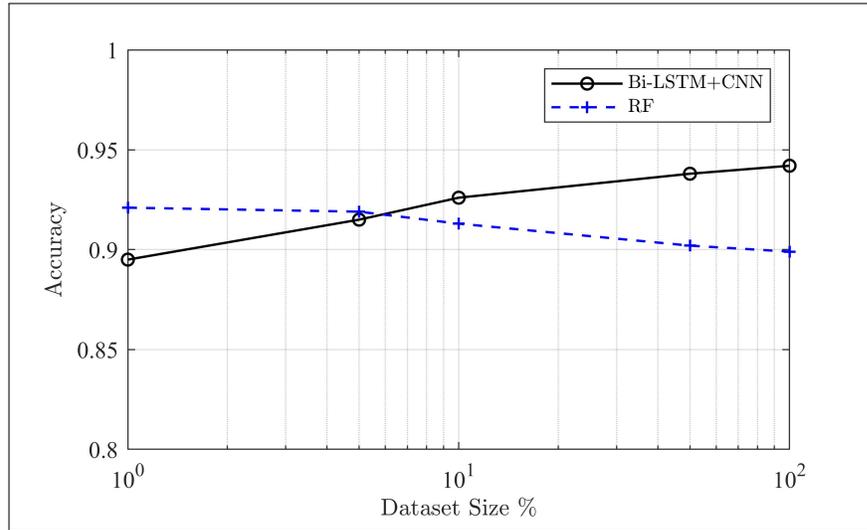


Fig. 6. Binary classification accuracy versus the percentage of the dataset used for Bi-LSTM+CNN and RF. Accuracy is averaged over 100 runs, where each run randomly resamples the dataset. Note that the x-axis is in the logarithmic scale.

is highly dependent on the quality of the dataset used. The aforementioned feature is important in the argument for the use of DL techniques instead of SL methods for improved accuracy [56].

Multi-label Classification. A similar experiment was conducted for the case of multi-label classification instead of binary classification. This was also done using the HARD dataset, but instead of classifying positive/negative reviews (see Pre-processing, Step 3) we included all rating values from 1 to 5. Thus the problem became a 5-label classification problem. Similar to the previous experiment, a comparison was made between the best-performing deep and SL models, based on the results of Sections 4.1 and 4.2; these were the CNN and RF classifiers. To ensure a fair comparison between the various labels, the dataset was balanced around the smallest label, which was the lowest rating of 1. Since rating 1 contained around 12,500 reviews, random samples of the same size (12,500) were taken from each label, adding up to a full dataset of 62,500 reviews. The data was then split (80-10-10) for training, testing and validation (note that cross-validation was used for RF). Finally the mean accuracy and standard deviation values were obtained from 100 runs (where each run randomly samples the dataset before it is split) and recorded in Table 10.

Figure 7 presents the mean accuracy trend for both CNN and RF for a varying dataset size. We can see that the RF classifier outperforms the CNN classifier for all dataset sizes. Furthermore, similar to Figure 6, it can be seen that as dataset size increases, so does the accuracy of CNN, whereas the performance of the RF classifiers remains consistent. Therefore, the CNN classifier will benefit from a larger dataset, allowing it to overtake the RF classifier in terms of classification accuracy.

4.4 Embeddings

In the previous sections we have seen the TRANS model underperforming when compared to other DL architectures such as the Bi-LSTM+CNN hybrid model for binary classification and the CNN model for 5-label classification. This does not match recent findings in the literature [30]. We believe this is due to the fact that default single-layer token

Table 10. The performance of CNN and RF versus dataset size variation for 5-label classification of HARD. Accuracy is averaged over 100 runs, where each run randomly resamples the dataset.

Dataset Size	Number of Reviews	CNN		RF	
		Mean Accuracy	Standard Deviation	Mean Accuracy	Standard Deviation
1%	625	0.481	0.1113	0.748	0.0538
5%	3,125	0.587	0.0457	0.777	0.0238
10%	6,250	0.615	0.0309	0.775	0.0194
50%	31,250	0.656	0.0152	0.750	0.0074
100%	62,500	0.668	0.0120	0.743	0.0060

