

Segmentation of existing methods for text segmentation



segmentation as to divide a text into a sequence of terms. Statistical approaches, such as N-gram Model 21, 22, 23, calculate the frequencies of words co-occurring as neighbors in a training corpus. When the frequency exceeds a predefined threshold, the corresponding neighboring words can be treated as a term. Vocabulary-based approaches 18, 19, 20 extract terms by checking for their existence or frequency in a predefined vocabulary. The most obvious drawback of existing methods for text segmentation is that they only consider surface features and ignore the requirement of semantic coherence within a segmentation. This might lead to incorrect segmentations as described in Challenge 1. To this end, we propose to exploit context semantics when conducting text segmentation. POS tagging.

POS tagging determines lexical types (i. e., POS tags) of words in a text. Rule-based POS taggers attempt to assign POS tags to unknown or ambiguous words based on a large number of hand-crafted 10, 11 or automatically learned 12, 13 linguistic rules. Statistical POS taggers avoid the cost of constructing tagging rules by building a statistical model automatically from a corpora and labeling untagged texts based on those learned statistical information. Mainstream statistical POS taggers employ the well-known Markov Model 14, 15, 16, 17 which learns both lexical probabilities and sequential probabilities from a labeled corpora and tags a new sentence by searching for tag sequence that maximizes the combination of lexical and sequential probabilities. Note that both rule-based and statistical POS taggers rely on the assumption that texts are correctly structured which, however, is not always the case for short texts. More importantly, existing methods only considers lexical features and ignores word semantics.

This might lead to mistakes, as illustrated in Challenge 3. Our work attempts to build a tagger which considers both lexical features and underlying semantics for type detection. Semantic labeling. Semantic labeling discovers hidden semantics from a natural language text. Named entity recognition (NER) locates named entities in a text and classifies them into predefined categories (e.

g., persons, organizations, locations, etc.) using linguistic grammar-based techniques as well as statistical models like CRF¹ and HMM². Topic models³ attempt to recognize “latent topics”, which are represented as probabilistic distributions on words, based on observable statistical relations between texts and words.

Entity linking^{5, 6, 7, 8} employs existing knowledge bases and focuses on retrieving “explicit topics” expressed as probabilistic distributions on the entire knowledge base. Despite the high accuracy achieved by existing work on semantic labeling, there are still some limitations. First, categories, “latent topics”, and “explicit topics” are different from human-understandable concepts.

Second, short texts do not always observe the syntax of a written language which, however, is an indispensable feature for mainstream NER tools. Third, short texts do not contain sufficient content to support statistical models like topic models. The work most related to ours are conducted by Song et al.¹⁹ and Kim et al.²⁰ respectively, which also represent semantics as concepts. ¹⁹ employs the Bayesian Inference mechanism to conceptualize instances

and short texts, and eliminates instance ambiguity based on homogeneous instances.

Kim et al. (2010) captures semantic relatedness between instances using a probabilistic topic model (i. e., LDA), and disambiguates instances based on related instances.

In this work, we observe that other terms, such as verbs, adjectives, and attributes, can also help with instance disambiguation. We incorporate type discernment into our framework for short text understanding of conduct instance disambiguation based on various types of context information.