# TOKENIZATION

# TEXT REPRESENTATION AND VECTORIZATION

Convert text to numbers

Computers can do only ONE thing, that is, COUNTING!

**SYRACUSE UNIVERSITY**
School of Information Studies

# TOKENIZATION

A tokenizer has a set of rules about grouping characters into tokens.

## Word Tokenization with Python NLTK

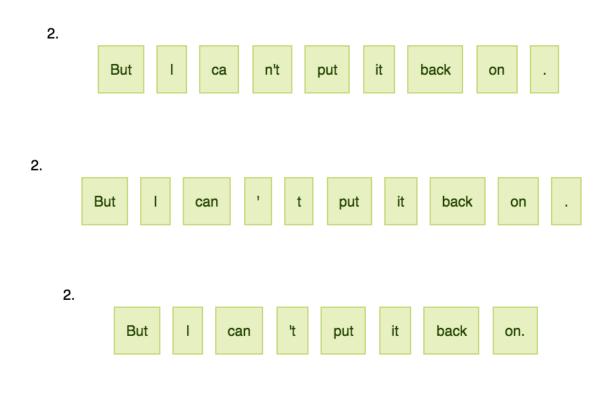This is a demonstration of the various **tokenizers** provided by NLTK 2.0.4.

**Tokenize Text**

**Enter text**

In Düsseldorf I took my hat off. But I can't put it back on.

Enter up to 50000 characters

Tokenize

### TreebankWordTokenizer

1.

| In | Düsseldorf | I | took | my | hat | off | . |

2.

| But | I | ca | n't | put | it | back | on | . |

SYRACUSE UNIVERSITY
School of Information Studies

# TOKENIZATION RULES

2.

| But | I | ca | n't | put | it | back | on | . |

2.

| But | I | can | ' | t | put | it | back | on | . |

2.

| But | I | can | 't | put | it | back | on. |

2.

| But | I | can't | put | it | back | on. |

# TOKENIZATION IS NOT EASY

Tokenizing URLs

Choosespain.com

# TOKENIZATION IS NOT EASY

Tokenize text strings with no white space

Chinese (New Year couplets):
养猪大如山老鼠头头死

Raise|pigs|big|as|mountain|rats|all|die
养|猪|大|如|山|老鼠|头头|死

Raise|pigs|big|as|mountain rats, all|die
养|猪|大|如|山老鼠| 头头|死

SYRACUSE UNIVERSITY
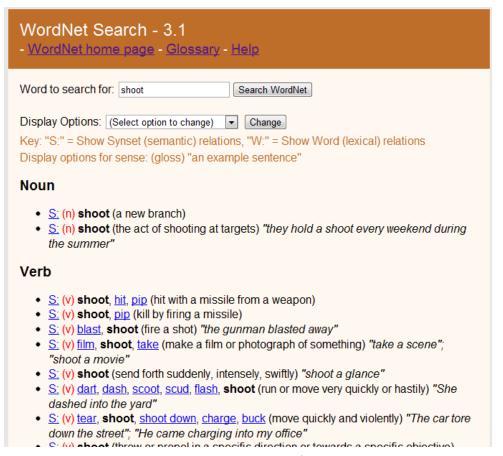School of Information Studies

# TOKENIZATION IS NOT EASY

Lowercase vs. uppercase

Words with inflected forms
 "dishwasher" vs. "dishwashers"

Words with multiple senses
 "There is a money **bank** near the river **bank."**

**SYRACUSE UNIVERSITY**
School of Information Studies

# WORDNET

http://wordnetweb.princeton.edu/perl/webwn

**SYRACUSE UNIVERSITY**
School of Information Studies

# WORD SENSE DISAMBIGUATION (WSD)

WSD techniques use word context to decide the word sense

Could introduce more errors to next steps

So far does not help search engines significantly

Not widely used in text mining

Text mining tends to use shallow features to process large amount of text data.