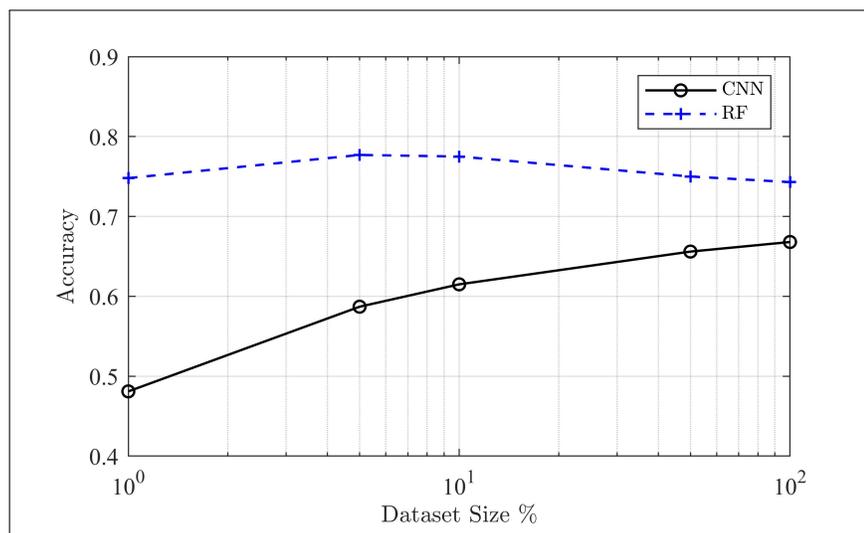


Fig. 7. 5-label classification accuracy versus the percentage of the dataset used for CNN and RF. Accuracy is averaged over 100 runs, where each run randomly resamples the dataset. Note that the x-axis is in the logarithmic scale.

embeddings were used (vector size 64), whereas the TRANS model benefits greatly from more complex embedding layers.

One of the main differences between SL and DL for NLP is the presence of an embedding layer in the front end of DL techniques and the lack thereof in SL techniques. For word embeddings, several options could be used, either learning word embeddings from the problem or reusing embeddings that are publicly available. Learning word embeddings from the problem can be done using three main methods. A model can learn the embeddings separately through a specific framework such as Word2Vec [53], fastText [40], GloVe [57] or the more recent araBERT [18]. Alternatively, it can learn embedding jointly using an initial embedding layer or it can use sentiment specific embeddings [12]. We chose to use the second technique in this work (learning embedding jointly using an initial embedding layer) and compare it with state-of-the-art araBERT model (v0.1). The authors in [18] obtained araBERT as a result of pre-training BERT for the Arabic language. They also found that it outperforms multilingual BERT (mBERT) [27] for SA performed on the HARD dataset. We are also comparing the default embedding layers with a custom embedding layer, utilizing

Table 11. The effect of embedding testing on 3K sentences for binary classification and 5K sentences for 5-label classification. The accuracy recorded is the mean accuracy measured over 100 re-runs, where each re-run randomly resamples the dataset.

Embedding Type	Binary Accuracy	5-Label Accuracy
Single Layer Default Token	0.892	0.609
Single Layer Custom Token	0.894	0.616
Two Layer Default Token + Position	0.916	0.638
Two Layer Custom Token + Position	0.914	0.650
araBERT	0.927	0.736

trigram multi-dialect embeddings taken from several websites [29, 32]. We used the TRANS architecture to compare the following:

- Using a single default embedding layer for tokens
- Using a single custom embedding layer for tokens
- Using two default embedding layers, one for tokens and another for token positions
- Using a custom embedding layer for tokens and a default embedding layer for token positions
- Using araBERT

The results can be seen in Table 11 for binary and 5-label SA using the HARD dataset. We note that the chosen embedding vector size was 300. The results show the superiority of the araBERT model for binary and multi-label classification. Furthermore, we note the fact that using a two-layer token and position embedding had a greater improvement in performance than using a single-token embedding layer. We also note that the custom embedding setup outperformed its default counterpart for multi-label classification. This could mean that the custom token embedding layer exploits label-specific phrases more effectively.

4.5 Misclassifications

Although the DL classifiers performed exceptionally well it is prudent to find the misclassifications of the model to establish a better understanding of the limitations of the model used and what could further be done to improve the classification. In Table 12, six misclassification examples of the CNN model are shown, including the review sentences in Arabic and their literal translation in English, along with their true rating and the prediction of the DL models. Keep in mind that these sentences are how they appear after the cleaning process, right before tokenization and padding.

The first case presents the important issue of mislabeling by the initial reviewer, where the sentence contains no positive aspects, however, the sentiment is labeled as positive. These errors can only be prevented by vetting the dataset used to remove all mislabeled data, which is not an easy process.

