

Natural language Processing

Saumya Bhatnagar

March 5, 2020

Table of contents I

Components of NLP

Topic Modeling/Document clustering

Sentiment Analysis

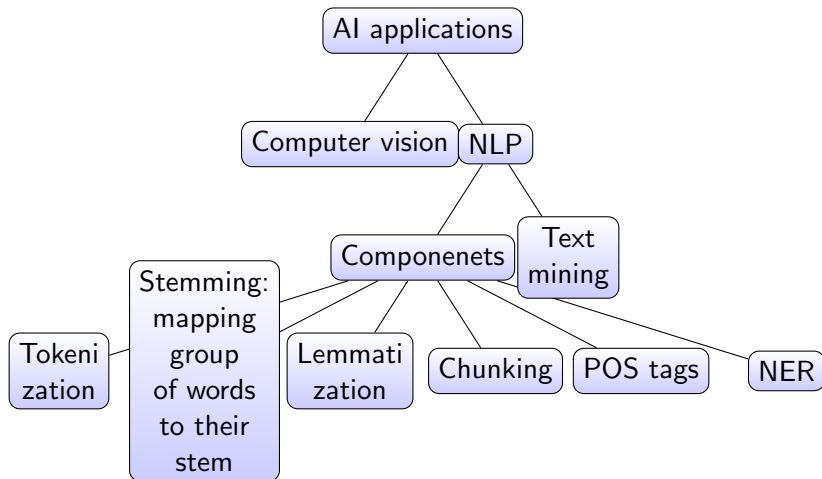
Information Retrieval

General

Text Mining: process of deriving high quality information from texts

Turn text into data for analysis via NLP

NLP (AI): deals with human languages
what is inflection?



Tokenizer

break sentences into words

Stemming and Lemmatization

Stemming

Might not be an actual word

Predefine steps

Speed is imp

Stemming: Stem word might not be an actual word

- ▶ Porter stemmer: based on 5 pre-defined rules, simple and fast, mainly used in IR (information retrieval) search queries
- ▶ Snow-ball stemmers: by NLTK, has non-english stemmers (french, german, english, italian,)
- ▶ Lancaster stemmer: Iterative, over-stemming may occur, shorter stem then Porter
- ▶ ISRI Stemmer
- ▶ RSLPS Stemmer

Lemmatization

Returns an actual word

Uses WordNet corpus

Language is imp

Word is imp

bayesian Modeling

content...

content...

content...

Thank You!