

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Classifying Text Messages for Emergency Response

Cornelia Caragea, Hyun-Woo Kim, Prasenjit Mitra, and John Yen
College of Information Sciences and Technology
Pennsylvania State University
University Park, PA-16801
{ccaragea,hxk263,pmitra,jyen}@ist.psu.edu

Abstract

In case of emergencies (e.g., earthquakes, flooding), rapid responses are needed in order to address victims' requests for help. Hence, the ability to classify tweets and text messages *automatically*, together with the ability to deliver the relevant information to the appropriate personnel are essential for enabling the personnel to timely and efficiently work to address the most urgent needs, and to understand the emergency situation better. The choice of *features* used to encode tweets and text message data is crucial for the performance of the learning algorithms. Here, we present a comparative study of four types of feature representations used to enable learning classifiers from such data. These feature representations are obtained using a "bag of words" approach, feature abstraction, feature selection, and Latent Dirichlet Allocation (LDA). The results of our experiments on a real-world text message data set show that feature abstraction can yield better performing models than those obtained by using a "bag of words", feature selection and LDA.

1 Introduction

The 7.0 Earthquake in Haiti has mobilized the entire world to support the relief effort, especially through novel uses of the cyberspace. Relief workers, reporters, and non-governmental organizations (NGOs) have used tweets and text messages extensively to spread and share information about the needs, events, and casualties in the Twitterworld. Regular citizens have also employed Twitter to rally others to support relief efforts. Both Haitians and relief workers have used mobile phones to send text messages regarding damages, resource needs, and security-related events. While there is useful information in these tweets and text messages, they are not well-organized to allow critical information (e.g., water, medical supply, food) to be delivered to those who need them in a timely and efficient fashion. Hence, the ability to classify tweets and text messages *automatically*, together with the ability to deliver the relevant information to the appropriate personnel are essential for enabling the personnel to timely and efficiently work to address the most urgent needs, and to understand the emergency situation better in the Emergency Response Sector.

Although tweets and text message classification can be performed with little or no effort by people, it still remains difficult for computers. Machine learning currently offers a promising approach to the design of algorithms for training computer programs to efficiently and accurately classify short text message data. Some of the main challenges in classifying such data are as follows: (i) tweets and text messages contain only a few words and, sometimes, require background information for accurate classification. For example, the message "*I live in Leogane, Route de Mellier Bongnotte #72, I need formula for my baby.*" requires knowledge that *formula* refers to *baby food*. The choice of features used to encode such data is crucial for the performance of the learning algorithms; (ii) tweets and text messages may belong to multiple categories, i.e., the *multi-label classification*; (iii) there may be possible errors in the manually generated labels (i.e., categories) of text messages, which can impact the performance of the learning algorithms; (iv) the training set is often limited in size.

054 In this study, we focused on the choice of features that are used to represent short text messages.
055 We used four types of feature representations to enable learning Naïve Bayes and Support Vec-
056 tor Machine classifiers to accurately classify text messages from Haiti earthquake, submitted to
057 Ushahidi-Haiti (<http://haiti.ushahidi.com>) through phone, e-mail, Twitter, or web. These feature
058 representations are obtained using: (i) a “bag of words”, i.e., all words in the vocabulary [11]; (ii)
059 feature abstraction methods, that find a partition of the set of words in the vocabulary by clustering
060 words based on the similarity between the class distributions that they induce [15]; (iii) feature se-
061 lection methods, that select a subset of features based on some chosen criteria [8]; and (iv) Latent
062 Dirichlet Allocation (LDA) [1], that finds hidden topics in the data. The topic words, i.e., the words
063 in each topic, can be seen as a set of discriminative features.