The second case presents a positively labeled sentiment that focuses on a single negative aspect of the review, in this instance, the Wi-Fi. Even though the user praises the hotel in the first part of the review, it is only a small part, around 20% of the total review. The reviewer’s focus on the negative aspect is why we believe the model misclassified this review.

The third case presents a review that highlights the shortcomings of the model used, as its prediction was based on the only positive feature in the sentence being the first word, which was misspelled.

The fourth case is the opposite of the second case, wherein the reviewer highlights the positive aspects for the majority of the sentence and refers to the negative aspect right at the end, which is the sentiment-deciding factor for the reviewer. The model misclassified this review as a positive review when in fact it is negative.

Table 12. Examples of misclassifications by the CNN model (binary classification).

#	Review	Rating	Model Prediction
1	رحله قصيره الواي فاي سي قوه الما في دورات المياه ضعيف جدا Short trip bad Wi-Fi power of water in washroom weak	1	0
2	استثنائي كل شي ممتاز وراقي من جميع النواحي عدا الواي فاي كلش ضعيف الواي فاي ضعيف جدا ولازم تعيد ربط الاتصال حتي اذا تنقلت من غرفه الي اخري في نفس الشقه وكذلك لازم تعيد الاتصال ويطلب كلمه السر في كل مره تغادر وترجع الشقه Exceptional everything is excellent and elegant in all aspects except for the Wi-Fi extremely weak Wi-Fi is very weak and you must reconnect even if you move between rooms in the same apartment and you must reconnect and it asks for password every time you leave and return to the apartment	1	0
3	جيد ولكن نبحث عن الافضل الموقع الخدمات لم تكن بالمستوي المطلوب مثل توفير مسلتزمات دورات المياه ايضا ارجو الاهتمام بالخدمات الفندقية المقدمه وشكرا Good but we are looking for the best the location the services were not up to par such as providing washroom amenities also please pay attention to the provided hotel services and thank you	0	1
4	الكويت مكان ممتاز عايلي ممتاز في الهدو والخدمه المتمازه والنظافه وخدمه الضيافه القهوه العربيه متازه تاخير في استلام الغرفه الحيك ان Kuwait is an excellent place familial excels in quietness and excellent service and cleanliness and hospitality Arabic coffee excellent late in receiving the room check in	0	1
5	ارجو منهم ان يحافظوا علي نظافه الفندق والاهتمام بالخدمه الفطار متنوع والخدمه داخل بوفيه الفطار ممتازين لقد قمت بالحجر مبكرا وعندما قدمت قال لي ناسف لا يمكننا ان نوفر لك سرير مزدوج فاعطاني سريرين منفصلين معنا هذه المره الثانيه لي اسكن في هذا الفندق فالخدمه بالمره الاولي كانت افضل بسبب عدم الزحمه بالفندق I implore them to keep the hotel clean and pay attention to the service breakfast was diverse and service at the breakfast buffet excellent I made reservations early and when I arrived he said we apologize we cannot provide you with a double bed so he gave me two separate beds even though this is the second time I live in this hotel the service the first time was better due to less congestion in the hotel	0	1
6	مثال واضح لموظفين رايعين في مكان غير مناسب طاقم عمل ممتاز جدا الاصوات حول الفندق وداخل الفندق سيه جدا لا يوجد فيش كهربا في مكان مناسب لاستخدام الحاسب الشخصي انترنت سيه جدا A clear example of fantastic service in unsuitable location the work crew very excellent noises around the hotel and in the hotel very bad unavailability of power plug in a suitable place to use personal computer internet very bad	1	0

Table 13. Examples of correct classifications of the CNN model (binary classification).

#	Review	Rating	Model Prediction
1	ضعيف الغرفه مزعجه السرير غير نظيف Weak the room annoying the bed not clean	0	0
2	مخيب للامل النظافه الخدمه الطلبات الفطور ملغي الانترنت معدوم في الغرف التلفزيون خربان ما نسمع الاذان للصلاه لا ارغب هذا الفندق في الرحلات القادمه وشكرا Disappointing cleanliness service requests breakfast canceled internet non-existent in rooms TV is non-functional we do not hear call to prayer I do not want to stay at this hotel in my next trips and thank you	0	0
3	رائع السكن مقارنه بالسعر ممتاز المكان حق نوم يعني تمام مرتاح نظافه وهدو مافي اي وساخه حتي يجيبون لك غسله لو مطول المكان تحفه مقارنه بالسعر Fantastic the residence comparing price excellent place for sleeping meaning you sleep comfortable quietness and cleanliness no dirt they even bring you a washing machine if staying long place is masterpiece comparing price	1	1
4	جميل وهادي وقريب من الخدمات وبه مسبح روعه وجيم واي فاي كل شي جميل تقريبا نظافه وهدو ومسبح رائع ماعدا اشيا صغيره المماطله في توفير بعض المستلزمات Beautiful and quiet and near services and has fantastic swimming pool and gym Wi-Fi everything beautiful approximately cleanliness and quietness and fantastic swimming pool except small things procrastination in providing some amenities	1	1

