Knowledge-Based Systems xxx (2015) xxx-xxx

Contents lists available at ScienceDirect

Knowledge-Based Systems

journal homepage: www.elsevier.com/locate/knosys

Please cite this article in press as: M. Franco-Salvador et al., Cross-domain polarity classification using a knowledge-enhanced meta-classifier, Knowl.



Cross-domain polarity classification using a knowledge-enhanced meta-classifier

Marc Franco-Salvador^{a,*}, Fermín L. Cruz^b, José A. Troyano^b, Paolo Rosso^a

^a Pattern Recognition and Human Language Technology (PRHLT) Research Center, Universitat Politècnica de València, Camino de Vera s/n, 46022 Valencia, Spain ^b Department of Languages and Computer Systems, University of Seville, Av. Reina Mercedes s/n, 41012 Sevilla, Spain

ARTICLE INFO

12 814Article history:15Received 30 September 201416Received in revised form 15 May 201517Accepted 18 May 201518Available online xxxx

19 *Keywords:* 20 Sentiment a

5 6

9

Sentiment analysis
Cross-domain polarity classification

22 Meta-learning

23 Word sense disambiguation

24 Semantic network

ABSTRACT

Current approaches to single and cross-domain polarity classification usually use bag of words, *n*-grams or lexical resource-based classifiers. In this paper, we propose the use of meta-learning to combine and enrich those approaches by adding also other knowledge-based features. In addition to the aforementioned classical approaches, our system uses the BabelNet multilingual semantic network to generate features derived from word sense disambiguation and vocabulary expansion. Experimental results show state-of-the-art performance on single and cross-domain polarity classification. Contrary to other approaches, ours is generic. These results were obtained without any domain adaptation technique. Moreover, the use of meta-learning allows our approach to obtain the most stable results across domains. Finally, our empirical analysis provides interesting insights on the use of semantic network-based features.

© 2015 Published by Elsevier B.V.

41 1. Introduction

40

Text classification (also known as text categorization) is the task 42 of assigning a category or categories to a text document from a set 43 44 of predefined categories. Although at first this topic was approached from a knowledge engineering perspective (manually 45 defining a set of rules encoding expert knowledge), in the 90s 46 machine learning became the main approach, and so it stands 47 today. A good survey on machine learning approaches to text 48 classification can be found in Sebastiani [51]. 49

50 The nature of the predefined categories in text classification 51 can be very heterogeneous. The most common task is that of 52 topic-based classification, attempting to classify documents according to their subject matter (e.g. Sports vs. Politics vs. 53 54 Economics). More recently, in the context of the Web 2.0 and social 55 media, it emerged the task of deciding whether a subjective text 56 (typically, a textual review of some product or a cultural or political issue) is positive or negative, depending on the overall sentiment 57 58 detected. This particular task is known as polarity classification or sentiment classification [54,42]. Although it can be defined in terms 59 60 of text classification (being positive and negative the predefined 61 categories) and tackled with similar approaches, polarity classification has been proved to be a more difficult task [42]: while topics 62

> * Corresponding author. E-mail address: mfranco@prhlt.upv.es (M. Franco-Salvador).

Based Syst. (2015), http://dx.doi.org/10.1016/j.knosys.2015.05.020

http://dx.doi.org/10.1016/j.knosys.2015.05.020 0950-7051/© 2015 Published by Elsevier B.V. are often identifiable by keywords alone, sentiment can be expressed in a more subtle manner, and even more when for instance irony is employed [48]. Therefore, solutions based only on bag-of-words representations of documents may not be enough.

In this work we are interested in single and cross-domain polarity classification. Since we are applying machine learning techniques, we start with a training set of documents to build some classifiers. In this context, single-domain classification is the aforementioned common text classification; it refers to training and testing classifiers on the same domain (e.g. movie reviews). Meanwhile, cross-domain classification refers to testing on a different domain (target domain) from that or those used in training (source domains), e.g. training on movie reviews and testing on books reviews. Because manually labeled documents are needed for training, the latter allows to work with domains where no labeled documents are available. The problem of cross-domain text classification was first tackled by Dai et al. [13], and the first results on cross-domain polarity classification were reported by Blitzer et al. [7].

In order to combine different approaches from the research literature and recent knowledge-based approaches, and also to measure the contributions of each one, we propose the use of a meta-learning scheme called Stacked Generalization [56]. The set of base classifiers to be combined using that scheme include solutions used in the past as a TF-IDF bag-of-words classifier, a TF-IDF word *n*-gram classifier, and a lexical resource for opinion

66

67

68

69

70

71

72

73

74

75

76 77

78

79

80

81

82

83

84

85

86

87

88

27

28

29

30

M. Franco-Salvador et al./Knowledge-Based Systems xxx (2015) xxx-xxx

mining-based classifier; but also two new proposals, a word sense disambiguation-based classifier and a vocabulary expansion-based classifier. The latter two classifiers are trained on the basis of knowledge graphs, a subset of a semantic network, i.e., BabelNet [38], focused on the concepts belonging to the text being classified.

The rest of the paper is structured as follows. In Section 2 we 94 95 describe the related work on single and cross-domain polarity clas-96 sification. In Section 3 we introduce our new knowledge-enhanced 97 meta-classifier. In Section 4 we evaluate our approach in the tasks 98 of single and cross-domain polarity classification, and compare it with other state-of-the-art approaches. In that section we evaluate 99 also the performance of our different base classifiers. Finally, in 100 Section 5 we draw the conclusions and mention directions for 101 future work. 102

103 2. Related work

104 The first experiments on single-domain polarity classification 105 using machine learning techniques were performed by Pang et al. 106 [42]. They used a movie review dataset extracted from IMDb.¹ 107 They concluded that polarity classification achieves worse results 108 than other text classification tasks when applying the standard 109 machine learning techniques. Another interesting conclusion was that using unigram presence instead of unigram frequency leads to 110 better results, contrary to observations in other works on text classi-111 112 fication [33].

