

Investigating the Effect of Preprocessing Arabic Text on Offensive Language and Hate Speech Detection

FATEMAH HUSAIN, Kuwait University, Kuwait

OZLEM UZUNER, George Mason University, USA

Preprocessing of input text can play a key role in text classification by reducing dimensionality and removing unnecessary content. This study aims to investigate the impact of preprocessing on Arabic offensive language classification. We explore five preprocessing techniques: conversion of emojis to Arabic textual labels, normalization of different forms of Arabic letters, normalization of selected nouns from dialectal Arabic to Modern Standard Arabic (MSA), conversion of selected hyponyms to hypernyms, hashtag segmentation, and basic cleaning such as removing numbers, kashidas, diacritics, and HTML tags. We also experiment with raw text and a combination of all five preprocessing techniques. We apply different types of classifiers in our experiments including traditional machine learning, ensemble machine learning, Artificial Neural Networks (ANN), and Bidirectional Encoder Representations from Transformers (BERT)-based models to analyze the impact of preprocessing. Our results demonstrate significant variations in the effects of preprocessing on each classifier type and on each dataset. Classifiers that are based on BERT do not benefit from preprocessing, while traditional machine learning classifiers do. However, these results can benefit from validation on larger datasets that cover broader domains and dialects.

CCS Concepts: • **Computing methodologies** → **Information extraction; Natural language processing.**

Additional Key Words and Phrases: artificial neural networks, offensive language detection, natural language processing, Arabic language, machine learning

ACM Reference Format:

Fatemah Husain and Ozlem Uzuner. 2020. Investigating the Effect of Preprocessing Arabic Text on Offensive Language and Hate Speech Detection. *ACM Trans. Asian Low-Resour. Lang. Inf. Process.* 1, 1, Article 1 (January 2020), 20 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Preprocessing is one of the critical steps in most text classification pipelines. The goal of preprocessing is to prepare textual input for accurate interpretation by a computational method. Earlier studies investigating the impact of preprocessing on text classification report different results for different preprocessing techniques and in different datasets [HaCohen-Kerner et al. 2020; Woo et al. 2020]. These findings indicate a need to better understand the impact of preprocessing on text classification. Moreover, the role of preprocessing in Arabic Natural Language Processing (NLP)

Authors' addresses: Fatemah Husain, f.husain@ku.edu.kw, Kuwait University, Kuwait City, Kuwait; Ozlem Uzuner, George Mason University, Fairfax, USA, ouzuner@gmu.edu.

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.

© 2020 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

53 has received limited attention [Duwairi and El-Orfali 2014; Saad 2010]. The Arabic language has unique characteristics
54 that need to be considered during preprocessing.

55 This study extends our previous work on offensive language detection and hate speech detection published within
56 the scope of the Fourth workshop on Open-Source Arabic Corpora and Corpora Processing Tools (OSACT) [Husain
57 2020b]. Our former study demonstrated the significant impact of intensive preprocessing on the performance of a
58 Support Vector Machine (SVM) classifier. Intensive preprocessing consisted of conversion of emojis to textual labels,
59 normalization of multiple shapes of Arabic letters to one shape for each letter, normalization of Arabic dialects to
60 Modern Standard Arabic (MSA) for a selected set of nouns, conversion of hyponyms to hypernyms for a selected set
61 of nouns, formatting hashtags to include textual content only, and basic cleaning such as removal of punctuation,
62 numbers, HTML tags, non-Arabic characters, etc. Our system ranked third in offensive language detection and ranked
63 first in hate speech detection, with macro-averaged F1 scores of 89% and 95%, respectively.

64 Similar intensive preprocessing has been adopted by other studies for offensive language detection but without fur-
65 ther investigating the contribution of different preprocessing techniques to the performance of classifiers. In [Husain
66 2020a], the author developed multiple ensemble traditional machine learning classifiers using all preprocessing tech-
67 niques defined by Husain [Husain 2020b] without reporting the change in performance before and after preprocessing.
68 Husain et al. [Husain et al. 2020] applied multiple preprocessing techniques with various deep learning classifiers with-
69 out showing the effect on performance.

70 This paper provides an in-depth investigation and evaluation of preprocessing on automatic offensive language
71 detection in Arabic. We utilize multiple classifiers and two different data sets. We assess the effect of various prepro-
72 cessing techniques individually and in combination. This study includes (1) traditional machine learning classifiers:
73 SVM and Logistic Regression (LR), (2) ensemble machine learning classifiers: Bagging and Random Forest (RF), (3)
74 Artificial Neural Network (ANN) classifiers: Recurrent Neural Networks (RNN) and Long Short Term Memory (LSTM),
75 and (4) Bidirectional Encoder Representations from Transformers (BERT) based models: AraBERT [Antoun et al. 2020]
76 and Arabic-BERT [Safaya et al. 2020].

77 This paper addresses the following questions:

- 78 (1) What are the most significant preprocessing techniques for Arabic offensive language detection?
- 79 (2) How does preprocessing affect traditional machine learning classifiers and the more advanced classifiers?
- 80 (3) Do the results of the advanced classifiers imply a need for change in traditional text classification methods?

81 The following section starts with an introduction to Arabic language, followed by a related work section. Then, we
82 discuss the methodology and the settings of the experiments, which include a description of the datasets, preprocessing
83 techniques, features, and classification models. The results of the experiments are presented after the methodology
84 section. We provide a detailed discussion and interpretation of our results in the discussion section. We include a
85 conclusion section to summarize our recommendations and findings. The paper ends with a presentation of future
86 work that highlights some possible ways to extend our work.

97 2 THE ARABIC LANGUAGE

98 Arabic is the official language for the majority of countries in the Middle East. It is also the original language of the
99 Quran and the Hadith, thus, it is widely studied by Muslims around the world. There are three main forms of the Arabic
100 language: (1) classical Arabic, which is the language used in Islamic manuscripts such as the Quran; (2) MSA, which
101 is the official language of schools, news outlets, Arabic books, etc.; and (3) dialectal Arabic, which is the actual spoken
102 language.

Arabic that varies from one area to another and from one society to another. According to [Habash 2010], the Arabic dialects can be divided into seven categories: Egyptian, Levantine, Gulf, North African, Iraqi, Yemeni, and Maltese.

Arabic is part of the Semitic language family. Scripts in Arabic are read and written in the opposite direction to English, i.e., from right to left. Arabic alphabet consists of 28 letters. It contains only two vowels: (أ/alif and ي/yaa'). Each letter in Arabic can be written in multiple shapes depending on the location within a word. For example, the letter "ف/faa" can be written as "ف / ف / ف" based on whether it is located at the beginning, in the middle, or at the end of a word [Gayar and Suen 2018; Habash 2010]. Figure 1 shows an example.

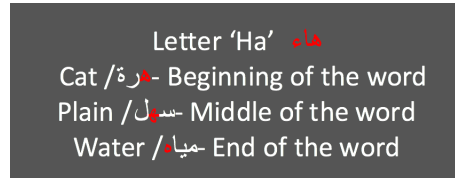


Fig. 1. Variations in the shape of Arabic letters based on the location of the letter within a word

The use of diacritics is very common in Arabic. Diacritics are often called Tashkil or Harakāt in Arabic, and they are used as a phonetic guide. In some cases, diacritics help reduce ambiguity as multiple words in Arabic can share the same spelling but have different pronunciation and meaning. For example, the Arabic word "شعر" means hair, "شعرٌ" means feel, and "شعرٌ" means poetry [Alhumoud et al. 2015; Boudad et al. 2018]. Figure 2 shows another example of the effect of diacritics.

