# Hierarchical emotion classification and emotion component analysis on chinese micro-blog posts

Hua Xu [a,1,*], Weiwei Yang [a,1], Jiushuo Wang [a,b,1]

[a] *State Key Laboratory of Intelligent Technology and Systems, Tsinghua National Laboratory for Information Science and Technology, Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China*
[b] *School of Information Science and Engineering, Hebei University of Science and Technology, Hebei 050018, China*

## ARTICLE INFO

## ABSTRACT

Text emotion analysis has long been a hot topic. With the development of social network, text emotion analysis on micro-blog posts becomes a new trend in recent years. However, most researchers classify posts into coarse-grained emotion classes, which cannot depict the emotions accurately. Besides, flat classification is mostly adopted, which brings difficulty for classifiers when given a large dataset. In this paper, by data pre-processing, feature extraction and feature selection, we classify Chinese micro-blog posts into fine-grained emotion classes, employing hierarchical classification to improve the performance of classifiers. Moreover, based on the regression values in classification procedure, we propose an algorithm to detect the principal emotions in posts and calculate their ratios.

## 1. Introduction

For years, researchers are trying to classify the emotions in text automatically. The source of text differs a lot, such as newspapers, magazines, stories, blog posts, etc. However, with the explosive development of social network services (SNS), more and more people begin to express their views and attitudes on the Internet. The views and attitudes, of course, often contain emotions.

Among all SNS platforms, the most popular one is micro-blog. Micro-blog posts directly reflect users' opinions. However, the posts are usually very short. Sina Weibo, for instance, the most widely used micro-blog platform in China, allows users to write no more than 140 Chinese characters in a single post. The length of posts brings challenges to emotion classification and requires more effective methods to extract features. Besides, Internet slang is not easy to cope with either because it does not follow language rules.

Emotion, definitionly, is a subjective thought or feeling like *happy, angry*, etc, while sentiment addresses the objective positive and negative attitudes. It is possible that a post contains sentiment but no emotions. Like the sentence in Example 1, it does not contain emotions as it is just telling us a fact, however, it contains the writer's negative attitudes towards the phone.

1.  The phone broke within two days.

Currently, most researchers are focusing on sentiment analysis and emotion classification on six basic coarse-grained emotion classes (Ekman, 1971), which consist of *happy, surprise, angry, disgusted, fear* and *sad*. However, coarse-grained emotions cannot depict the emotions in text perfectly. Example 2 clearly expresses the emotion *disappointed*, but it is not proper to classify it into any class above. In order to better describe emotions, fine-grained emotions need to be added to coarse-grained emotion categories, which forms hierarchy. Besides, adopting fine-grained emotions greatly increases the number of classes, which brings difficulty for flat classifiers, so hierarchical classification is required.

2.  This car is not so easy to drive as the ad says. I am so disappointed.

In addition to emotion granularity, corpus language also varies. So far, the corpus of most work is in English. Not many papers' results are based on Chinese.

In this paper, we hierarchically classify Chinese micro-blog posts into fine-grained emotion classes, employing the four-level hierarchy proposed in Xu, Lin, Pan, Ren, and Chen (2008). Psychological emotion dictionary, Internet slang dictionary and emoticon dictionary are employed to segment posts and form the feature space, which is then selected by a combination of $\chi^2$-test, word frequency and pointwise mutual information (PMI), in order to retain effective features. Finally, we employ support vector regression (SVR) and rule sets, which are generated by PMI values, to get the classification results, which, as reported later, are very encouraging.

Now many researchers are focusing on positive / negative or coarse-grained basic emotion classification with 6–7 classes, rather

* Corresponding author. Tel.: 86-1062796450, fax: 8610-62771792.
  *E-mail address:* xuhua@tsinghua.edu.cn, 147502394@qq.com (H. Xu).
[1] Indicates equal contributions from these authors

than fine-grained emotion classification or emotion component analysis. In this paper, a four-level fine-grained emotion hierarchy with 19 basic emotions is adopted. However, posts usually contain more than one kind of emotions. For example, when a fan criticizes a football team, his emotion can be a combination of disappointment, sadness and anger. So we propose an emotion component analysis (ECA) algorithm to detect the principal emotions in posts and calculate the corresponding ratios according to the classification results, which, more specifically, according to distances between regression values and class thresholds.

The rest of paper is organized as follows. Section 2 gives a brief view of related work. Section 3 introduces the hierarchy and hierarchical emotion classification algorithm. The ECA algorithm is introduced in Section 4. Section 5 presents the experimental results before the final section concludes the paper.

