

## 蒙古语词向量评测研究\*

# Research on Mongolian Word Vectors Evaluation

乌云塔那,王斯日古楞\*\*

Wuyuntana, Wangsiriguleng

(内蒙古师范大学计算机与信息工程学院,内蒙古呼和浩特 010022)

(Computer and Information Engineering College, Inner Mongolia Normal University, Hohhot, Inner Mongolia, 010022, China)

**摘要:**词向量具有良好的语义特性,可用于改善和简化许多自然语言信息处理应用。本研究利用 CBOW 和 Skip-gram 两种模型架构在不同数据和不同维度下训练蒙古语词向量,然后结合蒙古语特征设计一个语义语法综合测试集,并在此测试集上用语义和语法相似度来评测词向量质量。研究结果表明,蒙古语语义和语法相似性任务上, Skip-gram 模型优于 CBOW 模型, Skip-gram 模型的窗口大小为 5 的情况下,词向量质量最好,且随着词向量维度或训练数据的增大,词向量质量有明显的提高。

**关键词:**词向量 CBOW 模型 Skip-gram 模型 词向量质量 语义语法相似度

**中图分类号:**TP391.1 **文献标识码:**A **文章编号:**1002-7378(2018)01-0068-04

**Abstract:** The words vector has good semantic properties and can be used to improve and simplify many natural language processing applications. This study used CBOW (continuous Bag of words) and Skip-gram two model architectures to train the Mongolian word vectors in different data and different dimensions. Then we design a comprehensive semantic syntactic test set based on the Mongolia language features. And on this test set, we use semantic and syntactic similarity to estimate the quality of the word vectors. The results indicate that Skip-gram model is superior to CBOW model in Mongolian semantic and syntactic similarity tasks, and the word vectors quality is the best when the window size is 5. Moreover, with the increase of the word vectors dimension or training data, the quality of the word vectors is obviously improved.

**Key words:** word vectors, CBOW model, Skip-gram model, quality of the word vectors, semantic syntactic similarity

收稿日期:2017-11-01

修回日期:2017-12-15

作者简介:乌云塔那(1992-),女,硕士研究生,主要从事蒙古文信息处理研究, E-mail:724831848@qq.com。

\* 内蒙古自治区自然科学基金项目“基于条件随机场的蒙古文命名体识别研究”(2016MS0623)和国家自然科学基金项目“基于神经网络的蒙汉机器翻译研究”(61762072)资助。

\*\* 通信作者:王斯日古楞(1970-),女,博士,教授,主要从事蒙古文信息处理研究, E-mail:siriguleng@imnu.edu.cn。

## 0 引言

要将自然语言交给机器学习中的算法来处理,通常需要将语言数学化,词向量就是将语言中的词进行数学化的一种方式。词向量将某种语言中的每一个词映射成一个固定长度的短向量,将所有这些向量放在一起形成一个词向量空间,而每一词向量则为该空间中的一个点,在这个空间上引入“距离”,则可以根据词之间的距离来判断它们之间的语义、语法上的相似性。词向量具有良好的语义特性,可

用于改善和简化许多自然语言信息处理应用,并且词向量的质量影响自然语言信息处理应用的性能。因此,针对蒙古语词性标注、命名实体识别、短语识别、机器翻译方面的应用需求训练蒙古语词向量,研究词向量的评测具有重要的研究意义。词语作为连续向量的表示具有悠久的历史。已有许多学者用不同模型训练了词向量,比较经典的模型有神经网络语言模型(NNLM)<sup>[1]</sup>、双对数线性语言模型(LBL)<sup>[2]</sup>、循环神经网络语言模型(RNNLM)<sup>[3]</sup>、连续词袋模型(CBOW)和 Skip-gram 模型<sup>[4]</sup>等。另外,Mikolov 等<sup>[5]</sup>还提出了 Skip-gram 模型的几个扩展,即 Hierarchical Softmax 算法、负采样算法和欠采样技术,从而提高了词向量的质量和训练速度。针对形态丰富的语言,Bojanowski 等<sup>[6]</sup>提出了一种基于 Skip-gram 模型的新方法。此外,词向量可用于改善和简化许多 NLP 应用<sup>[7-8]</sup>。