064 We compared the performance of the trained classifiers using the above feature representations on
065 a real-world text message data set from Ushahidi-Haiti. The results of our experiments show that
066 feature abstraction generates features that can yield better performing classifiers than those obtained
067 by using a “bag of words”, features chosen by feature selection, and features as topic words output by
068 LDA. We also discuss the insights gained from these results and suggest directions of future research
069 to enhance the accuracy and the coverage for classifying tweets and text messages for improved
070 efficiency and coordination during the response, transition, and recovery of extreme events.
071

072 2 Methods

073
074 In this section, we describe the three feature representation methods, which are compared with the
075 “bag of words” approach: (1) feature abstraction; (2) feature selection; and (3) Latent Dirichlet
076 Allocation (LDA).
077

078 **Feature Abstraction.** Feature abstraction methods are potentially successful techniques for pro-
079 ducing appropriate features for classification [15]. They reduce the classifier input size by grouping
080 “similar” features to generate *abstract features* (also called abstractions). Silvescu et al. [15] pro-
081 posed an approach to simplifying the data representation used by a learner by grouping features
082 based on the Jensen-Shannon divergence [3] that result in minimal reduction in the mutual informa-
083 tion between features and the class variable.

084 Specifically, they used hierarchical agglomerative clustering to group the most “similar” features
085 at each step of the algorithm, based on the similarity between the conditional distributions of the
086 class variable given the features. The most “similar” features are identified as those that have the
087 smallest Jensen-Shannon divergence between the conditional distributions of the class that the fea-
088 tures induce. As an effect, abstract features that are predictive of the class variable are obtained. An
089 example of an abstract feature can be “food”, which is more general than the specific features “rice”
090 and “formula” (i.e., baby food). The abstract feature is identified by the group $\{rice, formula\}$.

091 Silvescu et al. [15] have shown that *abstraction* reduces the model input size and helps improve the
092 statistical estimates of complex models (especially when data are sparse) by reducing the number of
093 parameters to be estimated from data. In this study, we have applied the feature abstraction approach
094 of Silvescu et al. to generate the abstract feature representation.

095 **Feature Selection.** Feature selection methods attempt to remove redundant or irrelevant features
096 in order to improve classification performance of learning algorithms [5]. Feature selection selects
097 a subset of the available features based on some chosen criteria, and can substantially reduce the
098 number of model parameters. Kira and Rendell [8] proposed an algorithm for feature selection,
099 called Relief, which is not heuristic-based, is robust to noise and to interaction among features.

100 Relief is a weight-based algorithm. At each step, Relief samples from the training data an instance
101 x , and determines x 's *near-hit* (the closest instance from the same class as x) and *near-miss* (the
102 closest instance from the opposite class of x) in the training data, by using p -dimensional Euclidian
103 distance. A feature weight vector is updated for each such triplet to determine the relevance of all
104 features to the class variable. The algorithm terminates after k steps and returns those features whose
105 relevance level is above some user-specified threshold. We have used the Relief algorithm to select
106 a subset of features that are predictive to the class variable.

107 **Latent Dirichlet Allocation.** Latent Dirichlet Allocation (LDA) is an unsupervised method for
detecting hidden topics in the data proposed by Blei et al. [1]. LDA is a generative probabilistic

108 model of a collection of documents, which has been successfully used to perform dimensionality
109 reduction for text classification [20], where documents are multiple paragraphs and pages in length.

110 LDA [1] models each document in a collection as a mixture of topics (drawn from a conjugate
111 Dirichlet prior), and each topic as a distribution over words in the vocabulary. The topic distribution
112 of a document can be seen as a lower dimensional representation of the document (where the dimen-
113 sionality is equal to the number of topics). Furthermore, the union of the words with high probability
114 in each topic can be seen as a set of discriminative features for the collection of documents. We have
115 used these words to generate the topic words feature representation.

116 3 Experiments and Results

117
118
119
120 **Ushahidi Text Message Data Set.** The data set used in our experiments is the Ushahidi data set
121 (<http://haiti.ushahidi.com/>), which consists of 3598 text messages from Haiti earthquake. We used a
122 subset of 2116 text messages of the Ushahidi data set, for which the English translation is available.
123 While text messages are available in both Haitian Kreyol and English languages, we used only the
124 English version, as Munro and Manning [13] found no significant improvement from one language
125 to another on a similar task. The messages are classified into 10 categories: (1) *medical emergency*;
126 (2) *people trapped*; (3) *food shortage*; (4) *water shortage*; (5) *water sanitation*; (6) *shelter needed*;
127 (7) *collapsed structure*; (8) *food distribution*; (9) *hospital/clinic services*; (10) *person news*. Note
128 that a message may belong to multiple categories. For example, the message “*Good evening ONG,*
129 *I’m very happy for the aid you’re giving to the people, I thank you. But in my zone that’s to say*
130 *Lamenten 54 Rue St Juste we need shelter and food.*” belongs to both *shelter needed* and *food*
131 *shortage* categories.