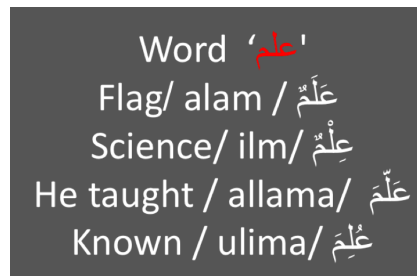


Fig. 2. The effect of diacritics on the meaning of Arabic words

Arabic words are justified by using elongation or Tatweel also called Kashida. Figure 3 shows an example of the use of elongation in Arabic.

The Arabic language has two genders for nouns: feminine and masculine. The base form is the masculine form. For example, the word "Qetta/قطة" refers to female cat and the word "Qett/قط" refers to male cat. In addition, Arabic has singular, dual, and plural forms of nouns, verbs, pronouns, and adjectives. For instance, for referring to two female cats, the word "Qettatan/قطان" is used, for two male cats, the word "Qettan/قطان" is used, and for plural cats, the word "Qettat/قطط" is used.

The aforementioned properties describe the general properties of MSA and dialectal Arabic. However, Arabic text in user-generated content is more diverse and has other properties that are important to consider as well. User-generated Arabic content is often transliterated in Latin characters, numbers, and punctuation. For example, the number "2" can

157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208

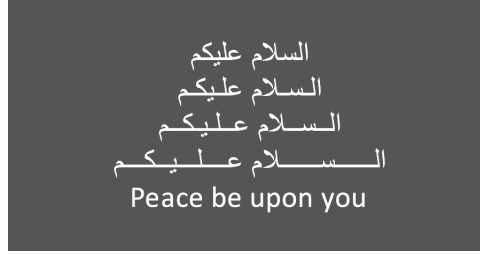


Fig. 3. Examples of elongation in Arabic

represent the letter “أ” (that sounds like “a” as in apple) and the number “3” can represent the letter “ع” (that is a guttural “aa”) [Darwish 2014]. Examples of words written in Arabizi, a form of user-generated Arabic, are “3a2albik/عالبك/on your heart” and “raw3a/روعة/wonderful” [Husain 2020b]. Moreover, code-switching among the various forms of Arabic and sometime English can be observed in the same post in user-generated content. Figure 4 demonstrates the diversity of user-generated content in Arabic language.

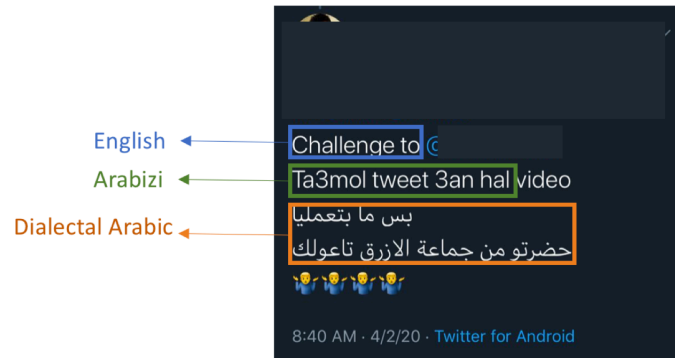


Fig. 4. An example user-generated post in Arabic

Some English terms are commonly used among Arabic speakers. Some Arabic text might show English text represented using Arabic letters. We show an example of a tweet for English text written in Arabic letters in Figure 5.

In some Arabic dialects, the pronunciation of the Arabic letters differs. Thus, users make use of letters from other languages to represent the dialectal sound of the letter rather than the standard Arabic letter. In this scenario, Urdu and Farsi alphabets are commonly used. Figure 6 shows an example of a tweet that is written in Iraqi Arabic, which replaces the letter “ك/kaa” in Arabic with the letter “گ/Gaa” in Farsi and Urdu.

3 RELATED WORK

There are few studies that focus on analyzing the effect of preprocessing on Arabic text classification. A pioneering study in this area was delivered by Saad in 2010 [Saad 2010]. Saad [Saad 2010] compared the effect of Arabic text preprocessing on the performance of different traditional machine learning classifiers, including: Decision Trees (DT), K-Nearest Neighbors (KNN), SVMs, Naïve Bayes (NB) and its variations (Multinomial Naïve Bayes (MNB), Complement (CNB), and Discriminative Multinomial Naïve Bayes (DMNB)). Preprocessing techniques that were applied

209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260



Fig. 5. An example for a post representing English in Arabic letters

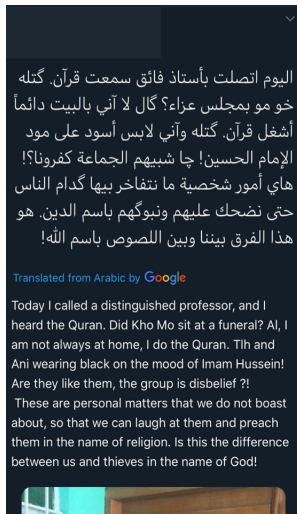


Fig. 6. An example of the use of non-Arabic letters in dialectal Arabic

include tokenizing strings to words, normalizing tokenized words, removing stop-words, different term weighting schemes, and Arabic morphological analysis. The weighting schemes consisted of the Boolean model that used 0 to refer to absence of a word or 1 to its presence, word count, normalized word count, term pruning, Term Frequency (TF), and Term Frequency-Inverse Document Frequency (TF-IDF). Arabic morphological analysis contained two stemming techniques: stemming and light stemming. Stemming substitutes words with their stems, while light stemming eliminates common affixes from words without substituting them with their stems. In [Saad 2010], seven Arabic datasets (BBC Arabic website bbcarabic.com corpus, CNN Arabic website cnnarabic.com corpus, Open Source Arabic Corpus (OSAC), Contemporary Corpus of Arabic (CCA), Aljazeera News corpus, Khaleej newspaper of the year 2004 corpus, and a corpus from multiple online Arabic newspapers) were used to evaluate the effect of preprocessing using Bag-of-Words (BOW). Results highlighted the significance of preprocessing Arabic text for feature reduction, which reduces

261 complexity of classifiers, consumes less storage, and saves processing time. Findings also showed the following: (1)
262 light stemming with term pruning achieves the highest performance score; (2) SVM and variations of NB perform
263 better than other algorithms; and (3) weighting schemes influence the performance of distance-based classifiers.
264

265 In another study [Duwairi and El-Orfali 2014], as part of building a sentiment analysis system for Arabic text, the
266 authors analyzed the effect of removing stop-words and stemming on the performance of classifiers. They used two
267 datasets: the first one contains reviews that are collected from political articles at Aljazeera Arabic news website and
268 the second one contains movie reviews written in MSA. They used BOW to represent their features and applied SVM,
269 NB, and KNN classifiers. The results showed that preprocessing improves performance over no preprocessing on the
270 Aljazeera dataset. However, for the movie reviews, stemming hurt performance. Authors attributed this result to the
271 use of the Rapidminer’s stemmer, which has a high error rate. The authors also pointed out the fact that the movie
272 reviews dataset is larger than the Aljazeera dataset, thus, the high error rate of the stemmer is more obvious and shows
273 sharper effects on the movie reviews dataset.
274
275

276 Effect of preprocessing on text classification has been investigated more extensively in English than in Arabic. In
277 this section, we also discuss two recent studies that focus on the effect of preprocessing on English text classification.
278 HaCohen-Kerner, Miller, and Yigal [HaCohen-Kerner et al. 2020] studied six preprocessing techniques: spelling cor-
279 rection, HTML tag removal, converting uppercase letters into lowercase letters, punctuation mark removal, reduction
280 of repeated characters, and stop-word removal. The experiments were performed on four datasets: the 4 Universities
281 Data Set (WebKB), the multi-labeled Reuters-21578 dataset (R8), the SMS Spam Collection dataset, and the Sentiment
282 Labelled Sentences dataset. They utilized BOW and three classifiers: a Bayes Networks (BN), a variant of SVM (SMO),
283 and a Random Forest (RF). The results demonstrates the importance of experimenting with various preprocessing
284 techniques in systematic combinations when developing a text classification system as different datasets show differ-
285 ent effects of preprocessing. For example, stop-word removal shows significant improvement in accuracy score when
286 applied to WebKB, R8, and SMS Spam Collection, but not when applied to Sentiment Labelled Sentences. Thus, the
287 main finding is that there are no preprocessing techniques which can always improve text classification performance
288 regardless of the dataset.
289
290

