

Understand Short Texts by Harvesting and Analyzing Semantic Knowledge by considering the crucial ambiguous base

¹Rajshri Dattu, ²Mr.S.M.Shinde

¹M.E.Student ,SVERI'sCollege of Engineering ,Pandharpur,
² Assistant Professor ,SVERI'sCollege of Engineering ,Pandharpur,

Abstract-Short messages don't generally watch the linguistic structure of a written language. So not applicable to traditional natural language processing tools, ranging from part-of-speech tagging to dependency parsing. Short texts usually do not contain sufficient statistical signals or signs to support for text mining such as topic modeling. Short texts are more ambiguous and noisy, and are generated in an enormous volume, which further increases the trouble to deal with them. Here focus on semantic information in order to better understand short texts. In this work, build a model framework for short text understanding which abuses semantic knowledge gave by a well-known knowledge base and automatically harvested from a web corpus. Our knowledge-intensive approaches disrupt traditional methods for tasks such as text segmentation, part-of-speech tagging, and concept labeling, in the sense that focus on semantics in all these tasks.

Index Terms - Short messages, segmentation, natural language processing.

I. INTRODUCTION

Understanding short texts is tough to many applications, but challenges remain in social network for messaging and chatting. Seeing short messages is hard to numerous applications. There are several challenges are as follows;

1. Short messages don't generally watch the linguistic structure of a written language.
2. Not applicable to traditional natural language processing tools, ranging from part-of-speech tagging to dependency parsing.
3. Short texts usually do not contain sufficient statistical signals or signs to support for text mining such as topic modeling.
4. Short texts are more ambiguous and noisy and are generated in an enormous volume, which further increases the trouble to deal with them.

Here work focus on Semantic Knowledge to better understand short texts. In this work, build a model framework for short text understanding which abuses semantic knowledge gave by a well-known knowledge base and automatically harvested from a web corpus. Our knowledge-intensive approaches disrupt traditional methods for tasks such as text segmentation, part-of-speech tagging, and concept labeling, in the sense that focuses on semantics in all these tasks.

Short texts refer to texts with limited context. Many applications, small-scale blogging services, and web search, etc., are required to deal with many short texts. A better understanding of short texts will deliver tremendous value. One of the essential tasks of text understanding is to recognize hidden semantics from texts. Lots of attempts have been committed to this field. For occurrence, named entity recognition locates named entities in text and classifies them into predefined topic models attempt to recognize "latent topics," which are represented as probabilistic distributions on words, from a text. Entity linking focuses on retrieving "explicit topics" expressed as probabilistic distributions on an entire knowledge base. However, categories, "latent topics," as well as "explicit topics" still have a semantic gap with humans' mental world. As stated in Psychologist Gregory Murphy's highly acclaimed book, "concepts are the glue that keeps our mental world collectively." Short text requirement is easy to understand and easy to implement.

II. RELEVANCE

In this work, it is argued that semantic knowledge is indispensable for short text understanding, which in turn benefits many real-world applications that need to handle a large number of short texts. According to the above discussion, three types of knowledge are required to cope with the challenges of short text understanding:

- 1) A comprehensive vocabulary;
- 2) Mappings between instances and concepts;

Semantic coherence between terms.

Will describing how can harvest this knowledge is based on the acquired knowledge, introduce knowledge-intensive approaches to understand short texts both effectively and efficiently.

III. OBJECTIVES

The different objectives of the proposed system are

- Develop a text segmentation module to find the most semantically coherent segmentation.
- Develop module to cluster similar concepts in the dictionary together.
- Designing Ambiguity removal in short texts and overcome the limitations of traditional approaches in handling them.
- Design technique to achieve better accuracy of short text understanding by harvesting semantic knowledge.
- Design approaches to facilitate online instant short text understanding.

IV. LITERATUREREVIEW

A. McCallum and W. Li, "Early results for named entity recognition with conditional random fields, feature induction, and web-enhanced lexicons,"

This paper describes WebListing, a method that obtains seeds for the lexicons from the labeled data, then uses the Web, HTML formatting regularities and a search engine service to significantly augment those lexicons. For example, based on the appearance of Arnold Palmer in the labeled data, we gather from the Web an extensive list of other golf players, including Tiger Woods (a phrase that is difficult to detect as a name without a good lexicon).

G. Zhou and J. Su, "Named entity recognition using an hmm-based chunk tagger,"

This paper proposes a Hidden Markov Model (HMM) and an HMM-based chunk tagger, from which a named entity (NE) recognition (NER) system is built to recognize and classify names, times and numerical quantities.

D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation."

This paper describes latent Dirichlet allocation (LDA), a generative probabilistic paradigm for groups of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a gathering is registered as a measurable mixture over an underlying set of topics. Each subject is, in turn, modeled as an infinite mixture over an underlying set of topic probabilities.

M. Rosen-Zvi, T. Gri_ths, M. Steyvers, and P. Smyth, "The author-topic model for authors and documents."

In this introduce the author-topic model, a generative model for documents that extends Latent Dirichlet Allocation to include signature information. Each author is affiliated with a multinomial distribution over topics, and each topic is connected with a multinomial pattern over words. A document with multiple contributors is modeled as a distribution over issues that is a mixture of the distributions associated with the authors.

R. Mihalcea and A. Csomai, "Wikify! Linking documents to encyclopedic knowledge,"

This paper proposes the use of Wikipedia as a resource for automatic keyword extraction and word sense disambiguation and records how this online encyclopedia can be used to achieve state-of-the-art results on both these jobs. The paper also explains how the two methods can be mixed into a system able to automatically improve a text with links to encyclopedic knowledge.

