

Building Datasets for Dialect Classifiers using Twitter

Anonymous ACL submission

Abstract

This paper outlines a scalable approach to creating manually annotated dialectal data sets from Twitter, at little to no cost, intended for a variety of machine learning models that may rely on dialect identification. We distinguish between two Arabic dialects, Egyptian and Gulf, using a variety of dialectal terms, and store an evenly balanced corpus of 300 thousand tweets. Experimental results show that benchmark classifiers trained using the approach outlined in this paper can determine in which dialect a given sentence was written with impressive accuracy, further revealing an important semantic and social connection.

1 Introduction

There are nearly 500 million Arabic speakers in the world, making it the largest member of the Semitic language family. Arabic is also the fourth most used language on the Internet with over 175 million users (Internet World Stats, 2017), and yet despite its significant influence, the Arabic language has garnered comparatively little attention in modern computational linguistics (Stanford NLP Group, 2017). The majority of published Arabic is in Modern Standard Arabic (MSA), which differs substantially from regional dialects (Zaidan, 2014). Aside from the speech domain, linguistic research tends to heavily rely on MSA data sets (Zaidan, 2014). Progress has been made to create Arabic dialectal data sets that can be used to train and evaluate classifiers for dialect identification. Such classifiers have also been found to outperform baselines that are solely reliant on MSA (Zaidan, 2014) and can be used for speech recognition systems, machine translation, personalizing search engine results and digital advertising. Dialectal Arabic is more

often present in informal settings such as blogs, comments, forums, and social media. The Arabic On-line Commentary Data Set (AOC) (Zaidan, 2014), and the Cross Lingual Arabic Blog Alerts (COLABA) project (Diab et al. 2010), are efforts to derive dialectal Arabic from such sources. As far as we know, there are no manually annotated Arabic dialectal data sets created entirely from social media sources. “In terms of language on Twitter, tweets in Arabic account for over 72% of all tweets generated in the Arab world.” (Arab Social Media Report, 2017).

We propose a novel way to build dialectal data sets from Twitter that can be used to train highly accurate models to determine in which dialect a given tweet was written. Using the Twitter API, we devise a flexible approach to identify, collect, and annotate tweets primarily written in two Arabic dialects, Egyptian and Gulf. Our method does not incur a huge cost and can be used to annotate thousands of tweets at a time. The approach utilizes non-polysemous, dialectal keywords to avoid capturing tweets containing terms that are spelled similarly but pronounced differently, or mean different things across dialects, thereby improving the likelihood of a collected tweet belonging exclusively to the dialect of interest. We initially conduct exploratory data analysis (EDA) and score tweets based on lexical diversity, analyzing the results to identify users more likely to frequently use relevant dialectal terms. We then inspect the resulting user timelines to collect up to 3,000 tweets (the API’s rate limit) and annotate them in the process. Once we have a large enough sample of annotated tweets, we combine the Egyptian and Gulf data sets into a single corpus and conduct topic modeling using Latent Dirichlet allocation (LDA) to reveal linear combination of terms that characterize or separate the two classes, visualizing the various topics through the tool

100 LDavis. With Egyptian and Gulf as the response
 101 variables (classes) and the collected tweets
 102 (documents) as the predictor variables, we finally
 103 use a benchmark Naive Bayes on a test set of
 104 75,735 manually annotated tweets to achieve a
 105 dialect prediction score over 92%.

107 Building the corpus

109 Using the Twitter Streaming API, we initially
 110 filter real-time tweets using a region specific
 111 keywords list created for each dialect, Egyptian
 112 and Gulf, referred to as EG and GLF.

- 113 ▪ GLF terms for Hijazi, Najadi, Omani,
 114 Kuwaiti, Emirati, Bahraini.
- 115 ▪ EG terms for various regional dialects,
 116 such as Alexandrian, Sa'idi, Isma'ilī,
 117 Badawi, Cairene.

120 Code from the book *Mastering Social Media*
 121 *Mining with Python* (Bonzanini, 2016), was used
 122 to stream tweets from the Twitter API. The
 123 dialectal keyword lists were created from a
 124 selection of terms available on *Mo3jam.com*, and
 125 additionally solicited from Native speakers of
 126 those dialects.

