# Web Data Mining
## Lecture 4: Text Mining

**Jaroslav Kuchař & Milan Dojčinovski**

jaroslav.kuchar@fit.cvut.cz, milan.dojchinovski@fit.cvut.cz

Czech Technical University in Prague - Faculty of Information Technologies - Software and Web Engineering
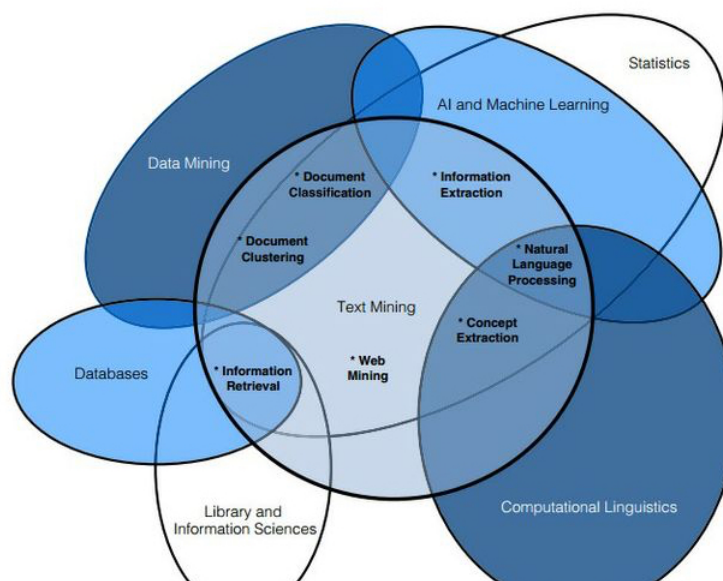
## Overview

- Introduction
- Text Processing
- Named Entity Recognition
- Relation Extraction

# Introduction

- ## What is Text Mining
  - *Extracting information and knowledge from text*
  - *Information and knowledge previously unknown to the user*
- ## Unknown information
  - *Information not known even for the writer*
  - *Rediscover (extract) information encoded in the text*
- ## Text Mining Process
  - *assembling large corpora of documents*
  - *performing preprocessing of documents*
  - *text transformation and feature generation*
  - *dimensions reduction - feature selection*
  - *pattern discovery (data mining)*
  - *results interpretation*

# Text Mining Overview

- S. M. Weiss, N. Indurkhya, T. Zhang, and F. Damerau, Text mining: predictive methods for analyzing unstructured information. Springer Science and Business Media, 2010.

## Main Challenges

- Text Mining is not easy!
- Main issues:
  - *High dimensionality*
  - *Different terms of same concepts*
    *→ e.g. car, vehicle, auto, automobile, ...*
  - *Ambiguity - same term can identify one or more concepts*
    *→ e.g. player (sportsman, musician, performer, reproducing device, ...)*
  - *Lack of structure*
    *→ The data (text) is not structured*
      *→ Compared to: spreadsheets, database tables, etc.*
      *→ Rows, columns, headings identify the meaning (semantics) of the content*
- But - The data is highly redundant!

## Text Mining Applications

- General applications
  - *Information retrieval (recommendation systems)*
  - *Information extraction*
- Social and Business
  - *Customer profile analysis*
  - *Social media data analysis*
  - *Trend analysis & Event tracking*
- Data Mining
  - *Document classification*
  - *Document clustering*
- Security and Crime
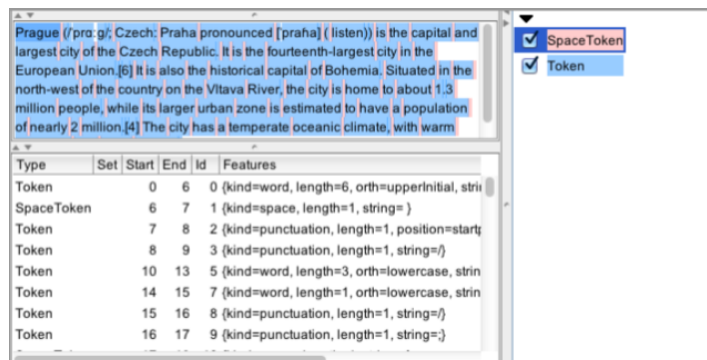  - *Spam filtering*
  - *Fraud detection*

## Overview

- Introduction
- Text Processing
- Named Entity Recognition
- Relation Extraction

