

Research on the Selection Method and Data Set Construction of Yi Language Stop Words

Junping Huang

Fangchenggang City, Guangxi Zhuang Autonomous Region, China

Abstract: This paper aimed to complete the research on the selection method and data set construction of the Yi language stop words which has positive significance for information retrieval, text classification, sentiment analysis and online public opinion analysis. On the basis of word segmentation, data cleaning and word frequency statistics of Yi language text, used statistical methods such as Document Frequency, Term Frequency, and Entropy, analyzed and selected the top 100 high-frequency words that were ranked by these three measures, through the intersection of three sets of high-frequency words, 89 words were identified as potential stop words in the Yi language, and according to the characteristics of stop words, 20 function words were selected as stop words. Based on the methods of selecting stop words in Chinese, Mongolian, Tibetan, Uyghur, etc., completed the construction of the Yi language stop words dataset.

Keywords: Yi Language, Stop Words, Selection Method, Data Set, Construction.

1. Introduction

Stop words refer to a type of word that appears frequently in a specific text and whose semantic information has little impact on Natural Language Processing. Stop words are pivotal in the Natural Language Processing of the Yi language, and can be used in tasks such as information retrieval, text classification, sentiment analysis, and online public opinion analysis. In Natural Language Processing, after word segmentation, eliminate stop words from the text before further processing will significantly reduce the text's dimensionality, enhance the computational efficiency of the computer, and consequently improve the speed, quality, and accuracy of text processing. This study explores the method of selecting stop words in Yi language from a certain amount of Yi language texts that have completed word segmentation, data cleaning, and word frequency statistics, and constructs a dataset of stop words in Yi language.

2. Related Work

Abroad, the study of stop words initially arose from research in information retrieval. Luhn [1], [2] found in his research on information retrieval that some words appear frequently, but these words seriously affect the retrieval performance, he proposed using noise to represent such words. Wilbur et al. [3] identified stop words in a collection using automated statistical testing. Ho [4] proposed a stop words localization and recognition strategy in adaptive text recognition. Khabibull et al. [5] used DF and TF statistical algorithms to construct a stop words dataset for Karakalpak language.

In the study of Chinese stop words, Hao et al. [6] proposed to use Chi-square test to complete text classification and automatic extraction of stop words lists. Zou et al. [7] proposed a stop words selection method based on statistical and information theory models. Jiang Zhaozhong [8] studied Chinese automatic word segmentation based on context and stop words drive, and pointed out that Chinese stop words refers to words with high frequency and no great retrieval significance. Xiong Wenxin et al. [9] studied the stop word

filtering of information retrieval users' query statements. Zhou Qinqiang et al. [10] used the program flow control to eliminate the single independent words, English characters, numbers, a series of mathematical symbols and the Chinese words containing these symbols, so that the pure Chinese words with more than two characters became the feature items representing text information. Yang et al. [11] conducted a comparative study on feature selection methods in the statistical learning of text classification, and carried out dimensionality reduction processing on text.

Research on stop words in Chinese minority languages, Gong Zheng et al. [12], [13] obtained a list of Mongolian stop words through statistical methods, in which there were often some entity nouns that were closely related to the subject but affected the retrieval accuracy, and used the Joint Entropy algorithm to preliminarily determine the Mongolian stop words, after removed the entity nouns and homonyms from the initial determination of Mongolian stop words, the list of Mongolian stop words was determined by comparing the parts of speech of English stop words and Mongolian stop words. Zhu Jie et al. [14] analyzed the selection of stop words in Tibetan through methods such as Term Frequency, Document Frequency, Entropy, etc., and proposed a Tibetan stop words selection method that combined Tibetan function words, special verbs, and automatic processing methods. SAIMAITI et al. [15] analyzed the characteristics of Uyghur stop words, used the methods of Document Frequency, Term Frequency and Entropy to conduct statistics on a large number of corpus, and analyzed the part of speech distribution of candidate stop words.

3. Stop Words Selection Method

Used statistical methods such as Document Frequency (DF), Term Frequency (TF), and Entropy (EN) to automatically select stop words in Yi language. All three methods apply statistical methods to assess the importance of a word to a text set, and the importance of a word increases in proportion to the number of times it appears in the text set.

3.1. Document Frequency Method

The Document Frequency method counts how many texts