291
292 Woo, Kim, and Lee [Woo et al. 2020] also analyzed the effect of preprocessing on text classification. Their study
293 applied two main types of preprocessing techniques: special character elimination, lemmatization, lowering upper
294 case, and punctuation splitting or merging were categorized as typical preprocessing techniques. The second type of
295 preprocessing techniques were more advanced and consisted of technical terminology preprocessing and sorting sen-
296 tences based on their complexity. Technical terminology preprocessing used a customized rule-based algorithm with
297 a technical term corpus. Rules were used to extract technical terminology and Wikipedia API was used to verify the
298 accuracy of the identified terminologies. When a technical terminology was identified, the terminology was reformed
299 to ensure proper segmentation during preprocessing to reduce any semantic change. The second advanced preprocess-
300 ing technique calculated the entropy of each sentence based on the distribution of syllables. Entropy was used as the
301 sentence complexity score. The Stanford Natural Language Inference (SNLI) corpus was used for training and testing
302 the model. In this study, the authors considered modern neural network models only: a Convolutional Neural Network
303 (CNN) model, Siamese LSTM model, and a transformer model. The experiments evaluated combinations of preprocess-
304 ing techniques. Results showed higher accuracy scores when two preprocessing techniques are applied, specifically,
305 lemmatization and punctuation splitting, lemmatization and lowering upper case letters, or lowering upper case let-
306 ters and punctuation splitting. If only one technique is applied, the authors recommend using lemmatization, lowering
307 upper case letters, or punctuation splitting. The results also demonstrate that preprocessing techniques which shorten
308
309
310
311
312

the lengths of sentences are not recommended. Special character elimination, normalization techniques, and ordering sentences based on their complexities do not improve accuracy.

This paper covers the gap in previous studies by providing an in-depth analysis of various Arabic text preprocessing techniques using several classifiers, including state-of-the-art language models, for text classification.

4 METHODOLOGY

Generally, we follow a traditional text classification pipeline in our experimental study, consisting of preprocessing, feature extraction, classification, and performance evaluation. We borrow the preprocessing procedures from Husain [Husain 2020b]. We apply each preprocessing technique separately and all preprocessing techniques together to create multiple versions of the input datasets in addition to the raw data. We apply multiple classifiers to each version of the datasets. We compare the results of classifiers that utilize input text preprocessed using different techniques and their combination against the raw data baseline in order to understand the contribution of each preprocessing technique to the end result.

4.1 Datasets

The datasets used in this study contain tweets. One dataset contains multiple dialects while the other is focused on the Levantine dialect.

The first dataset is a multi-dialect OSACT dataset [Mubarak et al. 2020], which supports two offensive language classification tasks: (a) detecting if a post is offensive or not, and (b) identifying if a post is hate speech or not. The providers of the OSACT dataset performed some preprocessing to ensure the privacy of users: Twitter user mentions were substituted by “@USER”, URLs were substituted by “URL”, and empty lines were replaced by “<LF>”. Moreover, the OSACT dataset is highly imbalanced, 19% of the entire dataset is offensive tweets, and 5% contains hate speech.

The second offensive language dataset is a Levantine Twitter Dataset for Hate Speech and Abusive Language (L-HSAB)¹, which contains 5,846 tweets labeled as hate speech (468 tweets), abusive language (1,728 tweets), or normal (3,650 tweets) [Mulki et al. 2019]. L-HSAB is the first Levantine hate speech Arabic dataset. The dataset is manually filtered to ensure the collection of only Syrian and Lebanese tweets without duplication. Then, the tweets are cleaned to remove some characters (such as @, RT, and #) before presenting them to the annotators.

For both datasets, we shuffle them and then randomly split them into 80% training dataset and 20% testing dataset. Table 1 summarizes classification tasks related to each dataset. For all models, the training portion is used to develop the models and the testing portion is used to evaluate the models.

4.2 Preprocessing techniques

The following are the preprocessing techniques that were included in our experiments.

4.2.1 Emoji Conversion. Emojis are often used to convey feelings and attitudes, which could be very valuable to offensive language detection. Preprocessing text for emojis has been shown to improve aggression detection [Orăsan 2018]. We convert emojis to Arabic textual labels that describe the content of the emojis. For this, we use BeautifulSoup4 4.8.2² to extract the entire set of emojis defined by Unicode.org³, which provides textual descriptions for emojis written in

¹<https://github.com/Hala-Mulki/L-HSAB-First-Arabic-Levantine-HateSpeech-Dataset>

²<https://pypi.org/project/beautifulsoup4/>

³<https://home.unicode.org>

Table 1. Classification tasks

| Dataset | Number of Classes | Class Labels | Training Dataset Size | Testing Dataset Size |
|---------|-------------------|-----------------|-----------------------|-----------------------|
| OSACT | 2 | Offensive | 1,581 offensive | 408 offensive |
| | | Not offensive | 6,417 not offensive | 1,590 not offensive |
| OSACT | 2 | Hate speech | 408 hate speech | 96 hate speech |
| | | Not hate speech | 7,590 not hate speech | 1,902 not hate speech |
| L-HSAB | 3 | Abusive | 1,399 abusive | 327 abusive |
| | | Hate | 365 hate | 101 hate |
| | | Normal | 2,909 normal | 739 normal |

Table 2. Examples of emojis and their Arabic labels with English translations

| Emoji | Arabic Label | English Translation |
|-------|-----------------------|-----------------------------|
| 😊 | وجه مبتسم قليلا | slightly smiling face |
| 😄 | وجه مبتسم بعيون كبيرة | grinning face with big eyes |
| 🐵 | وجه القرد | monkey face |
| 🐒 | قرد | monkey |
| 🎂 | كعكة عيد الميلاد | birthday cake |

English language. We translate the descriptions to Arabic using Translate 1.0.7 python package⁴. The final extracted emoji list contains 1,374 emojis. Table 2 shows examples of emojis with their Arabic textual labels along with their English translations.

4.2.2 Arabic Dialect Normalization. To normalize dialectal nouns, we convert to MSA. For example, the word cat from multiple dialects – Egyptian “أطة/Otta”, Levantine “بسة/Bisse”, Gulf “قطوة/Qatwa”, Moroccan “قطلة/Qetta”, Yemeni “دماه/demah”, and Iraqi “بازونة/Bazzuna” – are converted to MSA form “قطلة”. The selected dialectal nouns are based on a manual inspection of the training set.

4.2.3 Word Categorization. This step converts hyponyms to hypernyms for a selected set of nouns. From manual inspection of a sample of tweets from the training dataset, we notice that it is very common to mention names of animals among hate speech tweets. Thus, we manually created a list of common animal names used in different Arabic dialects, such as “كلب/dog”, “خنزير/pig”, “حية/snake”. The list considers the dialectal variations in animal names as well as the gender and quantity variations in animal names. We substitute all animal names with the Arabic word “حيوان”, which means animal.

4.2.4 Letter Normalization. Some Arabic letters can be written in various forms depending on their location within a word. We normalize three letters to convert all their forms into one form – Alif (أ، آ، إ to ا), Alif Maqsura (ئ، ي to ي), and Ta Marbuta (ة to ه). In addition, letters repeated more than two times within a word are reduced to two times only.

4.2.5 Miscellaneous & Hashtag Segmentation. In addition to the previous techniques, we apply some miscellaneous preprocessing techniques that are commonly performed in the field. We remove numbers, Kashida, diacritics, HTML tags, more than one space, three or more repetitions of any character, and some symbols or terms (e.g., “”, “...”, “!”, “?”).