The fifth case introduces a detailed review that highlights the positive feature (breakfast), with the mention of the word (ممتازين), meaning excellent. It however fails to emphasize the negative aspects, i.e., no mention of the words (bad, poor, negative). Subsequently, the model misclassified the sentiment of the sentence and confused it with a positive sentiment.

The sixth case combined both positive and negative aspects, where the reviewer labeled this review as positive the model classified it as negative.

In Table 13, four correct classification examples of the CNN model are shown, including the review sentences in Arabic and their literal translation in English, along with their true rating and the prediction of the DL models. We can see that these sentences clearly state the positive or negative aspects and emphasize the sentiment with words such as (رائع) which is the normalized version of (رائع), meaning fantastic, and (ضعيف) meaning weak.

5 CONCLUSION

In this study, a thorough comparison was conducted between deep and SL techniques for Arabic SA. This comparison is the first of its kind, as previous works only consider a fraction of the techniques proposed in this work. The three main archetypes of DL (CNN, LSTM and GRU), as well as all their possible hybrids, were evaluated against eight of

the most-widely used SL classifiers. Furthermore, some state-of-the-art techniques were included such as the TRANS architecture as well as the araBERT pre-trained model. The two datasets used are some of the largest available: the HARD and BRAD datasets. DL techniques outperformed SL techniques in all cases except in 5-label classification, where RF performed best. Furthermore, both techniques performed best when trained and tested using the HARD dataset as compared to the BRAD dataset, even though BRAD was 33% larger. The reasoning behind this is attributed to the significantly large vocabulary space in BRAD, which was more than twice that of HARD. The best-performing SL techniques were the ensemble learning methods, Random Forest, Decision Tree, and AdaBoost, whereas all the DL techniques performed somewhat similarly.

These results show that using DL techniques is more favorable to SL methods for SA of Arabic reviews. This clashes with what was proposed in [7], where it was found that SVM outperformed CNN and RNN in terms of binary sentiment classification. The performance of DL methods improved in our experiments due to the use of larger datasets for training and testing the DL models. It is observed that there is a direct relationship between the performance of DL techniques and the size of the dataset used. To verify this claim, we plotted the learning curve as a function of dataset size. While the increase in accuracy of the DL models diminishes as the dataset increases, it still has a significant effect on the mean accuracy obtained from multiple runs.

Thus we can conclude that for large Arabic datasets (containing tens of thousands of sentences), the use of DL techniques for SA is optimal, whereas for small datasets (less than 10,000 sentences), using a SL technique would be preferred. It must also be highlighted that the TRANS architecture benefits greatly from the use of more intricate embedding layers, particularly pre-trained models such as araBERT. With these advanced embeddings, the TRANS model is able to outperform all other deep and SL techniques for all sizes of datasets except for those that are very small (less than 3,000 sentences).

We believe that further treatment of Arabic dialects could enhance the performance of DL techniques for Arabic SA applications. While the performance achieved by the DL classifiers was exemplary, further work can be done to increase the accuracy from the mid-90s to the high-90s.