## Preprocessing

- Text needs to be preprocessed to be machine understandable
- Converting a raw text file into a well-defined sequence of linguistically-meaningful units
  - *Identifying words (tokens)*
  - *Special symbols - dots, commas, whitespaces*
  - *Identifying sentences*
- Needed for further text processing stages
  - *Part-Of-Speech tagging*
  - *Noun phrase chunking*
  - *Morphological analyzers*
  - *Named entity recognition*

## Tokenization

- Grouping sequences of characters into logical elements called tokens (symbols, words, numbers)
  - *tokens are usually separated by whitespace characters or punctuation*

- Tokenizer
  - *system component performing the process of tokenization*

## Tokenization (cont.)

- Trivial for a person familiar with the language structure

- More complicated for a computer program
  - *characters are sometimes token delimiters and sometimes not*

- Main groups
  - white spaces
    - → *space, tabs, new lines*
    - → *acting as delimiters and not as tokens*
  - ()<>!?"
    - → *delimiters and sometimes tokens*
  - .,:-'
    - → *delimiters only in specific situations*

# Tokenization (cont.)

- Main situations
  - **.,:**
    - → *not delimiter in numbers*
    - → *e.g.* 100,000.5 , 12:45:37, ...
  - **.**
    - → *part of abbreviations and other constructs*
    - → *e.g.* P.S., D.I.Y, Dr., ...
  - **'**
    - → *part of the token*
      - → *e.g.* isn't
    - → *terminator*
      - → *e.g.* Tess'
    - → *quoting*
  - **-**
    - → *part of a phone number, accounts, ...*
    - → *e.g.* 123-456-789
- Generally language-dependent

# Tokenization Example

- NLTK code:

```
1   import nltk
2   text = """Mr. Speaker, Mr. Vice President, Members of Congress,
3   the First Lady of the United States, and Citizens of America:
4   ...
5   In 9 years, the United States will celebrate the 250th annivers
6   of our founding -- 250 years since the day we declared
7   our Independence.
8   ...
9   In Chicago, more than 4,000 people were ...
10  """
11  tokens = nltk.word_tokenize(text)
12  print(tokens)
```

- Output:

```
1   ['Mr.', 'Speaker', ',', 'Mr.', 'Vice', 'President', ',', 'Members',
2    'of', 'Congress', ',', 'the', 'First', 'Lady', 'of', 'the', 'United',
3    'States', ',', 'and', 'Citizens', 'of', 'America', ':', '...', 'In',
4    '9', 'years', ',', 'the', 'United', 'States', 'will', 'celebrate',
5    'the', '250th', 'anniversary', 'of', 'our', 'founding', '--',
6    '250', 'years', 'since', 'the', 'day', 'we', 'declared', 'our', 'Independen
7    '4,000', 'people', 'were', '...']
```

# Frequency Analysis

- Basic analysis of textual data
- Term frequency
  - *compute frequency distributions as ranked lists of terms*
- Lexical diversity
  - *diversity of an individual's or group's vocabulary*

```python
import nltk
from collections import Counter
text = None
with open('speech.txt', 'r') as f:
    text = f.read()

tokens = nltk.word_tokenize(text)
c = Counter(tokens)
print(c.most_common()[:10])

print(1.0*len(set(tokens))/len(tokens))
```

```
[(',', 261), ('.', 253), ('the', 215), ('and', 176), ('to', 142), ('of', 142
0.28324761204996324
```

# Stemming and Lemmatization

- The goal
  - *Reduce the dimensionality!*
  - *Convert each of the tokens to a standard form.*
  - *Decreases the number of tokens and increases its frequency.*
- Lemmatization
  - *Finding lemma (canonical form) for a given word*
  - *words:* are, is, was, were*; lemma:* be
- Stemming
  - *Finding stem for a given word*
  - *Stem: root form of a word to which suffixes can be attached*
  - *word:* waiting*; stem:* wait