⁴<https://pypi.org/project/translate/1.0.7/>

469 4.4.2 *Ensemble Machine Learning*. Bagging and Random Forest represent ensemble machine learning classifiers. They
470 were implemented using Python scikit-learn library. Similar to the traditional machine learning classifiers, BOW with
471 2 to 5 characters is used in both.
472

473 4.4.3 *Artificial Neural Network (ANN)*. The ANN classifiers – RNN and LSTM – are also included in the study. We use
474 Keras to develop ANN models using 3000 for maximum features, 0.25 for drop-out, 100 units, 1000 for patch size, and
475 500 epochs. Only the ANN classifiers are implemented using two features: BOW with 2 to 5 characters and AraVec
476 word embeddings. Each feature is used separately from the other.
477

479 4.4.4 *BERT-Based Models*. We include BERT-based language models in our analysis. BERT is the current state-of-the-
480 art in text classification. We use AraBERT [Antoun et al. 2020] (first version of AraBERT, also called AraBERTv1-base
481 and bert-base-arabert) and Arabic-BERT [Safaya et al. 2020] to represent BERT-based classifiers. We apply each using
482 maximum length = 128, patch size = 16, epochs = 2, epsilon = 1e-8, and learning rate = 2e-5. We do not use feature
483 engineering; instead, we use the pooling layer from the encoder and feed it into a Feed Forward Neural Network
484 (FFNN). We use the HuggingFace library⁶ to apply both BERT models and develop our experiments using the PyTorch-
485 Transformers library.
486
487

488 5 RESULTS

489
490 We consider macro-averaged measurements for all metrics: precision, recall, accuracy, and F1. The choice of macro-
491 averaged metrics is motivated by the class imbalance in our datasets. Overall, results do not demonstrate any significant
492 patterns for specific Arabic text preprocessing techniques across all tasks. The variations in performance achieved by
493 different preprocessing techniques for each classifier are small.
494

495 Figure 7 and Figure 8 show results for offensive language detection using the OSACT dataset. As can be seen from
496 Figures (a) to (d) for the traditional machine learning and ensemble classifiers, the SVM and the Random Forest perform
497 better when applied to fully processed text, while the LR and Bagging perform better when applied to dialect normalized
498 text. The ANN classifiers show different results when applied using BOW features versus AraVec features, as can be
499 noticed from Figures (e) and (f). In general, preprocessing text reduces the performance of BOW-based ANN classifiers.
500 For Figures (g) and (h), miscellaneous preprocessing and hashtag segmentation improve the performance of AraVec-
501 based ANN classifiers. Both BERT-based models in Figures (i) and (j) show better performance when applied to dialect
502 normalized text, while AraBERT also reports better performance when applied to emoji converted text.
503

504 Results for hate speech detection using the OSACT dataset are presented in Figure 9 and Figure 10. The SVM (Figure
505 (a)) shows very high performance when applied to fully processed text; however, the LR reports the lowest performance
506 when applied to fully processed text. As can be seen from Figure (b), the LR performs better with miscellaneous pre-
507 processed text and hashtag segmentation. Figure (c) for the Random Forest demonstrates equal effects on performance
508 by raw text and fully processed text, which also report the highest performance. The Bagging classifier (Figure (d)) and
509 all ANN classifiers from Figures (e) to (h) improve in performance when applied with miscellaneous preprocessed text
510 and with hashtag segmentation. BERT-based models (Figures (i) and (j)) do not show improvement on preprocessed
511 text over raw text.
512

513 Figure 11 and Figure 12 highlight the results of abusive and hate speech detection using the L-HSAB dataset. Figure
514 (a) for the SVM demonstrates that there is no improvement from preprocessing text, while Figure (b) for the LR reports
515

516
517
518
519 ⁶<https://huggingface.co/>

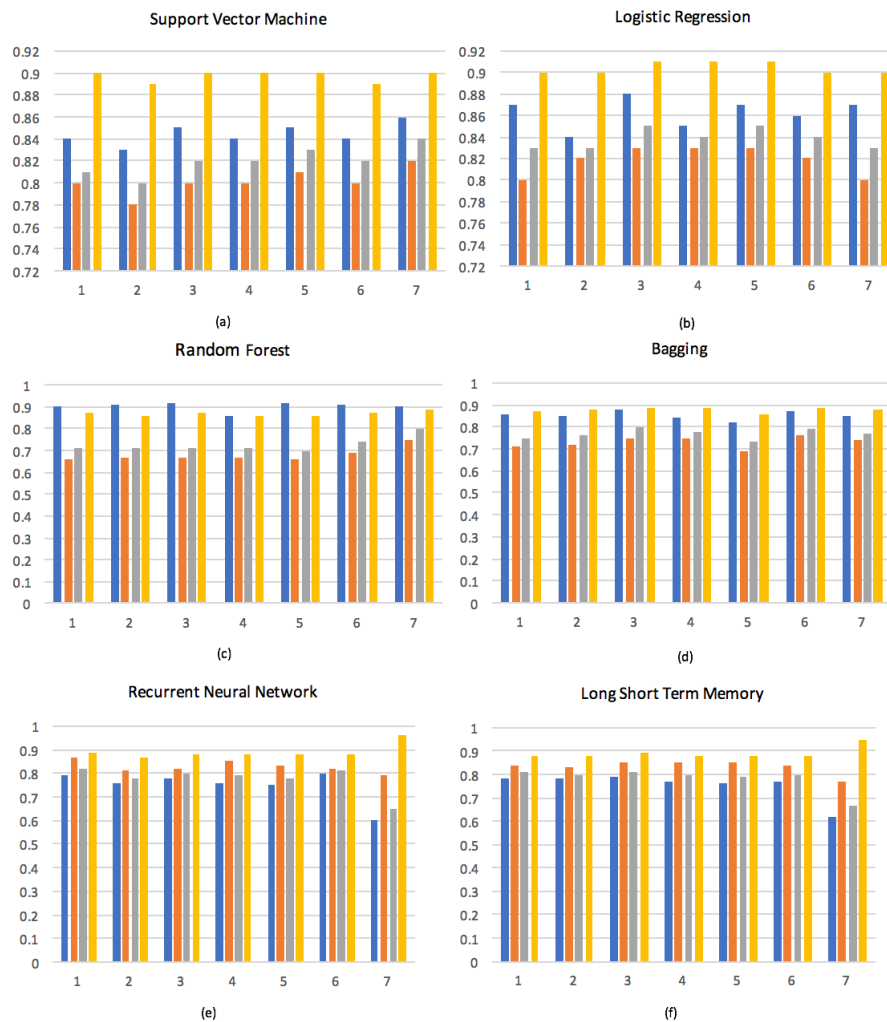


Fig. 7. Results of offensive language detection using the OSACT dataset

slightly better performance on dialect normalized text and letter normalized text over raw text. The Random Forest (Figure (c)) achieves better performance on dialect normalized text and fully processed text. On the other hand, the Bagging model in Figure (d) does not gain from preprocessing text. The RNN model with BOW (Figure (e)) does not gain from preprocessing as raw text records the highest performance. Figure (f) for the LSTM model with BOW shows the best performance on raw text, dialect normalized text, and miscellaneous preprocessed and hashtag segmented text. AraVec-based models show different results, in Figure (g) for the RNN, best performance is on miscellaneous preprocessed and hashtag segmented text, while in Figure (h) for the LSTM, the best performance is on raw text. Figure (i) for AraBERT shows the same highest performance score on four preprocessing techniques: emoji conversion, dialect normalization, word categorization, and miscellaneous preprocessing and hashtag segmentation. Arabic-BERT model

demonstrates slightly better performance on text preprocessed for emoji conversion and dialect normalization over raw text.