132 **Experimental Design.** Our experiments are designed to explore what feature representations of
133 short text messages, which are provided as input to machine learning classifiers, result in best classi-
134 fication performance. We used four types of feature representations to enable learning Naïve Bayes
135 and Support Vector Machines (SVM) classifiers on the Ushahidi text message data set:

- 136 • a bag of words representation, i.e., all words in the vocabulary. After stemming and re-
137 moving stop words, and words with document frequency less than 3, the vocabulary size is
138 1525 (BoW) [11];
- 139 • a bag of m words chosen using the RELIEF feature selection method (FS) [8];
- 140 • a bag of m abstractions over all words in the vocabulary, i.e., an m -size partition of the vo-
141 cabulary obtained by grouping words into m abstract terms based on the similarity between
142 the class distributions that they induce (FA) [15];
- 143 • a bag of m topic words output by Latent Dirichlet Allocation (LDA) as the top 20 words
144 from k topics (the number of topic words m is bounded by $20 \times k$) (TW) [1].

145 In our experiments, we used WEKA implementation [6] of Naïve Bayes Multinomial and SVM
146 with the default parameters, and MALLET implementation [10] of LDA. The LDA parameters are
147 set to default, except for the number of iterations of Gibbs sampling, which is set to 3,000, and the
148 random seed, which is set to 1. The number of topics k is set to 9 (chosen to be close to the number
149 of categories in the data set). This results in $m = 165$ topic words. Hence, we trained classifiers for
150 $m = 165$ for all of the above feature representations. In the case of feature abstraction, the 165-size
151 partition of the vocabulary produces classifiers that use smaller number of “features” compared to
152 the “bag of words” representation, i.e., 1525 words, and at the same time, the model compression is
153 not very stringent so as to lose important information in the data through *abstraction*).

154 Because a text message may belong to one or more categories, we trained 10 “one vs. others” binary
155 classifiers, one for each category. For all experiments, we report the average F1 Measure obtained
156 in a 5-fold cross-validation experiment.

157 **Results.** Table 1 shows the comparison of average F1 Measure (along with 95% confidence inter-
158 vals) using binary SVMs and Naïve Bayes trained on the Ushahidi text message data set for each of
159 the ten categories. The feature representations used to train the classifiers are as follows: (i) “bag of
160 words” (BoW); (ii) abstractions used as “features” in the classification model, which are obtained by
161 feature abstraction (FA); (iii) features selected by Relief feature selection (FS); and (iv) topic words,
output by LDA (TW).

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
209
210
211
212
213
214
215

Class	Support Vector Machines			
	BoW	FA	FS	TW
medical emergency	0.29 ± 0.06	0.27 ± 0.08	0.12 ± 0.07	0.11 ± 0.05
people trapped	0.68 ± 0.11	0.74 ± 0.09	0.64 ± 0.14	0.62 ± 0.23
food shortage	0.71 ± 0.02	0.73 ± 0.03	0.71 ± 0.06	0.72 ± 0.07
water shortage	0.66 ± 0.03	0.67 ± 0.02	0.63 ± 0.04	0.65 ± 0.03
water sanitation	0.91 ± 0.01	0.94 ± 0.01	0.96 ± 0.01	0.95 ± 0.01
shelter needed	0.52 ± 0.02	0.52 ± 0.05	0.44 ± 0.07	0.48 ± 0.04
collapsed structure	0.42 ± 0.08	0.33 ± 0.15	0.31 ± 0.16	0.39 ± 0.20
food distribution	0.27 ± 0.05	0.27 ± 0.03	0.18 ± 0.07	0.17 ± 0.09
hospital/clinic services	0.56 ± 0.04	0.59 ± 0.06	0.47 ± 0.08	0.51 ± 0.05
person news	0.55 ± 0.06	0.59 ± 0.04	0.39 ± 0.10	0.45 ± 0.04