## Stemming

- Part rule-based and part dictionary-based
  - *If token in dictionary return from dictionary*
  - *It token ends with* s *strip* s
  - *If token ends with* ies *replace by* y
  - *...*

```
1  word = "printing"
2  regexp = r'^(.*?)(ing|ly|ed|ious|ies|ive|es|s|ment)?$'
3  stem, suffix = re.findall(regexp, word)[0]
4  stem # print
```

```
1  import nltk
2  from nltk.stem.porter import PorterStemmer
3  text = """Mr. Speaker, Mr. Vice President, Members of Congress,
4  the First Lady of the United States, and Citizens of America:
5  """
6  stemmer = PorterStemmer()
7  tokens = nltk.word_tokenize(text)
8  stems = {token:stemmer.stem(token) for token in tokens}
9  print(stems)
```

```
1  {'the': 'the', 'Citizens': 'Citizen', 'President': 'Presid',
2  'Congress': 'Congress', 'of': 'of', 'Members': 'Member', 'and': 'and',
3  'Mr.': 'Mr.', 'First': 'First', ',': ',', 'United': 'Unit', 'Lady': 'Ladi',
4  'America': 'America', ':': ':', 'Speaker': 'Speaker', 'Vice': 'Vice',
5  'States': 'State'}
```

## Lemmatization

```
1   import nltk
2   from nltk.corpus.reader.wordnet import NOUN,VERB
3   from nltk.stem import WordNetLemmatizer
4   text = """Mr. Speaker, Mr. Vice President, Members of Congress,
5   the First Lady of the United States, and Citizens of America:
6   ...
7   In 9 years, the United States will celebrate the 250th anniversary of our fo
8   -- 250 years since the day we declared our Independence.
9   ...
10  In Chicago, more than 4,000 people were ...
11  """
12  lemmatizer = WordNetLemmatizer()
13  tokens = nltk.word_tokenize(text)
14  lemmas = {token:lemmatizer.lemmatize(token, pos=VERB) for token in tokens}
15  print(lemmas)
```

```
1   {'celebrate': 'celebrate', ':': ':', 'Congress': 'Congress', 'anniversary':
2   'of': 'of', 'we': 'we', '250': '250', 'Independence': 'Independence',
3   'First': 'First', ',': ',', 'since': 'since', 'United': 'United',
4   'Chicago': 'Chicago', 'America': 'America', 'Speaker': 'Speaker', 'In': 'In'
5   'day': 'day', 'people': 'people', '250th': '250th', 'the': 'the', '9': '9',
6   'President': 'President', 'our': 'our', 'Members': 'Members', 'and': 'and',
7   'declared': 'declare', 'Mr.': 'Mr.', '.': '.', 'Citizens': 'Citizens',
8   'will': 'will', 'were': 'be', 'Lady': 'Lady', '--': '--', '4,000': '4,000',
9   '...': '...', 'years': 'years', 'Vice': 'Vice', 'States': 'States',
10  'more': 'more', 'founding': 'found', 'than': 'than'}
```

# Other Dimensionality Reductions

- ## Stop words
  - *almost never provide any interesting information*
    - → *articles* a *and* the
    - → *pronouns such as* it *and* they

```
1  from nltk.corpus import stopwords
2  stops = stopwords.words('english')
3  text = """In Chicago, more than 4,000 people were ..."""
4  tokens = nltk.word_tokenize(text)
5  filtered_tokens = [token for token in tokens if token not in stops]
6  filtered_tokens
```

```
1  ['In', 'Chicago', ',', '4,000', 'people', '...']
```

- ## Frequency information
  - *most frequent words are often stopwords and can be deleted*
  - *rare words are often typos and can also be dismissed*
  - *tf-idf*

# Other Dimensionality Reductions (cont.)

- ## Synonyms vs Antonyms

```
1  from nltk.corpus import wordnet
2  token = "lady"
3  synonyms = []
4  antonyms = []
5  syns = wordnet.synsets(token)
6  for syn in syns:
7      print("%s - %s " % (syn.lemmas()[0].name(), syn.definition()))
8      for l in syn.lemmas():
9          synonyms.append(l.name())
10         if l.antonyms():
11             antonyms.append(l.antonyms()[0].name())
12 print(set(synonyms))
13 print(set(antonyms))
```
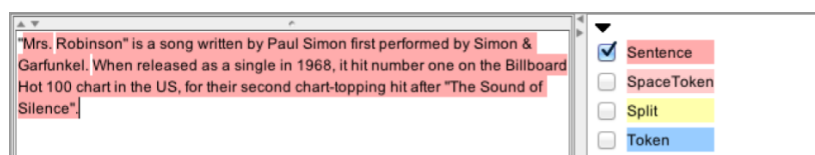
```
1  lady - a polite name for any woman
2  dame - a woman of refinement
3  Lady - a woman of the peerage in Britain
4  {'peeress', "ma'am", 'madam', 'dame', 'gentlewoman', 'Lady', 'lady',
5  {'Lord', 'nobleman'}
```