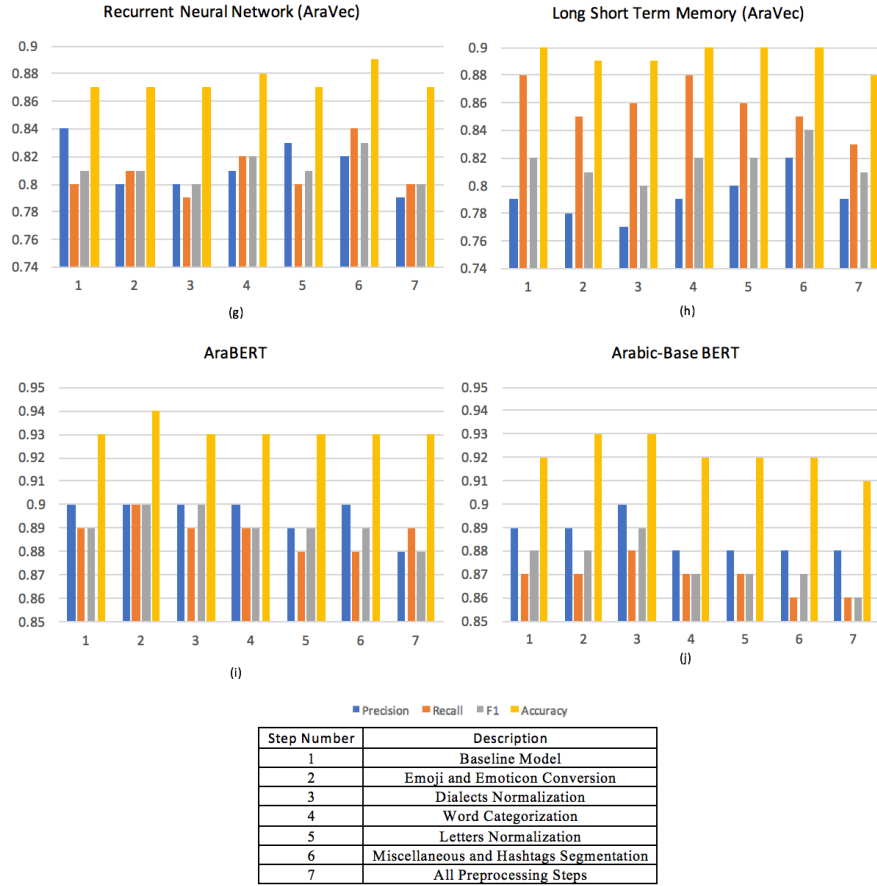


Fig. 8. Results for offensive language detection using the OSACT dataset (continue)

6 DISCUSSION AND FINDINGS

We perform a statistical analysis for the macro-averaged F-1 results to summarize results from the three offensive language detection tasks: offensive language detection on the OSACT dataset, hate speech detection on the OSACT dataset, and abusive and hate speech detection on the L-HSAB dataset. The statistical analysis includes minimum, maximum, mean, standard deviation, and range of macro-averaged F1 scores among the preprocessing techniques presented in the previous section, in addition to the raw text without preprocessing (baseline) and fully processed text using all preprocessing techniques on the same time. Tables 4, 5, and 6 show the statistical analysis for the results from the three tasks.

Generally, the results demonstrate the limited effect of Arabic text preprocessing on offensive language detection. From the statistical analysis, the traditional machine learning classifiers gain in performance more than the ensemble

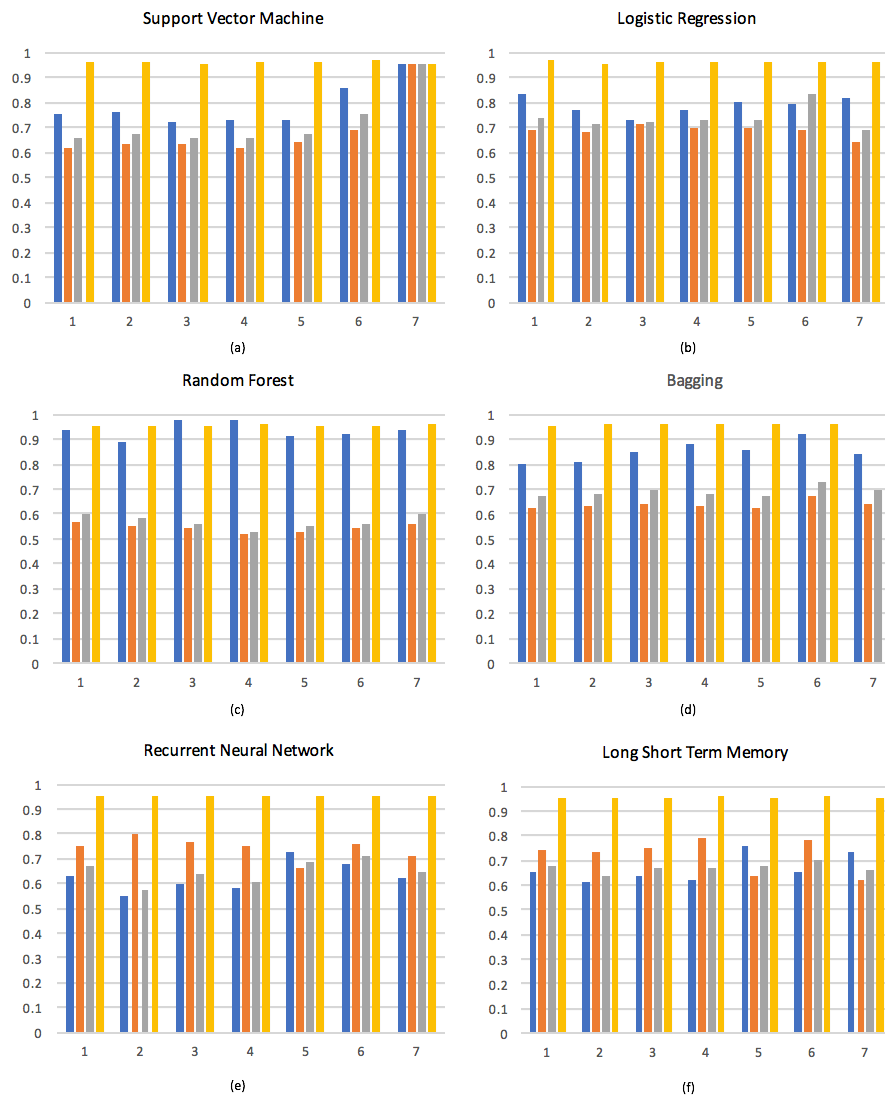


Fig. 9. Results for hate speech detection using the OSACT dataset

machine learning classifiers and deep learning classifiers. Only for hate speech detection on the OSACT dataset, the traditional machine learning, particularly the SVM, performs much better than all other classifiers including BERT-based classifiers.

The advanced BERT-based models gain the least from preprocessing. This finding implies that for systems that depend on advanced models, we can use the raw text without preprocessing. We expect that this is because BERT-based models learn the context rather than the individual words and their exact form, and are therefore able to effectively conflate the variations of words that appear in the same context. Preprocessing those words would be redundant in terms of conflating variants as the contextual representation which is already immune to the variations.