Class	Naïve Bayes			
	BoW	FA	FS	TW
medical emergency	0.29 ± 0.09	0.31 ± 0.06	0.25 ± 0.11	0.14 ± 0.03
people trapped	0.67 ± 0.06	0.71 ± 0.08	0.71 ± 0.13	0.65 ± 0.09
food shortage	0.77 ± 0.03	0.76 ± 0.02	0.73 ± 0.02	0.75 ± 0.04
water shortage	0.69 ± 0.02	0.67 ± 0.02	0.66 ± 0.03	0.66 ± 0.02
water sanitation	0.94 ± 0.01	0.94 ± 0.01	0.95 ± 0.01	0.95 ± 0.01
shelter needed	0.45 ± 0.09	0.46 ± 0.09	0.38 ± 0.10	0.44 ± 0.03
collapsed structure	0.35 ± 0.10	0.45 ± 0.12	0.43 ± 0.14	0.41 ± 0.11
food distribution	0.26 ± 0.05	0.27 ± 0.04	0.20 ± 0.07	0.20 ± 0.09
hospital/clinic services	0.59 ± 0.05	0.61 ± 0.04	0.54 ± 0.08	0.55 ± 0.04
person news	0.61 ± 0.05	0.61 ± 0.05	0.48 ± 0.10	0.52 ± 0.06

Table 1: Comparison of average F1 Measure (with 95% confidence intervals) obtained in 5-fold cross-validation experiments using binary Support Vector Machines and Naïve Bayes classifiers trained on the Ushahidi text message data set for each of the ten classes. The feature representations used to train the classifiers are as follows: (i) “bag of words” (BoW); (ii) abstractions used as “features” in the classification model, which are obtained by feature abstraction (FA); (iii) features selected by Relief feature selection (FS); and (iv) topic words, output by LDA (TW).

As can be seen from the table, FA significantly outperforms BoW for most of the categories from the Ushahidi text message data set, using both Naïve Bayes and SVM classifiers. For few categories, for example *shelter needed*, FA-based SVM matches the performance of BoW-based SVM with substantially smaller number of features, i.e., 165 and 1525 features are used for training FA-based SVM and BoW-based SVM, respectively.

Compared to FS and TW, FA significantly outperforms both of them for the same number of features used in the classification model, for both SVM and Naïve Bayes classifiers, on all categories except *water sanitation*. Although topic models have been successfully applied to documents that are multiple paragraphs and pages in length, we found that they do not work very well when applied to short text messages.

It is interesting to note that the performance of SVM is worse than that of NB for many categories using any of the feature representations used in this study. This could be due to *overfitting* (see [19] for a theoretical analysis of overfitting for the SVM algorithm). However, as already noted, FA-based SVM significantly outperforms BoW-based SVM for many of the categories. This suggests that FA can help minimize *overfitting* (through parameter smoothing).

4 Related Work

The problem of learning classifiers from short text messages has started to receive significant attention in the machine learning literature. Healy et al. [7], Hidalgo et al. [4], and Cormack et al. [2] have previously addressed the problem of identifying *spam* short messages, by employing various

216 machine learning algorithms (such as Naïve Bayes, SVM, Logistic Regression, and Decision Trees)
217 and various feature representations (such as “bag of words”, “bag of words” augmented by statisti-
218 cal features, e.g., the proportion of upper case letters or punctuation in the text, orthogonal sparse
219 word bigrams, character bigrams and trigrams). Gupta and Ratinov [14] have employed transfer
220 learning techniques to classify short online dialogs, by enriching the set of features using external
221 data sources. Munro and Manning [13] have focused on classifying medical text messages, written
222 in Chichewa language, that were received by a clinic in Malawi, and have shown that incorporating
223 morphological and phonological variation could improve classification performance. Furthermore,
224 Munro [12] has presented a brief survey about the crowdsourced translation to English of text mes-
225 sages written in Haitian Kreyol during the January 12 earthquake in Haiti. Collaborating online,
226 people around the world were able to translate more than 40,000 messages in a short time, which
227 led to saving hundreds of lives, and direct the food and medical aid to tens of thousands [12]. Star-
228 bird and Stamberger [16] introduced a Twitter hashtag syntax for reporting events related to crisis.