# Frequency Analysis 2

```
1  import nltk
2  from nltk.corpus import stopwords
3  from string import punctuation
4  stops = stopwords.words('english')
5  from collections import Counter
6
7  text = None
8  with open('speech.txt', 'r') as f:
9      text = f.read()
10
11  tokens = nltk.word_tokenize(text)
12  filtered_tokens = [token for token in tokens if token not in stops]
13  nopunc_tokens = [token for token in filtered_tokens if token not in punctuat
14
15  c = Counter(nopunc_tokens)
16  print(c.most_common()[:10])
17  print(1.0*len(set(nopunc_tokens))/len(nopunc_tokens))
```

```
1  [('--', 67), ('I', 38), ('We', 35), ('America', 30), ('American', 30),
2  ('must', 20), ("'s", 20), ('new', 19), ('us', 18), ('country', 18)]
3  0.5049261083743842
```

# Sentence Splitting

- Process of identifying sentences in text
- Not a trivial process
  - *Dot can have also other uses then sentence splitting symbol, for example, Mrs. Smith is …*
- Heuristics can be used to improve the sentence splitting
- Maintain a list of patterns containing dot or write regex
  - *Mrs. | Mr. | Dr. | Ph.D. | Ing. | Bc.*

## Sentence Splitting (cont.)

```
1   import nltk
2   text = """Mr. Speaker, Mr. Vice President, Members of Congress,
3   the First Lady of the United States, and Citizens of America:
4   ...
5   In 9 years, the United States will celebrate the 250th anniversary of
6   -- 250 years since the day we declared our Independence.
7   ...
8   In Chicago, more than 4,000 people were ...
9   """
10  sentences = nltk.sent_tokenize(text)
11  print(sentences)
```

```
1   [
2   'Mr. Speaker, Mr. Vice President, Members of Congress, \nthe First La
3   'In 9 years, the United States will celebrate the 250th anniversary o
4   '...',
5   'In Chicago, more than 4,000 people were ...\n'
6   ]
```

## N-Grams

- Set of co-occurring words within a given window
  – *Bigrams for N=2*
  – *Trigrams for N=3*

- Number of N-grams
  – #NG = X-(N-1) *where X is a number of tokens in the text.*

```
1   import nltk
2   text = """
3   In 9 years, the United States will celebrate the 250th anniversary of
4   -- 250 years since the day we declared our Independence.
5   """
6   from nltk.util import ngrams
7   for ng in ngrams(nltk.word_tokenize(text),3):
8       print(ng)
```

```
1   ('In', '9', 'years')
2   ('9', 'years', ',')
3   ('years', ',', 'the')
4   (',', 'the', 'United')
5   ('the', 'United', 'States')
6   ...
```

# N-Grams application

- Collocations
  - *Expressions of multiple words which commonly co-occur*

```
1  import nltk
2  from nltk.collocations import BigramCollocationFinder
3  from nltk.metrics import BigramAssocMeasures
4  text = None
5  with open('speech.txt', 'r') as f:
6      text = f.read()
7  bigram_finder = BigramCollocationFinder.from_words(
8      nltk.word_tokenize(text))
9  bigram_finder.apply_freq_filter(2)
10 bigrams = bigram_finder.nbest(BigramAssocMeasures.chi_sq, 10)
11 bigrams
```

```
1  [('Homeland', 'Security'), ('Middle', 'East'), ('United', 'States'),
2  ('middle', 'class'), ('43', 'million'), ('Rare', 'Disease'),
3   ('Republican', 'President'), ('civil', 'rights'),
4   ('illegal', 'immigrant'), ('inner', 'cities')]
```

# N-Grams application (cont.)

- Generating text using N-Grams

```
1  import nltk
2  import random
3  text = None
4  with open('speech.txt', 'r') as f:
5      text = f.read()
6  ng = ngrams(nltk.word_tokenize(text),2)
7  fd = nltk.ConditionalFreqDist(ng)
8
9  word = "the"
10 output = []
11 for i in range(15):
12     output.append(word)
13     word = random.choice(list(fd[word].keys()))
14 print(" ".join(output))
15 """
```

```
1  the courage to make childcare accessible and reaffirmed our uniform.
2  the Bible teaches us harm . So I 've saved
3  the day we focus on becoming lobbyists for drugs from
4  ...
```