677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728

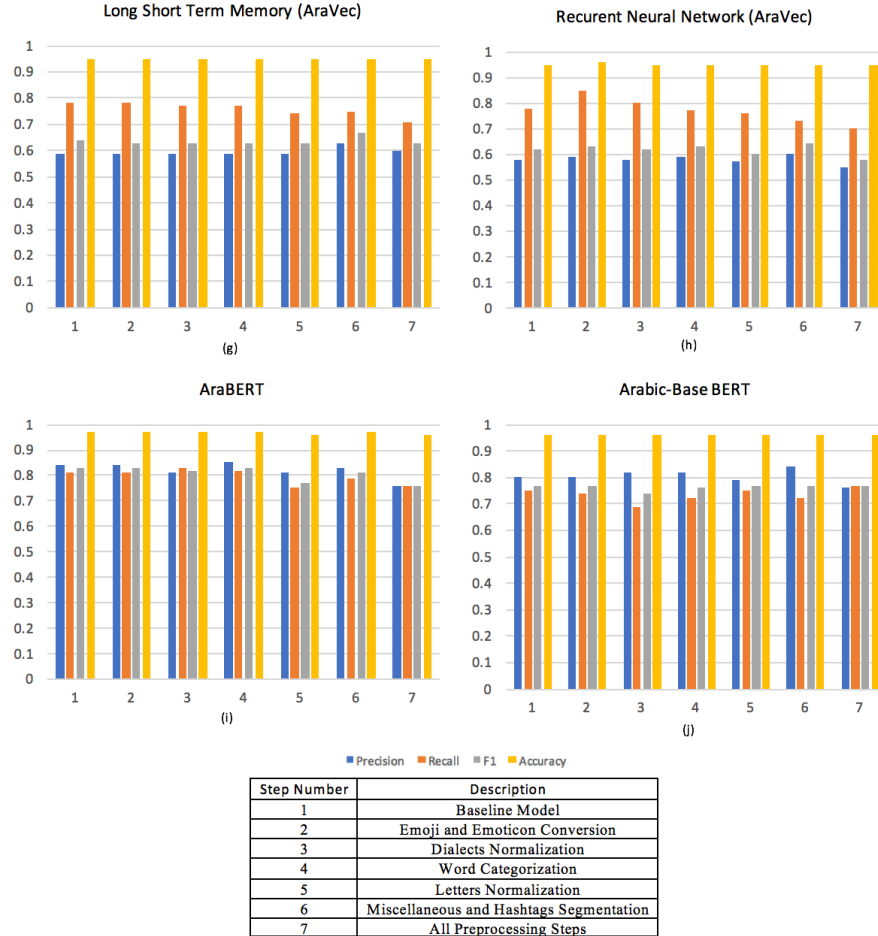


Fig. 10. Results for hate speech detection using the OSACT dataset (continue)

Examining the values of standard deviations from the statistical analysis tables show that the majority of experiments show very small values of 0.01, which indicates that macro-averaged F1 scores of the preprocessing techniques and the other two versions of the data – raw and fully processed using all techniques together – are clustered very closely around the mean. Among all tables, the largest standard deviation is 0.11, which is the result from the SVM classifier for hate speech detection. This finding reveals the high sensitivity of the SVM classifier to the various forms of the input text.

The statistical analysis tables also show that the lowest macro-averaged F1 scores for all tasks come from ensemble machine learning classifiers: Random Forest and Bagging. However, the classifier with the best performance varies among the tasks. For instance, AraBERT performs best on both offensive language detection on the OSACT dataset and on abusive and hate speech detection on the L-HSAB dataset, while SVM performs best on hate speech detection on the OSACT dataset.

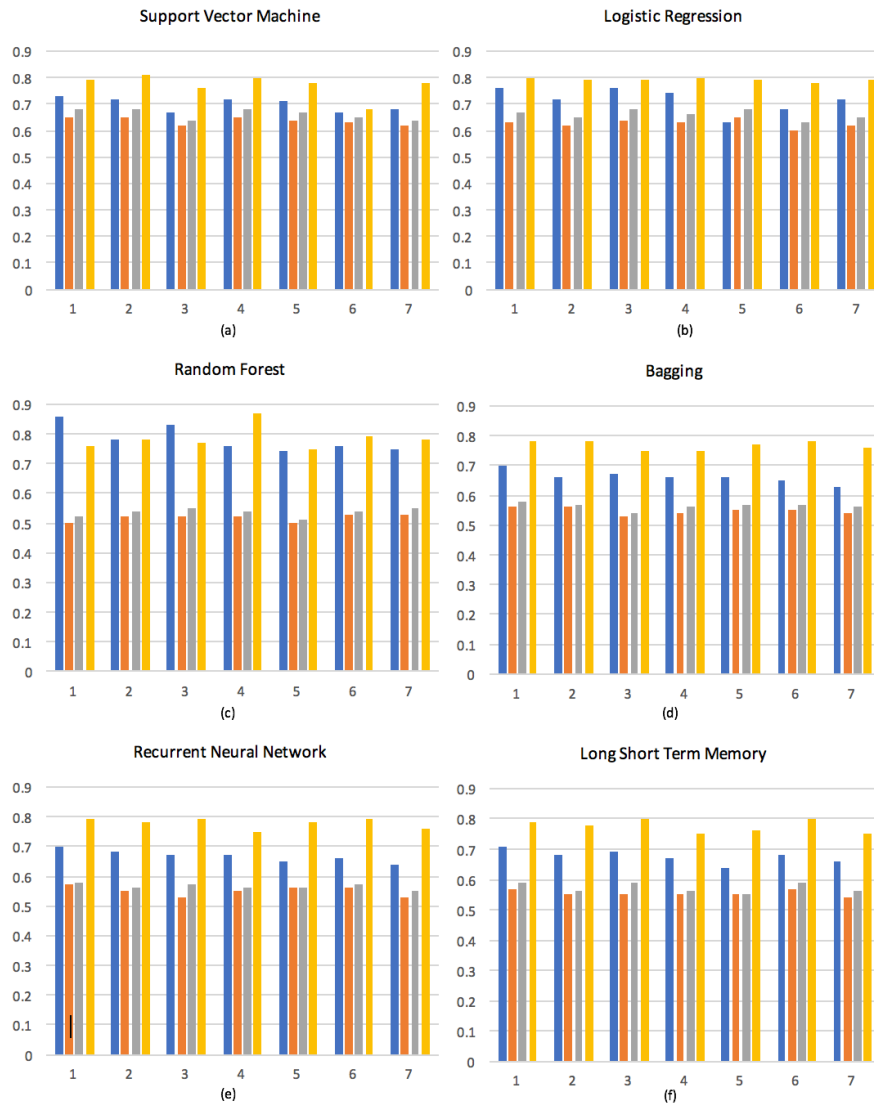
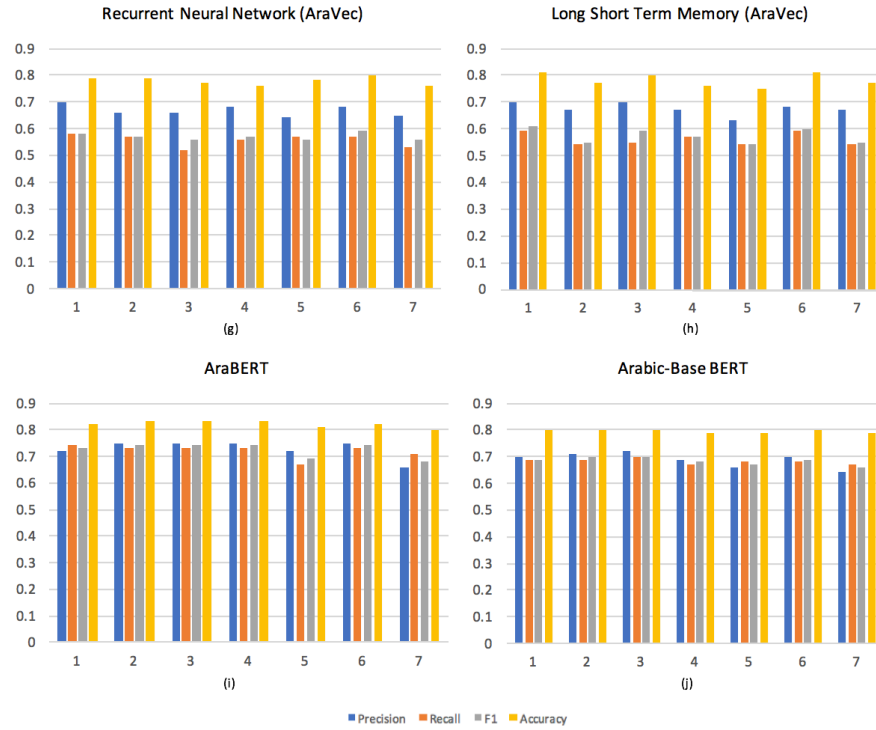