REFERENCES

- [1] Mahmoud Al-Ayyoub, Abed Allah Khamaiseh, Yaser Jararweh, and Mohammed N Al-Kabi. 2019. A comprehensive survey of arabic sentiment analysis. *Information Processing & Management* 56, 2 (2019), 320–342.
- [2] Mahmoud Al-Ayyoub, Aya Nuseir, Kholoud Alsmearat, Yaser Jararweh, and Brij Gupta. 2018. Deep learning for Arabic NLP: A survey. *Journal of computational science* 26 (2018), 522–531.
- [3] Sadam Al-Azani and El-Sayed El-Alfy. 2018. Emojis-based sentiment classification of Arabic microblogs using deep recurrent neural networks. In *2018 International Conference on Computing Sciences and Engineering (ICCSE)*. IEEE, 1–6.
- [4] Sadam Al-Azani and El-Sayed M El-Alfy. 2017. Hybrid deep learning for sentiment polarity determination of arabic microblogs. In *International Conference on Neural Information Processing*. Springer, 491–500.
- [5] Ahmad Al-Sallab, Ramy Baly, Hazem Hajj, Khaled Bashir Shaban, Wassim El-Hajj, and Gilbert Badaro. 2017. Aroma: A recursive deep learning model for opinion mining in arabic as a low resource language. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 16, 4 (2017), 1–20.
- [6] Ahmad Al Sallab, Hazem Hajj, Gilbert Badaro, Ramy Baly, Wassim El Hajj, and Khaled Bashir Shaban. 2015. Deep learning models for sentiment analysis in Arabic. In *Proceedings of the second workshop on Arabic natural language processing*. 9–17.
- [7] Mohammad Al-Smadi, Omar Qawasmeh, Mahmoud Al-Ayyoub, Yaser Jararweh, and Brij Gupta. 2018. Deep Recurrent neural network vs. support vector machine for aspect-based sentiment analysis of Arabic hotels' reviews. *Journal of computational science* 27 (2018), 386–393.
- [8] Abdulaziz M Alayba, Vasile Palade, Matthew England, and Rahat Iqbal. 2018. A combined CNN and LSTM model for arabic sentiment analysis. In *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*. Springer, 179–191.
- [9] Huda Abdulrahman Almuzaini and Aqil M Azmi. 2020. Impact of stemming and word embedding on deep learning-based arabic text categorization. *IEEE Access* 8 (2020), 127913–127928.