127 *Polysemy in the Arabic Dialects* (Ibrahim, 2005)
 128 was used to research terms that co-occur in
 129 various Arabic dialects and have different
 130 meanings. This was part of an ongoing effort to
 131 separate the EG and GLF sub-corpora as much as
 132 possible, by avoiding polysemous terms when
 133 streaming tweets to ensure a more truthful
 134 representation of the dialect being queried.
 135 Furthermore, collected tweets were inspected for
 136 homographs during EDA, in order to ensure
 137 sufficient distance between the sub corpora before
 138 the manual annotation of the two classes.

139 EG	GLF
140 السكاك (el sekkah)	الطریق (el tareeq)
141 ازیک (ezayek)	خربز (kharbez)
142 دوافر (dawafer)	ازبن (azben)
143 ماناخیر (manakheir)	ابخاس (abkhas)
144 كوبايا (kobaya)	شخل (shkhat)
ا زای (ezay)	البشكارة (el bashkarah)
تسبس (tahyeis)	مسلاحة (maghsalah)

145 Figure 1: Sample of non-polysemous terms used to
 146 query dialectal tweets for each class.

147 After collecting at least 100 tweets from each
 148 keyword-filtered stream, we identify users

150 that contribute the most tweets to those
 151 streams, as well as users with high lexical
 152 diversity scores. A simple function finds the
 153 length of unique terms in the text and divides
 154 it by the length of the entire text, and is used
 155 to mask tweets with lexical diversity scores
 156 over 0.4, while an additional mask is used to
 157 exclude tweets under 40 characters.

	cleaned_name	cleaned_text_score	cleaned_text	cleaned_text_len
7	h_sawires	0.411765	أزي اللبس كانت ينكره في المصادر قبل إلخراج لويفر	51
44	ttoozaahran	0.410714	مثل عراقة مصرى الساحة عشق العصر سكنى أزاي و الساحة تلويني	56
49	Abdelaziz_	0.431818	أنت غريب المقام كأزي التلبيخ مثلاه الله	44
56	Ejal3mri0	0.426230	الله تكلم كل ما شرطك عزيزاً غير ضلالي بعمره أزاي عذباً	61
73	mariam_red	0.511628	أزي العرف على الفضل مصطفى عذنك يا - جامده	43

159 Figure 2: From left to right (Twitter username,
 160 lexical diversity score, tweet content, tweet
 161 length).

162 To ensure that the tweets truthfully represent
 163 the dialect being annotated, we manually
 164 inspect a user's timeline based on simple
 165 criteria before adding their tweets to the sub-
 166 corpora. A user's own tweets (excluding
 167 retweets and mentions) should be heavily
 168 inundated with dialectal terms and should not
 169 contain too many retweets or spam-like
 170 behavior. Users timelines should also not
 171 contain high volumes of repetitive tweets or
 172 tweets that are cross-cultural.

173 For example, user timelines with high
 174 volumes of religious verses are avoided, as
 175 similarity in this content across dialects adds
 176 noise to the data. Prioritization was given to
 177 user timelines with tweets focused on specific
 178 topics, such as fashion, politics, humor,
 179 parenting, romance, and sensitive subjects.
 180 We then conduct a Timeline Search using
 181 Twitter's REST API and collect the most
 182 recent 3,200 tweets (rate limit for the API)
 183 from qualified users. Additional users whose
 184 tweets appear on those chosen timelines can
 185 also often be highly dialectical, and upon
 186 further inspection can be added to the
 187 timeline search.

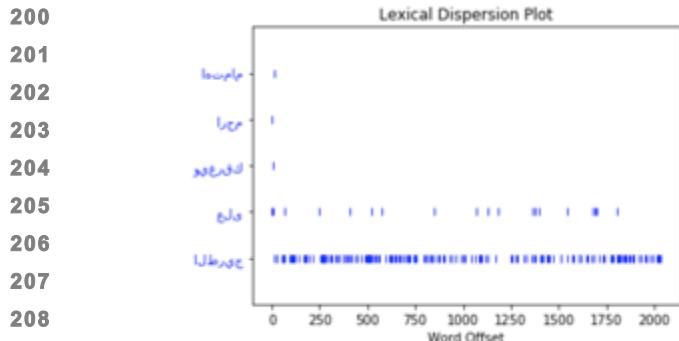


Figure 3: Plotting of the GLF subcorpus using EG and GLF terms reveals that GLF terms have a much higher lexical dispersion, further confirming the lack of impurity in the different subcorpora.