Fig. 11. Results for abusive and hate speech detection using the L-HSAB dataset

For offensive language detection on the OSACT dataset, the highest macro-averaged F1 score is 0.90, achieved by the AraBERT classifier. Table 4 shows that the highest variations among all forms of input text occur when applying the RNN and the LSTM with BOW features for offensive language detection. Table 5 presents the statistical analysis for hate speech detection on the OSACT dataset. It highlights the effectiveness of the SVM classifier. Moreover, the variation is relatively high among the input text forms for the SVM classifier in comparison to the rest of the classifiers. In Table 6 for abusive and hate speech detection on the L-HSAB dataset, the variations become more even, most of the classifiers have 0.01 standard deviation and BERT-based classifiers significantly exceed the others in performance.



| Step Number | Description |
|-------------|---|
| 1 | Baseline Model |
| 2 | Emoji and Emoticon Conversion |
| 3 | Dialects Normalization |
| 4 | Word Categorization |
| 5 | Letters Normalization |
| 6 | Miscellaneous and Hashtags Segmentation |
| 7 | All Preprocessing Steps |

Fig. 12. Results for abusive and hate speech detection using the L-HSAB dataset (continue)

This finding could be related to the diversity of the Arabic language forms used in the OSACT dataset versus the L-HSAB dataset, as the OSACT dataset covers broader dialects. Consequently, the preprocessing techniques we developed significantly affect performance. The L-HSAB dataset has only Levantine text; thus, preprocessing has less impact on the performance.

To further examine the generalizability and robustness of the classifiers on preprocessed text, we aggregate macro-averaged F1 scores from all tasks into Table 7. AraBERT model shows the highest minimum macro-averaged F1, the highest maximum macro-averaged F1, and the highest mean macro-averaged F1 scores. These results indicate the power of AraBERT in detecting offensive content regardless of preprocessing.

Table 8 summarizes the results from all experiments, including all tasks depending on macro-averaged F1 score per preprocessing technique to analyze the results for each technique separately. Results show the highest minimum and maximum macro-averaged F1 scores when using all preprocessing techniques at the same time. Text that has been preprocessed using miscellaneous and hashtag segmentation shows the highest average macro-averaged F1. All techniques report very similar standard deviation scores, which means that the variation in macro-averaged F1 among the

Table 4. Statistical analysis of macro-averaged F1 results for offensive language detection using the OSACT dataset

| Classifier | Minimum | Maximum | Mean | Standard Deviation | Range |
|---------------|---------|---------|------|--------------------|-------|
| SVM | 0.80 | 0.84 | 0.82 | 0.01 | 0.04 |
| LR | 0.83 | 0.85 | 0.84 | 0.00 | 0.02 |
| Random Forest | 0.70 | 0.80 | 0.72 | 0.03 | 0.10 |
| Bagging | 0.73 | 0.80 | 0.77 | 0.02 | 0.07 |
| RNN (BOW) | 0.65 | 0.82 | 0.77 | 0.06 | 0.17 |
| LSTM (BOW) | 0.67 | 0.81 | 0.78 | 0.05 | 0.14 |
| RNN (AraVec) | 0.80 | 0.83 | 0.81 | 0.01 | 0.03 |
| LSTM (AraVec) | 0.80 | 0.84 | 0.82 | 0.01 | 0.04 |
| AraBERT | 0.88 | 0.90 | 0.89 | 0.00 | 0.02 |
| Arabic-BERT | 0.86 | 0.89 | 0.87 | 0.00 | 0.03 |

Table 5. Statistical analysis of macro-averaged F1 results for hate speech detection using the OSACT dataset

| Classifier | Minimum | Maximum | Mean | Standard Deviation | Range |
|---------------|---------|---------|------|--------------------|-------|
| SVM | 0.66 | 0.95 | 0.72 | 0.11 | 0.29 |
| LR | 0.69 | 0.83 | 0.73 | 0.04 | 0.14 |
| Random Forest | 0.53 | 0.60 | 0.57 | 0.03 | 0.07 |
| Bagging | 0.67 | 0.73 | 0.69 | 0.02 | 0.06 |
| RNN (BOW) | 0.57 | 0.71 | 0.65 | 0.05 | 0.14 |
| LSTM (BOW) | 0.64 | 0.70 | 0.67 | 0.02 | 0.06 |
| RNN (AraVec) | 0.58 | 0.64 | 0.62 | 0.02 | 0.06 |
| LSTM (AraVec) | 0.63 | 0.67 | 0.64 | 0.01 | 0.04 |
| AraBERT | 0.76 | 0.83 | 0.80 | 0.03 | 0.07 |
| Arabic-BERT | 0.74 | 0.77 | 0.76 | 0.01 | 0.03 |

tasks is very small regardless of the preprocessing technique. The lowest range is 0.35, which is achieved by several preprocessing techniques, including emoji conversion, dialect normalization, and miscellaneous and hashtag segmentation. Overall, the findings from this table demonstrate a better performance (highest maximum macro-averaged F1) among all tasks when applying all preprocessing techniques at the same time.

7 CONCLUSION

In this study, we examined the impact of several preprocessing techniques on Arabic offensive language detection. We focus on five preprocessing techniques: conversion of emojis to Arabic textual labels, normalization of different forms of Arabic letters, normalization of selected nouns from Arabic dialects to MSA, conversion of selected hyponyms to hypernyms, hashtag segmentation, and basic cleaning processes, such as removing numbers, Kashida, diacritics, and

Table 6. Statistical analysis of macro-averaged F1 results for abusive and hate speech detection using the L-HSAB dataset

| Classifier | Minimum | Maximum | Mean | Standard Deviation | Range |
|---------------|---------|---------|------|--------------------|-------|
| SVM | 0.64 | 0.68 | 0.66 | 0.02 | 0.04 |
| LR | 0.63 | 0.68 | 0.66 | 0.02 | 0.05 |
| Random Forest | 0.51 | 0.55 | 0.53 | 0.01 | 0.04 |
| Bagging | 0.54 | 0.58 | 0.56 | 0.01 | 0.04 |
| RNN (BOW) | 0.55 | 0.58 | 0.56 | 0.01 | 0.03 |
| LSTM (BOW) | 0.55 | 0.59 | 0.57 | 0.02 | 0.04 |
| RNN (AraVec) | 0.56 | 0.59 | 0.57 | 0.01 | 0.03 |
| LSTM (AraVec) | 0.54 | 0.61 | 0.57 | 0.03 | 0.07 |
| AraBERT | 0.68 | 0.74 | 0.72 | 0.03 | 0.06 |
| Arabic-BERT | 0.66 | 0.70 | 0.68 | 0.01 | 0.04 |

Table 7. Statistical analysis of macro-averaged F1 results from all tasks

| Classifier | Minimum | Maximum | Mean | Standard Deviation | Range |
|---------------|---------|---------|------|--------------------|-------|
| SVM | 0.70 | 0.82 | 0.73 | 0.05 | 0.12 |
| LR | 0.72 | 0.79 | 0.74 | 0.02 | 0.07 |
| Random Forest | 0.58 | 0.65 | 0.65 | 0.02 | 0.07 |
| Bagging | 0.65 | 0.70 | 0.70 | 0.02 | 0.06 |
| RNN (BOW) | 0.59 | 0.70 | 0.66 | 0.04 | 0.11 |
| LSTM (BOW) | 0.62 | 0.70 | 0.67 | 0.02 | 0.08 |
| RNN (AraVec) | 0.65 | 0.69 | 0.67 | 0.01 | 0.04 |
| LSTM (AraVec) | 0.66 | 0.71 | 0.68 | 0.02 | 0.05 |
| AraBERT | 0.77 | 0.82 | 0.80 | 0.02 | 0.05 |
| Arabic-BERT | 0.75 | 0.79 | 0.77 | 0.01 | 0.03 |

HTML tags. The experiments utilize several classifiers, including traditional machine learning classifiers, ensemble machine learning classifiers, deep learning classifiers with BOW and AraVec, and BERT-based classifiers. The results show unequal effects of preprocessing techniques based on the classifier and demonstrate limited gains from preprocessing on offensive language detection. Specifically for the advanced BERT-based classifiers, preprocessing shows limited gains and can be omitted from text classification pipelines.