# Part-of-Speech Tagging

- Linguistic analysis
- Assigning linguistic categories to words in a text
  - *Number of categories is not fixed (from 6 or 7 to tens)*
  - *Difficulties:*
    - → bore *could be a noun, a present tense verb, or a past tense verb*
  - *can also take into consideration the context the words appear*

| Title | Abbreviation | Example |
|---|---|---|
| Noun | NN | computer |
| Proper Noun | PNP | Africa |
| Adjective | JJ | nice |
| Pronoun | PRP | he |
| Verb | VB | work |
| Adverb | RB | word "fast" in "He runs fast." |
| Proposition | IN | in |
| Conjunction | CC | and |

# Part-of-Speech Tagging

```
1  nltk.help.upenn_tagset()
```

```
1   ...
2   CC: conjunction, coordinating
3       & 'n and both but either et for less minus neither nor or plus so
4       therefore times v. versus vs. whether yet
5   ...
6   IN: preposition or conjunction, subordinating
7       astride among uppon whether out inside pro despite on by throughout
8       below within for towards near behind atop around if like until below
9       next into if beside ...
10  JJ: adjective or numeral, ordinal
11      third ill-mannered pre-war regrettable oiled calamitous first separable
12      ectoplasmic battery-powered participatory fourth still-to-be-named
13      multilingual multi-disciplinary ...
14  ...
15  NN: noun, common, singular or mass
16      common-carrier cabbage knuckle-duster Casino afghan shed thermostat
17      investment slide humour falloff slick wind hyena override subhumanity
18      machinist ...
19  ...
20  PRP: pronoun, personal
21      hers herself him himself hisself it itself me myself one oneself ours
22      ourselves ownself self she thee theirs them themselves they thou thy us
23  RB: adverb
24      occasionally unabatingly maddeningly adventurously professedly
25      stirringly prominently technologically magisterially predominately
26      swiftly fiscally pitilessly ...
27  UH: interjection
28      Goodbye Goody Gosh Wow Jeepers Jee-sus Hubba Hey Kee-reist Oops amen
29      huh howdy uh dammit whammo shucks heck anyways whodunnit honey golly
30      man baby diddle hush sonuvabitch ...
31  VB: verb, base form
32      ask assemble assess assign assume atone attention avoid bake balkanize
```

# Part-Of-Speech Tagging Example

- Performing POS using the GATE framework

# Simple POS Tagger

- Using large already tagged corpora and perform a mapping/classification task (with default to 'NN')
  - *e.g. ('the', 'DT'), ('day', 'NN'), ('celebrate', 'VB'), ...*

- Regular expressions

```
1   patterns = [
2       (r'.*ing$', 'VBG'),                 # gerunds
3       (r'.*ed$', 'VBD'),                  # simple past
4       (r'.*es$', 'VBZ'),                  # 3rd singular present
5       (r'.*ould$', 'MD'),                 # modals
6       (r'.*\'s$', 'NN$'),                 # possessive nouns
7       (r'.*s$', 'NNS'),                   # plural nouns
8       (r'^-?[0-9]+(.[0-9]+)?$', 'CD'),    # cardinal numbers
9       (r'.*', 'NN')                       # nouns (default)
10  ]
```

- Rules based approaches
  - *e.g. VB if the tag of the preceding word is 'TO'*

- N-grams

- …

- → Combinations of all existing approaches

## Part-Of-Speech Tagging Example

```python
1   import nltk
2   text = """
3   In 9 years, the United States will celebrate the 250th annivers
4   -- 250 years since the day we declared our Independence.
5   """
6   print(nltk.pos_tag(nltk.word_tokenize(text)))
```

```
1   [('In', 'IN'), ('9', 'CD'), ('years', 'NNS'), (',', ','),
2   ('the', 'DT'), ('United', 'NNP'), ('States', 'NNPS'),
3   ('will', 'MD'), ('celebrate', 'VB'), ('the', 'DT'),
4   ('250th', 'JJ'), ('anniversary', 'NN'), ('of', 'IN'),
5   ('our', 'PRP$'), ('founding', 'NN'), ('--', ':'),
6   ('250', 'CD'), ('years', 'NNS'), ('since', 'IN'),
7   ('the', 'DT'), ('day', 'NN'), ('we', 'PRP'),
8   ('declared', 'VBD'), ('our', 'PRP$'),
9   ('Independence', 'NN'), ('.', '.')]
```