- [10] Anwar Alnawas and Nursal Arici. 2019. Sentiment Analysis of Iraqi Arabic Dialect on Facebook Based on Distributed Representations of Documents. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 18, 3 (2019), 20.
- [11] Naomi S Altman. 1992. An introduction to kernel and nearest-neighbor nonparametric regression. *The American Statistician* 46, 3 (1992), 175–185.
- [12] A Aziz Altowayan and Ashraf Elnagar. 2017. Improving Arabic sentiment analysis with sentiment-specific embeddings. In *2017 IEEE International Conference on Big Data (Big Data)*. IEEE, 4314–4320.
- [13] Gilbert Badaro, Ramy Baly, Rana Akel, Linda Fayad, Jeffrey Khairallah, Hazem Hajj, Khaled Shaban, and Wassim El-Hajj. 2015. A light lexicon-based mobile application for sentiment mining of arabic tweets. In *Proceedings of the second workshop on arabic natural language processing*. 18–25.
- [14] Gilbert Badaro, Ramy Baly, Hazem Hajj, Wassim El-Hajj, Khaled Bashir Shaban, Nizar Habash, Ahmad Al-Sallab, and Ali Hamdi. 2019. A survey of opinion mining in Arabic: A comprehensive system perspective covering challenges and advances in tools, resources, models, applications, and visualizations. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 18, 3 (2019), 27.
- [15] Gilbert Badaro, Ramy Baly, Hazem Hajj, Nizar Habash, and Wassim El-Hajj. 2014. A large scale Arabic sentiment lexicon for Arabic opinion mining. In *Proceedings of the EMNLP 2014 workshop on arabic natural language processing (ANLP)*. 165–173.
- [16] Gilbert Badaro, Obeida El Jundi, Alaa Khaddaj, Alaa Maarouf, Raslan Kain, Hazem Hajj, and Wassim El-Hajj. 2018. Ema at semeval-2018 task 1: Emotion mining for arabic. In *Proceedings of The 12th International Workshop on Semantic Evaluation*. 236–244.
- [17] Gilbert Badaro, Hazem Hajj, and Nizar Habash. 2020. A Link Prediction Approach for Accurately Mapping a Large-scale Arabic Lexical Resource to English WordNet. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 19, 6 (2020), 1–38.
- [18] Fady Baly, Hazem Hajj, et al. 2020. AraBERT: Transformer-based Model for Arabic Language Understanding. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*. 9–15.
- [19] Ramy Baly, Gilbert Badaro, Georges El-Khoury, Rawan Moukalled, Rita Aoun, Hazem Hajj, Wassim El-Hajj, Nizar Habash, and Khaled Shaban. 2017. A characterization study of arabic twitter data with a benchmarking for state-of-the-art opinion mining models. In *Proceedings of the third Arabic natural language processing workshop*. 110–118.
- [20] Ramy Baly, Gilbert Badaro, Ali Hamdi, Rawan Moukalled, Rita Aoun, Georges El-Khoury, Ahmad Al Sallab, Hazem Hajj, Nizar Habash, Khaled Shaban, et al. 2017. Omam at semeval-2017 task 4: Evaluation of english state-of-the-art sentiment analysis models for arabic and a new topic-based model. In *Proceedings of the 11th international workshop on semantic evaluation (SEMEVAL-2017)*. 603–610.
- [21] Ramy Baly, Hazem Hajj, Nizar Habash, Khaled Bashir Shaban, and Wassim El-Hajj. 2017. A sentiment treebank and morphologically enriched recursive deep models for effective sentiment analysis in arabic. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 16, 4 (2017), 23.
- [22] Majdi Beseiso and Haytham Elmousalami. 2020. Subword attentive model for Arabic sentiment analysis: A deep learning approach. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 19, 2 (2020), 1–17.
- [23] Naaima Boudad, Rdouan Faizi, Rachid Oulad Haj Thami, and Raddouane Chiheb. 2017. Sentiment analysis in Arabic: A review of the literature. *Ain Shams Engineering Journal* (2017).
- [24] Leo Breiman. 2001. Random forests. *Machine learning* 45, 1 (2001), 5–32.
- [25] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259* (2014).
- [26] Jesse Davis and Mark Goadrich. 2006. The relationship between Precision-Recall and ROC curves. In *Proceedings of the 23rd international conference on Machine learning*. ACM, 233–240.
- [27] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [28] Hai Ha Dohaiha, PWC Prasad, Angelika Maag, and Abeer Alsadoon. 2018. Deep learning for aspect-based sentiment analysis: a comparative review. *Expert Systems With Applications* (2018).
- [29] Omar Einea, Ashraf Elnagar, and Ridhwan Al Debsi. 2019. Sanad: Single-label arabic news articles dataset for automatic text categorization. *Data in brief* 25 (2019), 104076.
- [30] Obeida ElJundi, Wissam Antoun, Nour El Droubi, Hazem Hajj, Wassim El-Hajj, and Khaled Shaban. 2019. hulmona: The universal language model in arabic. In *Proceedings of the Fourth Arabic Natural Language Processing Workshop*. 68–77.
- [31] Ashraf Elnagar. 2016. Investigation on sentiment analysis for Arabic reviews. In *2016 IEEE/ACS 13th International Conference of Computer Systems and Applications (AICCSA)*. IEEE, 1–7.
- [32] Ashraf Elnagar, Ridhwan Al-Debsi, and Omar Einea. 2020. Arabic text classification using deep learning models. *Information Processing & Management* 57, 1 (2020), 102121.
- [33] Ashraf Elnagar and Omar Einea. 2016. BRAD 1.0: Book reviews in Arabic dataset. In *2016 IEEE/ACS 13th International Conference of Computer Systems and Applications (AICCSA)*. 1–8.
- [34] Ashraf Elnagar, Yasmin S Khalifa, and Anas Einea. 2018. Hotel Arabic-reviews dataset construction for sentiment analysis applications. In *Intelligent Natural Language Processing: Trends and Applications*. Springer, 35–52.
- [35] Ashraf Elnagar, Leena Lulu, and Omar Einea. 2018. An Annotated Huge Dataset for Standard and Colloquial Arabic Reviews for Subjective Sentiment Analysis. *Procedia computer science* 142 (2018), 182–189.
- [36] Ashraf Elnagar, Sane M Yagi, Ali Bou Nassif, Ismail Shahin, and Said A Salloum. 2021. Systematic Literature Review of Dialectal Arabic: Identification and Detection. *IEEE Access* 9 (2021), 31010–31042.