229 Unlike these works, we focused on determining a subset of features that are most informative for
230 the target variable, either by selecting a subset of features from the entire vocabulary using feature
231 selection or LDA, or by constructing abstract features using feature abstraction. In addition, our text
232 message classification task is harder due to its multi-label nature (i.e., text messages may belong to
233 multiple categories).

234 The topics related to emergency response (ER) form an ontology that can be applied to emergency
235 response for a wide range of relief operations. Ontology development tool such as Protégé has
236 been widely used for developing ontology for different domains. Li et al. (2008) [9] proposed an
237 ontology for emergency response (ER). The top level concepts of the proposed ontology include:
238 aftermath-handling, emergency-rescue, emergency-response, and response-preparation. Each con-
239 cept is further refined by a set of subconcepts. Emergency-rescue, for instance, include medical-aid,
240 evacuation, and victim-assistance. Turoff et al. (2006) have designed a dynamic emergency response
241 management system: DERMIS [18], and have identified the characteristics of a good ER system.

242

243 5 Summary and Discussion

244

245 **Summary.** In this study, we compared four types of feature representations for learning Naïve Bayes
246 and SVM classifiers to *accurately* classify text messages about Haiti disaster relief (originating in
247 Haiti and elsewhere) so that they can be more easily accessed by NGOs, other relief workers, people
248 in Haiti, and their friends and families. These feature representations are: “bag of words”, abstract
249 features (or abstractions), features selected using feature selection, and topic words output by LDA.

250 The results of our experiments on the Ushahidi text message data set show that using *abstract fea-*
251 *tures* makes it possible to construct predictive models that use significantly smaller number of fea-
252 tures than those obtained using a bag of words representation. The resulting models are competitive
253 with, and often significantly outperform those that use the “bag of words” feature representation.
254 Moreover, *abstract features* yield better performing models than features selected by Relief feature
255 selection, and than topic words extracted using LDA.

256 **Discussion.** In learning from real-world text message data, other challenges may be encountered,
257 hence, making the learning problem harder. We point out some of these challenges and provide
258 potential solutions that will be addressed in future work: (i) Tweets and text messages may belong to
259 multiple categories. For example, the text message “*I live in the site Marassa 7. I ask some help like*
260 *water and food thank you*” belongs to both categories *Food Shortage* and *Water Shortage*. Hence,
261 the classification problem can be formulated as a *multi-label* problem [17], where a collection of
262 $|C|$ binary classifiers is trained (where $|C|$ is the number of categories). A test instance is classified
263 using all $|C|$ classifiers. In future work, we will use the *multi-label* formulation. However, in this
264 study, we performed binary classification for each category in order to determine which are the most
265 “difficult” categories to be classified. (ii) As with many real-world data, there may be possible errors
266 in labeling text messages. For example, the text message “*We in Canada turjo quote, we need food,*
267 *water and tents. count on your participation*” belongs to *Food distribution*. However, this example
268 is very similar to “*Good evening ONG, I’m very happy for the aid you’re giving to the people,*
269 *I thank you. But in my zone that’s to say Lamenten 54 Rue St Juste we need shelter and food.*”,
which belongs to *Food Shortage*. In future work, we plan to create a new category that will contain
examples from both *Food distribution* and *Food Shortage*.

270 Furthermore, possible errors in labeling may occur due to the presence of general terms in a text
271 message. For example, the text message “*We need help at Mahotiere 79. Since the catastrophe,*
272 *we have not seen anyone from the government*” is labeled as *Food distribution* and *Water sanitation*
273 in the Ushahidi data set. However, there is no indication of the type of help needed. For example,
274 people at Mahotiere 79 may need medical assistance or shelter. To distribute this message to food
275 and water departments may be very inefficient if the people have other more urgent needs. Instead,
276 we propose to use a general category, which consists of these types of messages. Hence, the general
277 department can *efficiently* determine what the people needs are and act accordingly. Further research
278 may also include: (i) Exploration of other types of abstraction based on semantically related words;
279 (ii) Classification of tweets about Haiti, provided by Twitter.