The findings from this study indicate that certain possible combinations of preprocessing techniques might improve the performance for some classifiers, such as deep learning systems with BOW features or SVMs for hate speech. This finding demonstrates the importance of further investigating the impact of preprocessing Arabic text for offensive language detection by conducting more comprehensive studies that include multiple dialects and social media platforms.

Table 8. Statistical analysis of macro-averaged F1 results from each preprocessing technique among all tasks

| Preprocessing technique | Minimum | Maximum | Mean | Standard Deviation | Range |
|--|---------|---------|------|--------------------|-------|
| No preprocessing | 0.52 | 0.89 | 0.71 | 0.10 | 0.41 |
| Emoji Conversion conversion | 0.54 | 0.90 | 0.70 | 0.11 | 0.35 |
| Dialect normalization | 0.54 | 0.90 | 0.70 | 0.11 | 0.35 |
| Word categorization | 0.53 | 0.89 | 0.70 | 0.11 | 0.36 |
| Letter normalization | 0.51 | 0.89 | 0.70 | 0.11 | 0.38 |
| Miscellaneous and hashtag segmentation | 0.54 | 0.89 | 0.72 | 0.10 | 0.35 |
| All preprocessing | 0.55 | 0.95 | 0.70 | 0.11 | 0.40 |

Our experiments utilize relatively small datasets with small number of offensive samples; thus, studying the effects with larger and more balanced datasets would support the generalizability of our findings.

8 FUTURE WORK

While preprocessing and its effects are important to understand, future work would benefit from inclusion of semantic information in addressing Arabic natural language processing tasks.

REFERENCES

- Sarah Alhumoud, Mawaheb Altuwaijri, Tarfa Albuhaireh, and Wejdan Alohaideb. 2015. Survey on Arabic Sentiment Analysis in Twitter, In World Academy of Science, Engineering and Technology. *International Journal of Social, Behavioral, Educational, Economic, Business and Industrial Engineering* 9, 1, 364–378.
- Wissam Antoun, Fady Baly, and Hazem Hajj. 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*. European Language Resource Association, Marseille, France, 9–15. <https://www.aclweb.org/anthology/2020.osact-1.2>
- Naaima Boudad, Rdouan Faizi, Rachid Oulad Haj Thami, and Raddouane Chiheb. 2018. Sentiment analysis in Arabic: A review of the literature. *Ain Shams Engineering Journal* 9, 4 (2018), 2479 – 2490. <https://doi.org/10.1016/j.asej.2017.04.007>
- Kareem Darwish. 2014. Arabizi Detection and Conversion to Arabic. In *Proceedings of the EMNLP 2014 Workshop on Arabic Natural Language Processing (ANLP)*. Association for Computational Linguistics, Doha, Qatar, 217–224. <https://doi.org/10.3115/v1/W14-3629>
- Rehab Duwairi and Mahmoud El-Orfali. 2014. A study of the effects of preprocessing strategies on sentiment analysis for Arabic text. *Journal of Information Science* 40, 4 (2014), 501–513. <https://doi.org/10.1177/0165551514534143> arXiv:<https://doi.org/10.1177/0165551514534143>
- Neamat El Gayar and Ching Suen. 2018. . Series on Language Processing, Pattern Recognition, and Intelligent Systems, Vol. 4. World Scientific. 1– 288 pages. <https://doi.org/10.1142/10693>
- Nizar Y. Habash. 2010. (1 ed.), Synthesis Lectures on Human Language Technologies, Vol. 3. Morgan Claypool Publishers. 1–187 pages. <https://doi.org/10.2200/S00277ED1V01Y201008HLT010>
- Yaakov HaCohen-Kerner, D. Miller, and Yair Yigal. 2020. The influence of preprocessing on text classification using a bag-of-words representation. *PLoS ONE* 15 (2020).
- Fatemah Husain. 2020a. Arabic Offensive Language Detection Using Machine Learning and Ensemble Machine Learning Approaches. arXiv:2005.08946 [cs.CL]
- Fatemah Husain. 2020b. OSACT4 Shared Task on Offensive Language Detection: Intensive Preprocessing-Based Approach. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools, with a Shared Task on Offensive Language Detection*. European Language Resource Association, Marseille, France, 53–60. <https://www.aclweb.org/anthology/2020.osact-1.8>
- Fatemah Husain, Jooyeon Lee, Samuel Henry, and Ozlem Uzuner. 2020. SalamNET at SemEval-2020 Task12: Deep Learning Approach for Arabic Offensive Language Detection. In *Proceedings of the 14th International Workshop on Semantic Evaluation*. Barcelona, Spain, 2133–2139.

- 989 Hamdy Mubarak, Kareem Darwish, Walid Magdy, Tamer Elsayed, and Hend Al-Khalifa. 2020. Overview of OSACT4 Arabic Offensive Language Detection
990 Shared Task. In *Proceedings of the 4th Workshop on Open-Source Arabic Corpora and Processing Tools (OSACT), with a Shared Task on Offensive Language*
991 *Detection*, Vol. 4. European Language Resource Association, Marseille, France. <https://www.aclweb.org/anthology/2020.osact-1.2>
- 992 Hala Mulki, Hatem Haddad, Chedi Bechikh Ali, and Halima Alshabani. 2019. L-HSAB: A Levantine Twitter Dataset for Hate Speech and Abusive
993 Language. In *Proceedings of the Third Workshop on Abusive Language Online*. Association for Computational Linguistics, Florence, Italy, 111–118.
994 <https://doi.org/10.18653/v1/W19-3512>
- 995 Constantin Orăsan. 2018. Aggressive Language Identification Using Word Embeddings and Sentiment Features. In *Proceedings of the First Workshop*
996 *on Trolling, Aggression and Cyberbullying (TRAC-2018)*. Association for Computational Linguistics, Santa Fe, New Mexico, USA, 113–119. <https://www.aclweb.org/anthology/W18-4414>
- 997 Motaz Saad. 2010. *The Impact of Text Preprocessing and Term Weighting on Arabic Text Classification*. Ph.D. Dissertation. [https://doi.org/10.13140/2.1.](https://doi.org/10.13140/2.1.4677.2164)
998 [4677.2164](https://doi.org/10.13140/2.1.4677.2164)
- 999 Ali Safaya, Moutasem Abdullatif, and Deniz Yuret. 2020. KUISAIL at SemEval-2020 Task 12: BERT-CNN for Offensive Speech Identification in Social
1000 Media. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*. International Committee for Computational Linguistics, Barcelona (online),
1001 2054–2059. <https://www.aclweb.org/anthology/2020.semeval-1.271>
- 1002 Abu Bakr Soliman, Kareem Eissa, and Samhaa R. El-Beltagy. 2017. AraVec: A set of Arabic Word Embedding Models for use in Arabic NLP. *Procedia*
1003 *Computer Science* 117 (2017), 256 – 265. <https://doi.org/10.1016/j.procs.2017.10.117> Arabic Computational Linguistics.
- 1004 HoSung Woo, JaMee Kim, and WonGyu Lee. 2020. Validation of Text Data Preprocessing Using a Neural Network Model. *Mathematical Problems in*
1005 *Engineering* 2020 (2020), 1958149. <https://doi.org/10.1155/2020/1958149>

1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040