词向量的评价方法有两种:一个是把词向量融入现有系统中,看能否提升现有系统<sup>[9]</sup>;另一个是从语言学的角度分析词向量,比如相似度。研究者发现相似的词不仅彼此接近,而且这个词具有多重相似度。这在早期的变形语言中已被观察到,例如,英语名词可以有多个单词结尾,如果在原始向量空间的子空间中搜索类似的单词,可以找到具有相似结尾的单词<sup>[10-11]</sup>。蒙古语是黏着性语言,其形态丰富,并存在许多词根相同的名词和动词的变形形式,这些形式表示相似的概念。比如动词“*хөдлөх, хөдлөж, хөдлөн, хөдлөх*”都有相同的词根“*хөд*”,表示“走”的不同形态。因此本研究从语言学角度的语义语法相似性来评价蒙古文词向量质量。蒙古语词向量评价方面不像英语一样有自己的语义语法测试集,因此,本研究结合蒙古语特征设计一个语义语法综合测试集,然后在此测试集上评测蒙古语词向量的质量。

## 1 蒙古语语义语法测试集的建立

词向量具有良好的语义特性,可以通过加减法操作来对应某种语义语法关系,并通过语义语法相似性来评价词向量。比如判断与  $C(\text{ᠠᠨᠠᠨᠠᠨ}) - C(\text{ᠠᠨᠠᠨᠠᠨ}) + C(\text{ᠠᠨᠠᠨᠠᠨ})$  接近的向量是不是  $C(\text{ᠠᠨᠠᠨᠠᠨ})$ 。

### 1.1 蒙古语语义测试集的建立

在大量数据上训练高维词向量时,所得到的向量可以回答诸如城市和它所属的国家之间的单词之间的微妙语义关系,例如,巴黎是法国的,柏林是德

国的。根据这种现象本研究建立了两种蒙古语语义关系集:首都-国家关系和男-女关系。每个类别的两个实例如表 1 所示,相同类别的两个单词对连在一起构成一个语义问题,共有 100 个语义问题。本测试集中只包含一个词构成的单词,不包含多单词实体(如 *ᠲᠤᠨ ᠤᠯᠤᠰ*)。

表 1 语义测试集中的两种语义关系实例

Table 1 Examples of two types of semantic questions in Semantic test set

Type of relationship	Word pair 1	Word pair 2
Capital city-Country	<i>ᠶ᠋ᠢᠵᠢᠨᠠᠨ</i> <i>ᠶ᠋ᠢᠵᠢᠨᠠᠨ</i>	<i>ᠶ᠋ᠤᠯᠤᠰ</i> <i>ᠶ᠋ᠤᠯᠤᠰ</i>
Man-Woman	<i>ᠠᠮ</i> <i>ᠠᠮ</i>	<i>ᠮᠤ</i> <i>ᠮᠤ</i>

### 1.2 蒙古语语法测试集的建立

蒙古语名词有格、数、领属等范畴的形态变化。蒙古语的格是通过名词后面缀接格附加成分来表示,例如“*ᠠᠨᠠᠨᠠᠨ*”(房子的),其中“*ᠠᠨᠠᠨ*”(房子)是名词,“*ᠠᠨ*”是格附加成分,共有 7 种格附加成分。蒙古语的复数是在名词单数形式上缀接复数附加成分来表示。例如名词单数形式“*ᠮᠠᠭᠢ*”(老师)上缀接复数附加成分“*ᠮᠠᠭᠢ*”,表示“老师们”的意思。蒙古语数词的变化也是在基数词上缀接各种附加成分表示,比如基数词“*ᠡᠨᠢ*”(一)上缀接“*ᠡᠨᠢ*”附加成分形成“*ᠡᠨᠢ*”,表示“第一”的意思。

本研究结合蒙古语特征建立了关于蒙古语名词格、复数、数词和代词的 4 种语法关系集。每个类别的两个实例如表 2 所示,相同类别的两个单词对连在一起构成一个语法问题,共有 544 个语法问题。

表 2 语法关系测试集中 4 种类型语法问题实例

Table 2 Examples of four types of syntactic questions in Syntactic test set

Type of relationship	Word pair 1	Word pair 2
Numerals form change	<i>ᠡᠨᠢ</i> <i>ᠡᠨᠢ</i>	<i>ᠡᠨᠢ</i> <i>ᠡᠨᠢ</i>
Plural nouns	<i>ᠮᠠᠭᠢ</i> <i>ᠮᠠᠭᠢ</i>	<i>ᠮᠠᠭᠢ</i> <i>ᠮᠠᠭᠢ</i>
Noun case	<i>ᠠᠨᠠᠨᠠᠨ</i> <i>ᠠᠨᠠᠨᠠᠨ</i>	<i>ᠠᠨᠠᠨᠠᠨ</i> <i>ᠠᠨᠠᠨᠠᠨ</i>
Pronouns form change	<i>ᠠᠨᠠᠨᠠᠨ</i> <i>ᠠᠨᠠᠨᠠᠨ</i>	<i>ᠠᠨᠠᠨᠠᠨ</i> <i>ᠠᠨᠠᠨᠠᠨ</i>