## Phrase Detection

- Grouping tokens into units
  - *Also known as chunking*
  - *Utilizes grammar categories identified within the POS tagging task*
  - *Can be a prerequisite for Named Entity Recognition tasks*
- Using pattern-based classifications
- Noun Phrase Chunking (NP-chunking)
  - *Based on a chunk grammar*
    - → *set of rules that indicate how sentences should be chunked*
    - → *e.g. determiner followed by a number/adjective and terminated by a noun*

## Phrase Detection (cont.)

```
1   import nltk
2   text = """
3   the first lady of the united states.
4   """
5   text_pos = nltk.pos_tag(nltk.word_tokenize(text))
6   grammar = "NP: {<DT>?<JJ>*<NN|NNS>}"
7   cp = nltk.RegexpParser(grammar)
8   result = cp.parse(text_pos)
9   print(result)
10  result.draw()
```

```
1   (S
2     (NP the/DT first/JJ lady/NN)
3     of/IN
4     (NP the/DT united/JJ states/NNS)
5     ./.)
```

## Overview

- Introduction
- Text Processing
- Named Entity Recognition
- Relation Extraction

# Named Entity Recognition

- Entity identification in free texts, and their classification
  - *turning verbose text data into a more compact structural form*
  - *special kind on phrase detection*
    - → *proper noun phrases*

- NER sub-tasks
  - *recognition*
    - → *spotting text fragments with entity mentions*
  - *classification*
    - → *assigning class to an entity mention*
  - *disambiguation via linking*
    - → *assigning URI (e.g., Wikipedia URI) describing the entity*
    - → *also known as Entity Linking/Disambiguation task*

- NER systems are usually trained on a large corpus
  - *using methods such as Conditional Random Fields*

- NER based on previously defined rules (grammars) or gazetteers
  - *Gazetteers - list of countries, persons, geo locations*

# Named Entity Recognition

- Results utilization
  - *dimension reduction*
  - *content enrichment*

- Various NER tools/APIs available
  - *Entityclassifier.eu*
  - *AlchemyAPI*
  - *OpenCalais*
  - *Wikimeta*
  - *DBpedia Spotlight*
  - *NERD*
  - *and many others*

# Entity Recognition using Entityclassifier.eu

**Extraction, Disambiguation and Classification of Entities and Named Entities**

**Input text**

The Charles Bridge is a famous historic bridge that crosses the Vltava river in Prague, Czech Republic.

**Settings**

Request timeout (in seconds): 60

Language of the input text
☑ English ☐ German ☐ Dutch

Provenance of types
☑ THD ☑ DBpedia ☑ Yago

Knowledge base (THD)
☑ Linked Hypernyms Dataset
☐ Local Wikipedia mirror
☐ Live Wikipedia

Types of entities to extract
☑ Named Entities ☐ Common Entities ☐ Both

**Run!**

**Detailed results for entity: Charles Bridge** ×

THD types

1. Bridge for entity disambiguated as Charles Bridge ACC: 0.85 +- 2.5%
2. route of transportation for entity disambiguated as Charles Bridge ACC: >= 0.85 +- 2.5%
3. infrastructure for entity disambiguated as Charles Bridge ACC: >= 0.85 +- 2.5%

DBpedia types

1. Place for entity disambiguated as Charles Bridge
2. ArchitecturalStructure for entity disambiguated as Charles Bridge

YAGO types

1. e 102898711 for entity disambiguated as Charles Bridge
2. Bridges completed in 1402 for entity disambiguated as Charles Bridge

**Results**

The Charles Bridge is a famous historic bridge that crosses the Vltava river in Prague, Czech Republic.