## 2. Related work

There are a number of research reports and papers about emotion classification on texts. Although probability-based algorithms are quite useful (Boiy, Hens, Deschacht, & Moens, 2007; Paltoglou, Gobron, Skowron, Thelwall, & Thalmann, 2010; Strapparava & Mihalcea, 2008), machine learning approach is more preferred by researchers nowadays (Wei & Gulla, 2012; Yang & Yu, 2013). Alm, Roth, and Sproat (2005) report the classification results on 1580 sentences picked from 22 fairy tales. They extract the common features such as emotion words, part-of-speech (POS), special punctuation, as well as some corpus-specific features like story type, story progress, direct speech etc. With all the features fed to the classification algorithm, they reach an accuracy of 0.69 on six basic emotion classes (Ekman, 1971). However, the dataset is imbalance as almost 60% of the sentences are *neutral*. To overcome this problem, (Ghazi, Inkpen, & Szpakowicz, 2010) adopt a two-level hierarchy and run the algorithm again on the same dataset. It turns out that hierarchical classification brings some improvements. Huang, Peng, Li, and Lee (2013) propose a multi-task multi-label classification model that performs classification based on both emotions and topics concurrently. Support vector machine (SVM) is an efficient algorithm. Beyond that the Í-Support Vector Regression (Gu, Sheng, Wang, Ho, Osman, & Li, 2015b) and Support vector ordinal regression (Gu, Sheng, Tay, Romano, & Li, 2015a) are also effective regression learning algorithms. In Cho and Kang (2012), they classify the emotion behind tweets by extracting feature vectors and using the SVM classification method. Li, Xiao, and Xue (2012) propose an unsupervised domain independent method for sentiment classification, to build a emotion vocabulary list for text emotional classification. In He (2013), the researcher compares the performance of the three methods (i.e., Naive Bayesian, SVM, and SMO) in micro-blog emotion classification. Liu and Chen (2015) propose a multi-label classification based approach for emotion analysis, including text segmentation, feature extraction and multi-label classification. Meanwhile, SVM can apply to other fields, such as it can be trained with some features so as to identify a test image (Xia, Wang, Sun, & Wang, 2014) and so on.

In order to better classify text, researchers spend time constructing and improving emotion lexicons. Mohammad and Turney (2010, 2012) report the construction of emotion lexicons using common words and crowdsourcing, while some others adopt classifier-based approach (Das, Poria, & Bandyopadhyay, 2012). Emotion lexicons bring magnificent improvement to emotion classification on text. As Aman and Szpakowicz (2007) Mohammad (2012) report, the WordNet Affect Lexicon (Strapparava & Valitutti, 2004; Strapparava, Valitutti, & Stock, 2006) significantly benefits the sentence-level emotion classification.

In addition to classification algorithms and emotion lexicons, corpus is also an option. Some researchers try to classify emotions on blog posts. The classification results on 624,905 blog posts are reported in Keshtkar and Inkpen (2009) Mishne and Gilad (2005). Mishne and Gilad (2005) extract frequency counts, length-related features and semantic orientation features from the posts and classify into 132 emotion classes using support vector machine (SVM). The highest accuracy among all emotions is 65.75%. Keshtkar and Inkpen (2009), on the other hand, adopt a five-level hierarchy. Bag-of-Words (BoW) and sentiment oriented features are fed to SVM and an average accuracy of 55.24% is achieved, much better than that of flat classification (24.73%). The flat classification can classify the examples directly relative to the hierarchical classification. While the hierarchical classification classifies the examples from top to bottom according to the pre-determined multi-layer classification system and gets the final classification result in the bottom. The flat classification is mostly adopted, which brings difficulty for classifiers to distinguish between the examples belong to its class and other classes when given a large dataset.

Recent years, as micro-blog is used more and more widely, micro-blog posts become a new source of corpus for emotion classification. Go, Bhayani, and Huang (2009) try to classify twitter posts into positive and negative emotions. They adopt unigram, bigram and POS as features and compare the results generated by Naive Bayes, MaxEnt and SVM. Their corpus-specified preprocessing to the posts is very instructive, such as removing @username, short links and repeated letters. Experiments on Chinese micro-blog posts are also reported in Liu, Feng, and Huang (2012) Tang and Chen (2011) Yuan and Purver (2012). Tang and Chen (2011) run experiments on the corpus on Plurk while other two reports get the dataset from Sina Weibo. Li and Xu (2014) propose and implement a novel method for text-based emotion classification using emotion cause extraction in microblog posts. All of the three reports classify posts into positive/negative emotions or 6 coarse-grained emotion classes.

Besides, experiments on other kinds of corpus are also reported, e.g. e-mails (Mohammad & Yang, 2011), novels (Mohammad, 2011) and Japanese dialog systems (Tokuhisa, Inui, & Matsumoto, 2008).

Our contributions are different. We hierarchically classify Chinese micro-blog posts into 19 fine-grained emotion classes with machine learning approach and propose an ECA algorithm based on the regression values. As for the contributions of this research in the application area of social management, the government can find some existing problems by analyzing public emotions in social media. Meanwhile we can also apply our algorithm to consumer behavior analysis for making the right decision. In the process of segmentation, a psychological emotion dictionary is adopted in this paper for improving the effect of the algorithm, which has important scientific values both on social network knowledge discovery and data mining.

## 3. Emotion classification

### 3.1. Hierarchy

The four-level hierarchy is presented in Fig. 1. According to Xu et al. (2008), it is developed on the basis of seven basic emotion classes proposed in Ekman (1971). This hierarchy contains 19 fine-grained emotion classes at the bottom level and 20 leaf nodes if considering *neutral*, which denotes the non-emotional class.

### 3.2. Preprocessing

We remove four kinds of elements, which are explained as follows, as they definitely does not contain emotions.

**Usernames.** If user A wants to share something with user B, he will include B's username by adding an @ symbol before it in his post, so that B will be notified. However, this part is surely non-emotional, so we take it away by detecting @ symbol and remove it together with the username.