We check for GLF terms in the EG subcorpus and vice versa by exploring them in a Lexical Dispersion plot. In Figure 3, low word offset for EG terms in the GLF subcorpus is a way of confirming whether or not the classes are sufficiently separated.

2 Text analysis and visualization

After manually annotating the data according to dialect class, we conduct EDA on the entire corpus.

- Corpus size: 300k (evenly split and labeled EG/GLF).
- Vocabulary size: 356,726

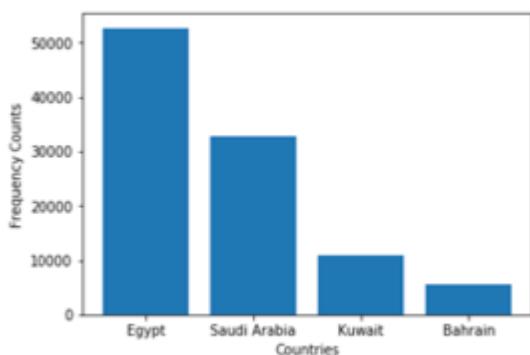


Figure 4: The geographic origins of tweets show a balance between classes (i.e.: The number of tweets originating from Egypt is roughly equal to all the Gulf countries combined).

Dialectal terms found during EDA but not included in the initial keywords list were used

to gather more tweets from the real-time streamer. Alternatively, polysemous terms appearing in high frequencies can be added to a stop words list before vectorizing the data.

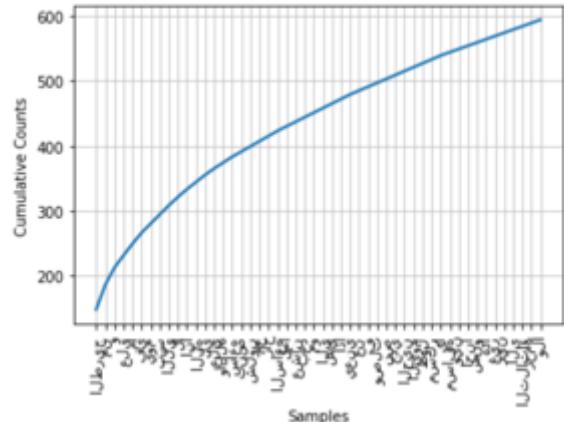


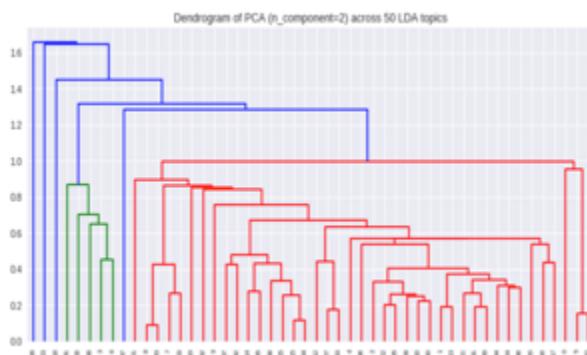
Figure 5: By determining the most frequent terms in our corpus, using the NLTK library we can determine the cumulative value counts of said terms.

3 Topic Modeling

Stop words can have a significant impact on topic modeling and the performance of a classifier. With that in mind, stop words were removed from the corpus before conducting topic modeling and the train test split for classification, decreasing the total corpus size by -18.93%. The stop words list was obtained from the paper *Effects of Stops Words Elimination for Arabic Information Retrieval: A Comparative Study*, which examines and compares the effects of term weighting and stop words on Arabic retrieval. The paper found that the overall performance of the 1590-word General Stoplist generally performed better than other lists developed for the paper (El-Khair, 2017). We pass the resulting corpus through Gensim's Dictionary representation to conduct topic modeling. Additionally, we remove terms that occur too frequently (over 60%) or too rarely (less than 20 occurrences in total).