- [37] Yoav Freund and Robert E Schapire. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of computer and system sciences* 55, 1 (1997), 119–139.
- [38] Yoav Freund, Robert E Schapire, et al. 1996. Experiments with a new boosting algorithm. In *icml*, Vol. 96. Citeseer, 148–156.
- [39] Edmund A Gehan. 1965. A generalized Wilcoxon test for comparing arbitrarily singly-censored samples. *Biometrika* 52, 1-2 (1965), 203–224.
- [40] Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning Word Vectors for 157 Languages. In *Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)*.
- [41] Klaus Greff, Rupesh K Srivastava, Jan Koutník, Bas R Steunebrink, and Jürgen Schmidhuber. 2016. LSTM: A search space odyssey. *IEEE transactions on neural networks and learning systems* 28, 10 (2016), 2222–2232.
- [42] Ali Hamdi, Khaled Shaban, and Anazida Zainal. 2018. Clasenti: A class-specific sentiment analysis framework. *ACM Transactions on Asian and Low-Resource Language Information Processing (TALLIP)* 17, 4 (2018), 1–28.
- [43] Sepp Hochreiter and Jürgen Schmidhuber. 1996. Bridging long time lags by weight guessing and “Long Short-Term Memory”. *Spatiotemporal models in biological and artificial systems* 37, 65-72 (1996), 11.
- [44] Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. 2015. An empirical exploration of recurrent network architectures. In *International Conference on Machine Learning*. 2342–2350.
- [45] Jawad Khan, Aftab Alam, Jamil Hussain, and Young-Koo Lee. 2019. EnSWF: effective features extraction and selection in conjunction with ensemble learning methods for document sentiment classification. *Applied Intelligence* (2019), 1–23.
- [46] Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882* (2014).
- [47] Yann LeCun, D Touresky, G Hinton, and T Sejnowski. 1988. A theoretical framework for back-propagation. In *Proceedings of the 1988 connectionist models summer school*, Vol. 1. CMU, Pittsburgh, Pa: Morgan Kaufmann, 21–28.
- [48] Lin Li, Tiong-Thye Goh, and Dawei Jin. 2018. How textual quality of online reviews affect classification performance: a case of deep learning sentiment analysis. *Neural Computing and Applications* (2018), 1–29.
- [49] Andy Liaw, Matthew Wiener, et al. 2002. Classification and regression by randomForest. *R news* 2, 3 (2002), 18–22.
- [50] Weibo Liu, Zidong Wang, Xiaohui Liu, Nianyin Zeng, Yurong Liu, and Fuad E Alsaadi. 2017. A survey of deep neural network architectures and their applications. *Neurocomputing* 234 (2017), 11–26.
- [51] Leena Lulu and Ashraf Elnagar. 2018. Automatic Arabic Dialect Classification Using Deep Learning Models. *Procedia Computer Science* 142 (2018), 262 – 269.
- [52] Gjorgji Madjarov, Dragi Koccev, Dejan Gjorgjevič, and Sašo Džeroski. 2012. An extensive experimental comparison of methods for multi-label learning. *Pattern recognition* 45, 9 (2012), 3084–3104.
- [53] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [54] Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In *Proceedings of the 12th international workshop on semantic evaluation*. 1–17.
- [55] Ali Bou Nassif, Ashraf Elnagar, Ismail Shahin, and Safaa Henno. 2020. Deep learning for Arabic subjective sentiment analysis: Challenges and research opportunities. *Applied Soft Computing* (2020), 106836.
- [56] Andrew Ng. 2015. What data scientists should know about deep learning. URL <https://www.slideshare.net/ExtractConf> 44 (2015).
- [57] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1532–1543.
- [58] F Rosenblatt. 1960. *Perceptrons and the Theory of Brain Mechanisms: Cornell Aeronautical Laboratory*. Technical Report. Report No. VG-1196-G-8.
- [59] S Rasoul Safavian and David Landgrebe. 1991. A survey of decision tree classifier methodology. *IEEE transactions on systems, man, and cybernetics* 21, 3 (1991), 660–674.
- [60] Jürgen Schmidhuber. 2015. Deep learning in neural networks: An overview. *Neural networks* 61 (2015), 85–117.
- [61] Amira Shoukry and Ahmed Rafea. 2012. Sentence-level Arabic sentiment analysis. In *2012 International Conference on Collaboration Technologies and Systems (CTS)*. IEEE, 546–550.
- [62] Mohammad Tubishat, Mohammad AM Abushariah, Norisma Idris, and Ibrahim Aljarah. 2019. Improved whale optimization algorithm for feature selection in Arabic sentiment analysis. *Applied Intelligence* 49, 5 (2019), 1688–1707.
- [63] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.

6 APPENDIX

After various iterations of trial and error, the following configurations consistently performed best and thus were used as the basis of this study.