V. PROPOSED WORK

MODELS: OFFLINE PROCESSING

A prerequisite to short text understanding is the knowledge about semantic relatedness between terms. In this offline Processing model design and construct the co-occurrence network and quantify semantic coherence. Apply this score for the indexing strategy to allow for approximate term extraction on the vocabulary, as well as the approach to determine instance ambiguity.

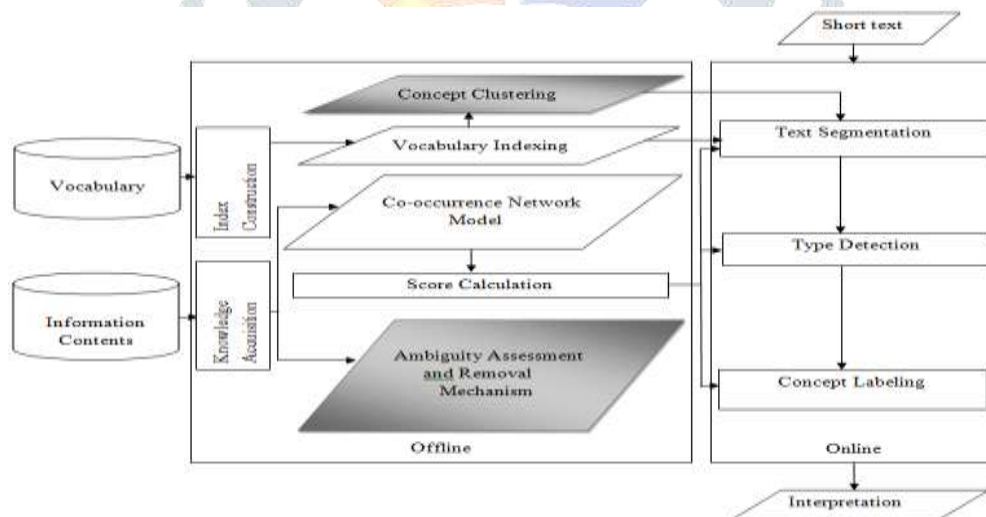


Figure1. Block diagram of the proposed system

Constructing Co-occurrence Network

A co-occurrence network used to model semantic relatedness. The co-occurrence network regarded as an undirected graph, where nodes are typed-terms and edge weight formulates the strength of semantic relatedness between typed terms.

It is observed that Terms of different types occur in different contexts. The more frequently two typed-terms co-occur in a sentence, the higher the semantic relatedness will be. The closer two typed-terms appear in a sentence, the higher the semantic relatedness will be. Common terms (e.g., "item" and "object") which co-occur with almost every other term are meaningless in modeling semantic relatedness; thus the corresponding edge weights should be penalized.

Scoring Semantic Coherence

In this module calculate a score to measure semantic coherence between typed-terms. Two types of coherence used as similarity and relatedness (co-occurrence). considering two typed-terms are coherent if they are semantically similar or they often co-occur on the web. Indexing Vocabulary for Approximate Term Extraction Approximate term extraction aims to locate substrings in a text which are similar to terms contained in a predefined vocabulary. To quantify the relationship between two strings, many similarity functions have been suggested

meantime. We present a Chain Model and a Pairwise Model which consolidate lexical furthermore, semantic highlights to direct good location. They accomplish superior precision over conventional POS taggers on the named benchmark.

VIII. REFERENCES

- [1] Wen Hua, Zhongyuan Wang, Haixun Wang, Kai Zheng, and Xiaofang Zhou, Senior "Understand Short Texts by Harvesting and Analyzing Semantic Knowledge" IEEE Trans. On Knowledge and data Engg., vol. 29, no. 3, Mar 2017
- [2] A. McCallum and W. Li, "Early results for named entity recognition with conditional random fields, feature induction, and web-enhanced lexicons," in Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003 - Volume 4, ser. CONLL '03, Stroudsburg, PA, USA, 2003, pp. 188–191.
- [3] G. Zhou and J. Su, "Named entity recognition using an hmm-based chunk tagger," in Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ser. ACL '02, Stroudsburg, PA, USA, 2002, pp. 473–480.
- [4] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," J. Mach. Learn. Res., vol. 3, pp. 993–1022, 2003.
- [5] M. Rosen-Zvi, T. Gri_ths, M. Steyvers, and P. Smyth, "The author-topic model for authors and documents," in Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence, ser. UAI '04, Arlington, Virginia, United States, 2004, pp. 487–494.
- [6] R. Mihalcea and A. Csomai, "Wikify! Linking documents to encyclopedic knowledge," in Proceedings of the sixteenth ACM conference on Conference on information and knowledge management, ser. CIKM '07, New York, NY, USA, 2007, pp. 233–242.
- [7] D. Milne and I. H. Witten, "Learning to link with Wikipedia," in Proceedings of the 17th ACM conference on Information and knowledge management, ser. CIKM '08, New York, NY, USA, 2008, pp. 509–518.
- [8] S. Kulkarni, A. Singh, G. Ramakrishnan, and S. Chakrabarti, "Collective annotation of Wikipedia entities in web text," in Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, ser. KDD '09, New York, NY, USA, 2009, pp. 457–466.
- [9] X. Han and J. Zhao, "Named entity disambiguation by leveraging Wikipedia semantic knowledge," in Proceedings of the 18th ACM conference on Information and knowledge management, ser. CIKM '09, New York, NY, USA, 2009, pp. 215–224.
- [10] "Structural semantic relatedness: A knowledge-based method to named entity disambiguation," in Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, ser. ACL '10, Stroudsburg, PA, USA, 2010, pp. 50–59.

