

# Tagging: descriptors to the given tokens is

Life, Friendship



Tagging: The descriptors are called the tags and the automatic assignment of the descriptors to the given tokens is called tagging. POS Tagging The process of assigning one of the parts of speech to the given word is called Parts Of Speech tagging, commonly referred to as POS tagging. Parts of speech include nouns, verbs, adverbs, adjectives, pronouns, conjunction and their sub-categories. POS Tagger A Part-Of-Speech Tagger (POS Tagger) is a software that reads text and then assigns parts of speech to each word (and other token), such as noun, verb, adjective, etc., It uses different kinds of information such as dictionary, lexicons, rules, etc. because dictionaries have category or categories of a particular word, that is a word may belong to more than one category. For example, run is both noun and verb so to solve this ambiguity taggers use probabilistic information. There are mainly two types of taggers: Rule-based- Uses hand-written rules to distinguish the tag ambiguity. Stochastic taggers are either HMM based - chooses the tag sequence which maximizes the product of word likelihood and tag sequence probability, or cue-based, using decision trees or maximum entropy models to combine probabilistic features.

HMM Hidden Markov Model (HMM) is a statistical Markov model in which the system being modeled is assumed to be a Markov process with unobserved (i. e. hidden) states.

In Markov models, the state is directly visible, and thus the state transition probabilities are the only parameters, while in the hidden Markov model, the state is not directly visible, but the output dependent on the state is visible. Each state has a probability distribution over the possible output tokens.

Therefore, the sequence of tokens generated by an HMM gives some

<https://assignbuster.com/tagging-descriptors-to-the-given-tokens-is/>

information about the sequence of states. Accuracy achieved The European group developed CLAWS, a tagging program that did exactly this, and achieved accuracy in the 93–95% range. Many machine learning methods have also been applied to the problem of POS tagging. Methods such as SVM, maximum entropy classifier, perceptron, and nearest-neighbor have all been tried, and most can achieve accuracy above 95%. A more recent development is using the structure regularization method for part-of-speech tagging, achieving 97.36% on the standard benchmark dataset.

Tagset A set of tags from which the tagger chooses a relevant tag for the word. Data set A merged Bhojpuri dataset containing sentences of Bhojpuri and the corresponding labels to the words. Natural Language Processing (NLP) with Python NLTK is a leading platform for building Python programs to work with human language data.

It provides easy-to-use interfaces to over 50 corpora and lexical resources such as WordNet, along with a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, wrappers for industrial-strength NLP libraries, and an active discussion forum. It has many libraries to work on natural language. Using we can tokenize and tag some text, identify some named entities and display a sparse tree. ACKNOWLEDGEMENT I express my profound and sincere gratitude to my mentor Dr. Anil Kumar Singh for providing me with all the facilities and support during my winter internship period. I would like to thank my guide Mr. Rajesh Mundotiya for their valuable guidance, constructive criticism, encouragement and also for making the

requisite guidelines enabling me to complete my work with utmost dedication and efficiency. At last, I would like to acknowledge my family and friends for their motivation, inspiration and support in boosting my moral without which my efforts would have been in vain.