300 **3.1 Latent Dirichlet allocation (LDA)**
 301
 302 We use Gensim’s doc2bow method to get a
 303 bag of words representation, and then
 304 initialize a Latent Dirichlet allocation to
 305 explore the linear combinations of the
 306 features that separate the two dialects. We
 307 interactively visualize 50 topic models that
 308 constitute each dialect by utilizing the
 309 LDAvis tool (LDA topic models available
 310 [here](#)).
 311

312 The topics are further broken down with
 313 Principle Component Analysis (PCA) into 2
 314 principle components, and plotted with a
 315 dendrogram representing 50 LDA topics.
 316
 317



328 Figure 6: Topics appearing closely together are
 329 more likely to represent the same dialect.
 330

332 **3.2 Latent Semantic Analysis (LSA)**

335 To further analyze dialectal relationships
 336 across the different tweets, we initialize a
 337 term frequency-inverse document frequency
 338 (TF-IDF) model using Gensim to transform
 339 the corpus for use with LSA. We then
 340 initialize an LSA transformation with 500
 341 topics to find cosine similarities between the
 342 LSA topics and a query. Using a dialectal
 343 tweet belonging to the EG class as a test, we
 344 turn the tweet it into a bag of words to be
 345 used as a query with LSA and project it onto a
 346 2-d vector space (topics), where 0 represents
 347 EG and 1 represents GLF.
 348
 349

350 مليش دعوه بحد بصراره ولا نادى ليه دعوه بصفقات حد "350
 351 السوق"
 352

353 (0, 0.256), (1, -0.065)
 354

355 We then perform a similarity query against
 356 the first 10 entries in the corpus, which have
 357 been annotated with the class EG.
 358

359 [(0, 0.97553933),
 360 (1, 0.979222),
 361 (2, 0.99163806),
 362 (3, 0.97607386),
 363 (4, 0.99301887),
 364 (5, 0.96669626),
 365 (6, 0.98321593),
 366 (7, 0.56111413),
 367 (8, 0.56111413),
 368 (9, 0.99779427)]
 369
 370

371 Figure 7: The first 10 documents (tweets) are
 372 labeled with class EG and the query itself is
 373 written in EG. The results show high similarity
 374 with most documents and two possible
 375 misclassifications (documents 7, 8).
 376

377 **4 Classification (EG & GLF)**

378 We label encode the classes EG and GLF and
 379 use the General Stoplist (El-Khair, 2017) to
 380 remove stop words from the corpus. A
 381 train/test split is conducted and the value of
 382 the test is set to 0.25. The resulting text is run
 383 through a pipeline that includes a Count
 384 Vectorizer, TF-IDF Transformer, and a
 385 Multinomial Naive Bayes classifier with
 386 default (benchmark) sklearn hyperparameter
 387 settings.
 388

389 **4.1 Classification Report**

390 A heatmap was created to illustrate the results
 391 of our Classification Report in Figure 8,
 392 showing the performance of our models’
 393 recall, precision, and f1-score (the harmonic
 394 mean of precision and recall), all of which are
 395 over 90%.
 396
 397

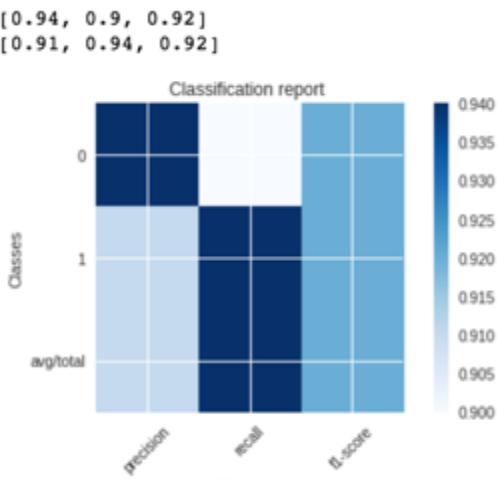


Figure 8: EG shown as 0 and GLF shown as 1. With 0.94 precision and 0.9 recall score for EG, and 0.91 precision and 0.94 recall for GLF. The F1 score for both classes is 0.92.

4.2 Confusion Matrix (raw values)

The raw values of the classification results are further explored through a confusion matrix to quantify the classifier’s performance.

Predicted	0	1	All
True	0	1	All
0	34328	3631	37959
1	2293	35483	37776
All	36621	39114	75735

Figure 9: A total of 75,735 tweets are used for the test set, with 34,328 predicted accurately for EG and 35,483 predicted accurately for GLF.