280 Acknowledgments

281
282 This research is partially supported by NSF RAPID grant IIS 1026763. We would like to acknowl-
283 edge Ushahidi for making the data available for this research, NetHOPE for general interests about
284 this research. We also wish to thank members of EMERSE research team for their contribution
285 to related discussions: Lee Giles, Jim Jansen, Andrea Tapia, Dinghao Wu, Anuj Jaiswal, Gregory
286 Traylor, Anthony Maslowski, Nathan McNeese, and Louis-Marie Ngamassi Tchouakeu.

287 References

- 289 [1] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*,
290 3:993–1022, 2003.
- 291 [2] G. V. Cormack, J. M. Gómez Hidalgo, and E. P. Sánz. Spam filtering for short messages. In *Proceedings*
292 *of the 16th ACM CIKM '07*, pages 313–320, 2007.
- 293 [3] T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley, 1991.
- 294 [4] J. M. Gómez Hidalgo, G. C. Bringas, E. P. Sánz, and F. C. García. Content based sms spam filtering. In
295 *DocEng '06: Proceedings of the 2006 ACM symposium on Document engineering*, pages 107–114, 2006.
- 296 [5] I. Guyon and A. Elisseeff. An introduction to variable and feature selection. *J. Mach. Learn. Res.*,
297 3:1157–1182, 2003.
- 298 [6] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten. The weka data mining
299 software: An update. *SIGKDD Explorations*, 11, 2009.
- 300 [7] M. Healy, S. J. Delany, and A. Zamolotskikh. An assessment of case base reasoning for short text message
301 classification, 2005.
- 302 [8] K. Kira and L. A. Rendell. The feature selection problem: Traditional methods and a new algorithm. In
303 *Proceedings of AAAI*, pages 122–134, San Jose, CA, 1992. Morgan Kaufmann Publishers Inc.
- 304 [9] X. Li, G. Liu, A. Ling, J. Zhan, N. An, L. Li, and Y. Sha. Building a practical ontology for emergency
305 response systems. In *Proceedings of IEEE International Conference on CSSE '08*, pages 222–225, 2008.
- 306 [10] A. McCallum. Mallet: A machine learning for language toolkit. <http://mallet.cs.umass.edu>, 2002.
- 307 [11] A. McCallum and K. Nigam. A comparison of event models for naive bayes text classification, 1998.
- 308 [12] R. Munro. Crowdsourced translation for emergency response in haiti: the global collaboration of local
309 knowledge. In *Relief 2.0 in Haiti*, 2010.
- 310 [13] R. Munro and C. D. Manning. Subword variation in text message classification. In *Human Language*
311 *Technologies: The Annual Conference of the North American Chapter of the ACL*, pages 510–518, 2010.
- 312 [14] R. Gupta and L. Ratinov. Text categorization with knowledge transfer from heterogeneous data sources.
313 In *Proceedings of AAAI 2008*, 2008.
- 314 [15] A. Silvescu, C. Caragea, and V. Honavar. Combining super-structuring and abstraction on sequence
315 classification. In *ICDM*, pages 986–991, 2009.
- 316 [16] K. Starbird and J. Stamberger. Tweak the tweet: Leveraging microblogging proliferation with a prescrip-
317 tive syntax to support citizen reporting. In *Proceedings of the 7th Intl. ISCRAM Conf. '10*, 2010.
- 318 [17] G. Tsoumakas and I. Katakis. Multi label classification: An overview. *International Journal of Data*
319 *Warehousing and Mining*, 3(3):1–13, 2007.
- 320 [18] M. Turoff, M. Chumer, B. Van de Walle, and X. Yao. Design of a dynamic emergency response manage-
321 ment information system (dermis). *Journal of Information Technology Theory and Application*, 2004.
- 322 [19] V. Vapnik. *Statistical Learning Theory*. John Wiley & Sons, N.Y., 1998.
- 323 [20] X. Wei and W.B. Croft. LDA-based document models for ad-hoc retrieval. In *Proceedings of ACM SIGIR*,
pages 178–185, 2009.