### 1.3 词向量评测方法

使用类比方式在建立的语义语法测试集上评测蒙古语词向量质量。具体操作方法:设语义语法测试集中每个问题的 4 个词依次对应 a、b、c、d。已知 a 之于 b 犹如 c 之于 d。先给出 a、b、c,再看  $C(a) - C(b) + C(c)$  最接近的词是否是  $C(d)$ 。如果计算出来的向量与测试集中的词 d 完全相同,则认为是正



于 CBOW 模型(表 6)。

表 6 模型架构在语 X 语法测试集上的准确性比较

Table 6 The accuracy comparison of model architecture in semantic-syntactic test set

Model architecture	Semantic accuracy (%)	Syntactic accuracy (%)
CBOW	5.12	32.85
Skip-gram	15.12	42.95

综上所述,蒙古语语义准确率总体比较低,第一个原因是测试中忽略了同义词的概念,比如,就测试集中的“ᠠᠨᠠᠨᠠ ᠠᠨᠠᠨᠠ ᠠᠨᠠᠨᠠ ᠠᠨᠠᠨᠠ”来说,“ᠠᠨᠠᠨᠠ”有同义词“ᠠᠨᠠᠨᠠ”,而本研究所采用的测评方法规定,只有计算出来的向量与测试集中的词“ᠠᠨᠠᠨᠠ”完全相同时才认为是正确答案,所以同义词“ᠠᠨᠠᠨᠠ”被忽略。第二个原因是训练语料库中缺少地名相关的词,因此训练出来的地名词向量质量比较差,无法表达出其语义。

#### 4 结束语

本研究分别使用 CBOW 模型和 Skip-gram 模型训练蒙古语词向量,并在自己建立的语义和语法测试集上评测了词向量的质量。研究表明,利用 Skip-gram 模型且窗口为 5 的情况下蒙古语词向量质量最好。随着词向量维度或训练数据的增大,词向量质量有明显的提高。

蒙古语是个形态丰富的语言,有着丰富的数、格、时、体、态等形态变化。这导致了蒙古语的词汇量庞大。词向量训练时使用固定大小的词汇表,这使得罕见词语无法向量化。因此,在后续工作中将研究基于子词单元(比如词素级、字符级)的词向量表示来提升形态丰富语言的性能。词向量训练好以后,通常会作为各种神经网络结构的初始值,Word2vec 模型是很浅层的神经网络,词向量经预训练后做为其初始值,通常可以提升任务上的效果。因此,后续研究将会把训练好的词向量作为初始值,运用到蒙汉机器翻译任务上,以提升其翻译效果。

#### 参考文献:

[1] BENGIO Y, DUCHARME R, VINCENT P, et al. A

neural probabilistic language model[J]. Journal of Machine Learning Research, 2003, 3: 1137-1155.

[2] MNH A, HINTON G. Three new graphical models for statistical language modelling[C]. 24th Annual International Conference on Machine Learning (ICML), Corvallis, 2007.

[3] MIKOLOV T. Statistical language models based on neural networks [D]. Lausanne: Brno University of Technology, 2012.

[4] MIKOLOV T, YIH W, ZWEIG G. Linguistic regularities in continuous space word representations[C]. Proceedings of NAACL-HLT, 2013.

[5] MIKOLOV T, SUTSKEVER I, CHEN K. Distributed representations of words and phrases and their compositionality [C]//BURGES C J C, BOTTOU L, WELLING M, et al (eds.). Advances in neural information processing systems 26 (NIPS 2013). Nevada: [S. n.], 2013.

[6] BOJANOWSKI P, GRAVE E, JOULIN A, et al. Enriching word vectors with subword information[Z]. Facebook AI Research, 2016.

[7] COLLOBERT R, WESTON J, BOTTOU L. Natural language processing (almost) from scratch[J]. Journal of Machine Learning Research, 2011, 12: 2493-2537.

[8] KIM Y. Convolutional neural networks for sentence classification[C]. Empirical Methods in Natural Language Processing. [S. l. ], 2014.

[9] TURIAN J, RATINOV L, BENGIO Y. Word representations: A simple and general method for semi-supervised learning[C]. Proc Association for Computational Linguistics. [S. l. ], 2010.

[10] MIKOLOV T. Language modeling for speech recognition in Czech [D]. Lausanne: Brno University of Technology, 2007.

[11] MIKOLOV T, KOPECKY, BURGET L, et al. Neural network based language models for highly inflective languages[C]. International Conference on Acoustics, Speech and Signal Processing. [S. l. ], 2009.

(责任编辑:陆雁)