*Results processed in 0.407 seconds.*

---

# Entities Recognition Task

- First and main NER task
- Find entity mentions in text
- Detect start and end offset for each entity mention in a text
- Example:
  - *"The* Charles Bridge *is a famous historic bridge that crosses the* Vltava *river in* Prague, Czech Republic.*"*
- Mentions:
  - *substring* "Charles Bridge"*, start: 4, end: 18*
  - *substring* "Vltava"*, start: 64, end: 70*
  - *substring* "Prague"*, start: 80, end: 86*
  - *substring* "Czech Republic"*, start: 88, end: 102*

# Entity Recognition using NLTK

```python
import nltk
text = """Mr. Speaker, Mr. Vice President, Members of Congress,
the First Lady of the United States, and Citizens of America:
...
In 9 years, the United States will celebrate the 250th anniversary of our fo
-- 250 years since the day we declared our Independence.
...
In Chicago, more than 4,000 people were ...
"""
tokens = nltk.word_tokenize(text)
tagged = nltk.pos_tag(tokens)

ne_chunked = nltk.ne_chunk(tagged, binary=True)
def extractEntities(ne_chunked):
    data = {}
    for entity in ne_chunked:
        if isinstance(entity, nltk.tree.Tree):
            text = " ".join([word for word, tag in entity.leaves()])
            ent = entity.label()
            data[text] = ent
        else:
            continue
    return data
extractEntities(ne_chunked)
```

```
{'America': 'NE', 'Chicago': 'NE', 'Citizens': 'NE', 'Congress': 'NE',
 'Members': 'NE', 'Mr. Speaker': 'NE', 'Mr. Vice': 'NE', 'United States': 'NE
```

# Custom NER Implementation

```python
import nltk
text = """Mr. Speaker, Mr. Vice President, Members of Congress,
the First Lady of the United States, and Citizens of America:
...
In 9 years, the United States will celebrate the 250th anniversary of our founding
-- 250 years since the day we declared our Independence.
...
In Chicago, more than 4,000 people were ...
"""
tokens = nltk.word_tokenize(text)
tagged = nltk.pos_tag(tokens)

entity = []
for tagged_entry in tagged:
    if(tagged_entry[1].startswith("NN") or (entity and tagged_entry[1].startswith("IN"))):
        entity.append(tagged_entry)
    else:
        if(entity) and entity[-1][1].startswith("IN"):
            entity.pop()
        if(entity and " ".join(e[0] for e in entity)[0].isupper()):
            print(" ".join(e[0] for e in entity))
        entity = []
```

```
Mr. Speaker
Mr. Vice President
Members of Congress
First Lady
United States
Citizens of America
United States
Independence
Chicago
```

## Entites Classification Task

- Task:
  - *for each mention assign a class (type)*
- The set of classes can be pre-defined (fixed) or be dynamically created
- Example:
  - *"The Charles Bridge is a famous historic bridge that crosses the Vltava river in Prague, Czech Republic."*
- Mentions:
  - *Mention: "Charles Bridge", class:* LOC
  - *Mention: "Vltava", class:* LOC
  - *Mention: "Prague", class:* LOC
  - *Mention: "Czech Republic", class:* GPE

## Classification Set of Classes

- Possible fixed set of classes:
  - *LOC (location), GPE (Geo-political entity), ORG (organization), PER (person), MISC (miscellaneous entity, anything else)*
- The set of classes can be also defined using an ontology
  - *Advantages: each class is identified with an unique URI*
    - → *Example: DBpedia Ontology*
      - → *see DBpedia Ontology 3.9*
    - → *covers 529 classes*
    - → *used to describe articles in Wikipedia (DBpedia)*
    - → *e.g., http://dbpedia.org/resource/Capital_city*

## Entity Classification using NLTK

```python
import nltk
text = """Mr. Speaker, Mr. Vice President, Members of Congress,
the First Lady of the United States, and Citizens of America:
...
In 9 years, the United States will celebrate the 250th anniversary of our fo
-- 250 years since the day we declared our Independence.
...
In Chicago, more than 4,000 people were ...
"""
tokens = nltk.word_tokenize(text)
tagged = nltk.pos_tag(tokens)

ne_chunked = nltk.ne_chunk(tagged, binary=False)
def extractEntities(ne_chunked):
    data = {}
    for entity in ne_chunked:
        if isinstance(entity, nltk.tree.Tree):
            text = " ".join([word for word, tag in entity.leaves()])
            ent = entity.label()
            data[text] = ent
        else:
            continue
    return data
extractEntities(ne_chunked)
```

```
{'America': 'GPE', 'Chicago': 'GPE', 'Citizens': 'ORGANIZATION', 'Congress':
'First Lady': 'ORGANIZATION', 'Members': 'ORGANIZATION', 'Mr. Speaker': 'PER
```