113 Recent works on polarity classification use the Multi-Domain 114 Sentiment Dataset [7] for evaluation. In its last version, the 115 resource is composed by Amazon product reviews of 25 product 116 types, though most works report results on only the four domains used by Blitzer et al. [7]: Books, Electronics, DVDs and Kitchen 117 118 appliances. Focused on single-domain polarity classification, 119 Dredze et al. [16] presented a new online learning method named 120 confidence-weigthed learning. The method is based on measuring the confidence of each parameter of the classifier; less confident 121 parameters are updated more aggressively than more confident 122 ones. They performed experiments on standard datasets related 123 to different text classification tasks, reporting very good results 124 125 for the Multi-Domain Sentiment Dataset. Another approach, pro-126 posed by Li and Zong [30], use *n*-grams combined with Binormal 127 Separation [22], an alternative to TF-IDF to select the optimal set 128 of features. They reported interesting results in single domain 129 classification.

130 Cross-domain polarity classification has gained popularity 131 thanks to the advances in domain adaptation [14,6,4]. These 132 techniques make use of labeled data from a source domain, and 133 unlabeled data from source and target domains to train their clas-134 sifiers. Using the different domains available in the Multi-Domain 135 Sentiment Dataset, Blitzer et al. [7] was also the first to report 136 results on cross-domain classification proposing two algorithms: 137 structural correspondence learning (SCL), and its variant using 138 mutual information (SCL-MI). The SCL model selects pivot (uni-139 gram and bigram) features frequently appearing in both source 140 and target domains. Then it learns to predict those pivot features 141 in the unlabeled data from both domains. Later, a singular value 142 descomposition is performed to reduce dimensions, and a binary 143 classifier is trained to determine the polarity. Similarly, interesting results on cross-domain polarity classification have been reported 144 145 by spectral feature alignment (SFA) [41]. Using unigram and 146 bigram features, the model exploits the mutual information 147 between each feature and the domain label to differentiate 148 domain-specific and domain-independent features. Next, a bipar-149 tite graph is constructed by dividing both types of features. An edge connects features from different types if there exists co-occurrence. Finally, a spectral clustering is performed to generate feature clusters and a binary classifier is built for the polarity classification. More recently, Bollegala et al. [8,9] used a cross-domain lexicon creation to generate a sentiment-sensitive thesaurus (SST) that groups different words expressing the same sentiment, using also unigram and bigram features as representation. This approach

also obtained competitive results in single-domain polarity classification.

Note that all cross-domain approaches use domain adaptation techniques extracting relevant features from the source domains, in order to obtain important features to classify the target domain. In contrast, we do not use unlabeled data from the target domain. Our approach is focused on proposing new knowledge-based features which allows for training models using the source domains that are able to be directly applied to the target domain. In Section 4.4 we compare our approach in the task of single-domain polarity classification against SST and the state-of-the-art approaches proposed by Dredze et al. [16] and Li and Zong [30]. Next, in Section 4.5 we compare our approach in the task of cross-domain polarity classification against SCL-MI, SFA and SST models.

3. Knowledge-enhanced meta-classifier

We propose the use of a meta-learning scheme for combining 174 different classical approaches, i.e., bag of words, *n*-grams or lexical 175 resource-based classifiers. Key to our approach is adding also other 176 knowledge-based classifiers. By using a semantic network, we per-177 form word sense disambiguation and generate new independent 178 classifiers for the main part-of-speech tags: disambiguated adjec-179 tives, nouns, verbs and adverbs. Using the disambiguated terms, 180 the semantic network allows us to obtain a vocabulary 181 expansion-based classifier. In Section 3.1 we present the semantic 182 network, and the word sense disambiguation and vocabulary 183 expansion methods. Then, in Section 3.2 we describe the base 184 classifiers that compose our system. Finally, in Section 3.3 185 we define the Stacked Generalization that we use to combine those 186 classifiers. 187

3.1. Word sense disambiguation and vocabulary expansion via a semantic network

A semantic network [53] is a (un)directed graph consisting of vertices, which represent concepts, and edges, which represent semantic relations between them. Concepts are usually organized into a taxonomic hierarchy. Fig. 1 shows a simple example of semantic network.

In this work we use the semantic network graph to: (i) perform word sense disambiguation, and (ii) perform a vocabulary expansion using the disambiguated words. Despite having the WordNet Semantic Network [21], which is an historical resource including 117,000 synsets² in English, in this work we are interested in employing a larger size wide-coverage lexical knowledge resource. Among those, we can find knowledge bases extracted automatically from Wikipedia such as DBPedia [5] or YAGO [27]. However, due to its WordNet-based internal structure combined with Wikipedia, the high amount of synsets included, and the lexicalizations of its concepts available in multiple languages,³ we chose the BabelNet Multilingual Semantic Network.

Please cite this article in press as: M. Franco-Salvador et al., Cross-domain polarity classification using a knowledge-enhanced meta-classifier, Knowl. Based Syst. (2015), http://dx.doi.org/10.1016/j.knosys.2015.05.020

2

89

90

91

92

93

172

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

173

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

² Set of word synonyms.

³ While this work is exclusively evaluated on English, this multilinguality allows us to perform at multilingual level.

Download English Version:

https://daneshyari.com/en/article/6862400

Download Persian Version:

https://daneshyari.com/article/6862400

Daneshyari.com