Table 14. TRANS

Layer (type)	Output Shape
Embedding	(None, 100, 64)
Transformer	(None, 100, 64)
Global Max Pooling	(None, 64)
Dropout 1	(None, 64)
Dense 1	(None, 32)
Dropout 2	(None, 32)
Dense 2	(None, 1)

Table 15. CNN

Layer (type)	Output Shape
Embedding	(None, 100, 64)
Convolutional 1 (1-D)	(None, 100, 300)
Convolutional 2 (1-D)	(None, 100, 64)
Convolutional 3 (1-D)	(None, 100, 32)
Global Max Pooling	(None, 32)
Dropout	(None, 32)
Dense 1	(None, 32)
Dense 2	(None, 1)

Table 16. LSTM

Layer (type)	Output Shape
Embedding	(None, 100, 64)
Bidirectional LSTM	(None, 64)
Dense 1	(None, 32)
Dense 2	(None, 1)

Table 17. GRU

Layer (type)	Output Shape
Embedding	(None, 100, 64)
GRU	(None, 100, 64)
Global Max Pooling	(None, 64)
Dense 1	(None, 32)
Dense 2	(None, 1)

Table 18. CNN+LSTM

Layer (type)	Output Shape
Embedding	(None, 100, 64)
Convolutional 1 (1-D)	(None, 100, 64)
Convolutional 2 (1-D)	(None, 100, 32)
Bidirectional LSTM	(None, 32)
Dense	(None, 1)

Table 19. LSTM+CNN

Layer (type)	Output Shape
Embedding	(None, 100, 64)
Bidirectional LSTM	(None, 100, 64)
Convolutional 1 (1-D)	(None, 100, 64)
Convolutional 2 (1-D)	(None, 100, 32)
Global Max Pooling	(None, 32)
Dropout	(None, 32)
Dense	(None, 1)

Table 20. GRU+CNN

Layer (type)	Output Shape
Embedding	(None, 100, 64)
GRU	(None, 100, 64)
Convolutional 1 (1-D)	(None, 100, 64)
Convolutional 2 (1-D)	(None, 100, 32)
Global Max Pooling	(None, 32)
Dropout	(None, 32)
Dense	(None, 1)

Table 21. CNN+GRU

Layer (type)	Output Shape
Embedding	(None, 100, 64)
Convolutional 1 (1-D)	(None, 100, 64)
Convolutional 2 (1-D)	(None, 100, 32)
GRU	(None, 100, 32)
Global Max Pooling	(None, 32)
Dense	(None, 1)

Table 22. LSTM+GRU

Layer (type)	Output Shape
Embedding	(None, 100, 64)
Bidirectional LSTM	(None, 100, 64)
GRU	(None, 32)
Dense 1	(None, 1)

Table 23. GRU+LSTM

Layer (type)	Output Shape
Embedding	(None, 100, 64)
GRU	(None, 100, 64)
Bidirectional LSTM	(None, 32)
Dense 1	(None, 1)

Table 24. Shallow Learning Classifiers

Name	Parameters
Decision Tree	ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features=None, max_leaf_nodes=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, presort='deprecated', random_state=None, splitter='best'
Random Forest	bootstrap=True, ccp_alpha=0.0, class_weight=None, criterion='gini', max_depth=None, max_features='auto', max_leaf_nodes=None, max_samples=None, min_impurity_decrease=0.0, min_impurity_split=None, min_samples_leaf=1, min_samples_split=2, min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=None, oob_score=False, random_state=None, verbose=0, warm_start=False
SVM (RBF)	C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='scale', kernel='rbf', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False
SVM (Linear)	C=1.0, class_weight=None, dual=True, fit_intercept=True, intercept_scaling=1, loss='squared_hinge', max_iter=1000, multi_class='ovr', penalty='l2', random_state=None, tol=0.0001, verbose=0
KNN	algorithm='auto', leaf_size=30, metric='minkowski', metric_params=None, n_jobs=None, n_neighbors=5, p=2, weights='uniform'
MLP	activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9, beta_2=0.999, early_stopping=False, epsilon=1e-08, hidden_layer_sizes=(100,), learning_rate='constant', learning_rate_init=0.001, max_fun=15000, max_iter=200, momentum=0.9, n_iter_no_change=10, nesterovs_momentum=True, power_t=0.5, random_state=None, shuffle=True, solver='adam', tol=0.0001, validation_fraction=0.1, verbose=False, warm_start=False
Gaussian NB	priors=None, var_smoothing=1e-09
AdaBoost	algorithm='SAMME.R', base_estimator=None, learning_rate=1.0, n_estimators=50, random_state=None