

NLP Basics

- * Tokenization
- * Lemmatization
- * Stemming

1. Feature Engineering

2. Language Processing

Tokenization

→ String breaking

A big dog. } Tokenization
word based based on
space.

tokenization } Sub word
based tokenization

A b i g d o g : } letter
based
tokenization

Stemming (Idea is to have ^{least #} independent features)

organize, organizes, organizing | democracy democratic

~~~~~  
how to deal with similar features?

Lemmatization [same idea  $\Rightarrow$  to have min # features]

lexeme (am, are, is)  $\Rightarrow$  (be) lemma  
car, cars, car's, cars'  $\Rightarrow$  (car) lemma

## Stemming

reducing words to their root.

1. Porter Stemmer
  2. Snowball Stemmer
- } part of NLTK

generate  $\rightarrow$  generat  
generation  $\rightarrow$  generat

stemming  
 $\downarrow$   
but,  
generat  
doesn't use context

but in lemmatization  $\Rightarrow$  they both  
becomes "generate"  
probably.