4.3 Confusion Matrix (percentages)

The raw values of the confusion matrix are converted to percentages in figure 10, showing that the classifier misclassified EG for GLF 4.64% of the time and GLF for EG only 3.13%.

Predicted	0	1	All
True	0	1	All
0	46.87	4.64	25.06
1	3.13	45.36	24.94
All	50	50	50

Figure 10: Out of all EG tweets, 46.9% were correctly classified, while 45.4% of GLF tweets were correctly classified.

4.4 AUROC

To explore which models are best at predicting the classes, we plot the Area Under the Receiver Operating Characteristic curve (AUROC). This allows us to combine the True positive rate (TPR) and the False positive rate (FPR) into a single metric, an AUROC score of 0.97.

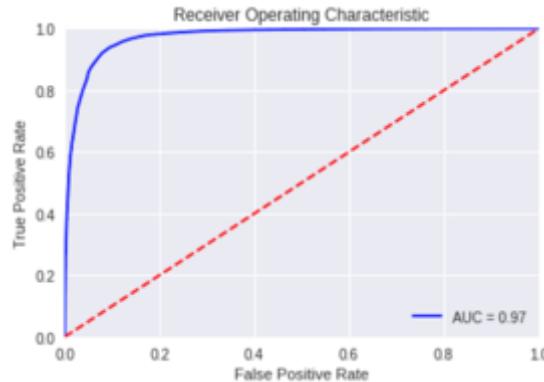


Figure 11: An optimal threshold can be achieved under a FPR of 0.2.

5 Conclusions and Future Work

Using a benchmark Naive Bayes on a test set of 75,735 manually annotated test tweets achieved a 92% dialect prediction score, with high AUC/ROC (0.97) and high f1 scores (.92 for both classes), suggesting that the classifier can do a good job for many other values of

500	threshold. These results show that dialect	550
501	trained models can improve accuracy when	551
502	compared to a system trained solely on an	552
503	MSA-English parallel corpus.	553
504		554
505	Social media is a good source of dialectal data	555
506	that can be used to train dialect identification	556
507	classifiers. By refining models created for this	557
508	paper, we intend to prototype dialect identifier	558
509	tools that will be available online and	559
510	incorporate a user feedback feature. As the	560
511	corpus size was intended to build a proof of	561
512	principle, we will continue expanding the	562
513	corpus and adding new dialects and	563
514	languages. Experimenting with different term	564
515	weights, models, and stop words, as well as	565
516	conducting deep learning text classification	566
517	using word embeddings (Keras) can further	567
518	improve results. Gensim also includes	568
519	numerous features that were not utilized,	569
520	including use word2vec, doc2vec, BM25.	570
521	Tagging the words with The Stanford Arabic	571
522	Parser, which includes a built-in Arabic Tree	572
523	Bank (Maamouri, 2005) and Arabic	573
524	Segmenter can also be used before term	574
525	weighting to only keep nouns and verbs,	575
526	which can further improve results.	576
527		577
528		578
529		579
530	References	580
531		581
532		582
533	Ibrahim Abu El-Khair. 2017. <i>Effects of Stop Words</i>	583
534	<i>Elimination for Arabic Information Retrieval: A</i>	584
535	<i>Comparative Study</i> . International Journal of	585
536	Computing & Information Sciences 4 (3), 119-	586
537	133. http://adsabs.harvard.edu/cgi-bin/bib_query?arXiv:1702.01925	587
538	Marco Bonzanini. 2016. <i>Mastering Social Media</i>	588
539	<i>Mining with Python</i> . Packt Publishing.	589
540	Birmingham, United Kingdom.	590
541	Mohamed Maamouri (project head). 2005. <i>Arabic</i>	591
542	<i>Treebank</i> . Linguistic Data Consortium, PA.	592
543	Mona Diab and Nizar Habash. Natural Language	593
544	Processing of Arabic and its Dialects. In <i>EMNLP</i>	594
545	<i>2014, Doha, Qatar</i> .	595
546	Mo3jam. 2017. User-generated dictionary of	596
547	colloquial Arabic. https://ar.mo3jam.com/	597
548		598
549		599