



ELSEVIER

Contents lists available at ScienceDirect

## Neurocomputing

journal homepage: [www.elsevier.com/locate/neucom](http://www.elsevier.com/locate/neucom)

# Semantic expansion using word embedding clustering and convolutional neural network for improving short text classification

Peng Wang<sup>a,\*</sup>, Bo Xu<sup>a</sup>, Jiaming Xu<sup>a</sup>, Guanhua Tian<sup>a</sup>, Cheng-Lin Liu<sup>a,b</sup>, Hongwei Hao<sup>a</sup>

<sup>a</sup> Institute of Automation, Chinese Academy of Sciences, Beijing 100190, PR China

<sup>b</sup> National Laboratory of Pattern Recognition (NLPR), Beijing 100190, PR China

## ARTICLE INFO

## Article history:

Received 4 May 2015

Received in revised form

22 June 2015

Accepted 30 September 2015

Communicated by Jinhui Tang

Available online 9 October 2015

## Keywords:

Short text

Classification

Clustering

Convolutional neural network

Semantic units

Word embeddings

## ABSTRACT

Text classification can help users to effectively handle and exploit useful information hidden in large-scale documents. However, the sparsity of data and the semantic sensitivity to context often hinder the classification performance of short texts. In order to overcome the weakness, we propose a unified framework to expand short texts based on word embedding clustering and convolutional neural network (CNN). Empirically, the semantically related words are usually close to each other in embedding spaces. Thus, we first discover semantic cliques via fast clustering. Then, by using additive composition over word embeddings from context with variable window width, the representations of multi-scale semantic units<sup>1</sup> in short texts are computed. In embedding spaces, the restricted nearest word embeddings (NWEs)<sup>2</sup> of the semantic units are chosen to constitute expanded matrices, where the semantic cliques are used as supervision information. Finally, for a short text, the projected matrix<sup>3</sup> and expanded matrices are combined and fed into CNN in parallel. Experimental results on two open benchmarks validate the effectiveness of the proposed method.

© 2015 Elsevier B.V. All rights reserved.

## 1. Introduction

The classification of short texts, such as search snippets, micro-blogs, product reviews, and short messages, plays important roles in user intent understanding, question answering and intelligent information retrieval [1]. Since short texts do not provide enough contextual information, the data sparsity problem is easily encountered [2]. Thus, the general methods based on bag-of-words (BoW) model cannot be directly applied to short texts [1], because the BoW model ignores the order and semantic relations between words. How to acquire effective representations of short texts to enhance the categorization performance has been an active research issue [2,3].

Conventional text classification methods often expand short texts using latent semantics, learned by latent Dirichlet allocation (LDA) [4] and its extensions. Phan et al. [3] presented a general framework to expand short and sparse texts by appending topic names, discovered using LDA over Wikipedia. Sahami and Heilman [5] enriched text representation by web search results using the

\* Corresponding author.

<sup>1</sup> Semantic units are defined as  $n$ -grams which have dominant meaning of text. With  $n$  varying, multi-scale contextual information can be exploited.

<sup>2</sup> In order to prevent outliers, a Euclidean distance threshold is preset between semantic cliques and semantic units, which is used as restricted condition.

<sup>3</sup> The projected matrix is obtained by table looking up, which encodes Unigram level features.

short text segment as a query. Furthermore, Yan et al. [6] presented a variant of LDA, dubbed biterm topic model (BTM), especially for short text modeling to alleviate the data sparsity problem. However, these methods still consider a text as BoW. Therefore, they are not effective in capturing fine-grained semantics for short texts modeling.

More recently, deep learning based methods have drawn much attentions in the field of natural language processing (NLP), which mainly evolved into two branches. One is to learn word embeddings by training language models [7–10], and another is to perform semantic composition to obtain phrase or sentence level representation [11,12]. Word embeddings, also known as distributed representations, typically represent words with dense, low-dimensional and real-valued vectors. Each dimension of the vectors encodes a different aspect of words. In embedding spaces, semantically close words are likely to cluster together and form semantic cliques. Moreover, the embedding spaces exhibit linear structure that the word embeddings can be meaningfully combined using simple vector addition [9].

In this paper, we aim to obtain the semantic representations of short texts and overcome the weakness of conventional methods. Similar to Li et al. [13] that cluster indicators learned by non-negative spectral clustering are used to provide label information for structural learning, we develop a novel method to model short texts using word embeddings clustering and convolutional neural network (CNN). For concision, we abbreviate our methods to



Download English Version:

<https://daneshyari.com/en/article/411614>

Download Persian Version:

<https://daneshyari.com/article/411614>

[Daneshyari.com](https://daneshyari.com)