## Entities Disambiguation Task

- Human language is not exact
  - *Same text can refer to totally different entities*
- Example: the entity "Maradona"
  - *can refer to the football player* "Diego Maradona"
  - *or, to the football coach and former player* "Hugo Maradona"
  - *or, to the movie about the football player* "Diego Maradona"
- Classification does not help with the ambiguity
  - *there can exist two diferent entities of same type*
  - "Diego Maradona" *and* "Hugo Maradona", *type:* PER
- Solution: use unique URIs to solve the ambiguity
  - *perform Entity Linking*

## Entity Linking Task

- Task:
  - *identifying entities using URIs (unique)*
- We can use knowledge base URIs to uniquely identify each entity
  - *e.g., Wikipedia or DBpedia or YAGO knowledge base URIs*
- Solution for the "Maradona" ambiguity:
  - *the **football player** "Diego Maradona"*
    $\rightarrow$ *URI: http://dbpedia.org/resource/Diego_Maradona*
  - *the **football coach and former player** "Hugo Maradona"*
    $\rightarrow$ *URI: http://dbpedia.org/resource/Hugo_Maradona*
  - *the **movie** about the football player "Diego Maradona"*
    $\rightarrow$ *URI:*
    *http://dbpedia.org/resource/Maradona_by_Kusturica*

## Coreference Resolution

- In text, different text fragments can refer to same entity
- Main objective
  - *mentioned subjects, pronouns and other referring expressions must be connected to the right individuals*
  - *match proper names and their variants in a document*
- Example:
  - Mary Smith *and* Mrs. Smith
    $\rightarrow$ *should be matched as same person*
  - International Business Machines Ltd. *and* IBM
    $\rightarrow$ *should be matched as same company*
- State of the art algorithms have accuracy of around 75%

## Overview

- Introduction
- Text Processing
- Named Entity Recognition
- Relation Extraction

## Relation Extraction

- Extracting semantic relation between entities
- Examples
  - *PERSON works for ORGANIZATION*
  - *PERSON attends EVENT*
  - *PERSON lives in LOCATION*
- Using domain-specific patterns
- Naive approach:
  - *fixed patterns*
    - → *X(Subject) works for Y(Object)*
  - *does not scale*
- Approaches based on linguistic analysis
  - *based on linguistic features (e.g. POS tagging)*

# Pattern based Relation Extraction

```
1   import nltk
2   import re
3   from nltk.sem import extract_rels,rtuple
4
5   text = None
6   with open('speech.txt', 'r') as f:
7       text = f.read()
8
9   sentences = nltk.sent_tokenize(text)
10  tokenized_sentences = [nltk.word_tokenize(sentence) for sentence in sentences]
11  tagged_sentences = [nltk.pos_tag(sentence) for sentence in tokenized_sentences]
12
13  OF = re.compile(r'.*\bof\b.*')
14
15  for i, sent in enumerate(tagged_sentences):
16      sent = nltk.ne_chunk(sent)
17      rels = extract_rels('PERSON', 'GPE', sent, corpus='ace', pattern=OF, window=5)
18      for rel in rels:
19          print('{0:<5}{1}'.format(i, rtuple(rel)))


1   1    [PER: 'Jewish/NNP Community/NNP Centers/NNPS']
2                                    'and/CC vandalism/NN of/IN' [GPE: 'Jewish/JJ']
3   123  [PER: 'Matt/NNP Bevin/NNP'] 'of/IN' [GPE: 'Kentucky/NNP']
4   198  [PER: 'Carryn/NNP Owens/NNP'] ',/, the/DT widow/NN of/IN a/DT' [GPE: 'U.S./NNP']
```

# Pattern based Relation Extraction 2

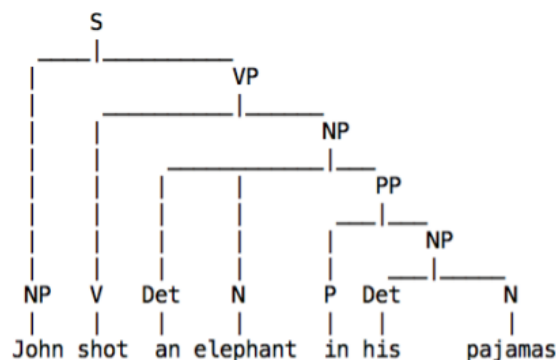- https://en.wikipedia.org/wiki/Donald_Trump

```
1   import nltk
2   import re
3   from nltk.sem import extract_rels,rtuple
4   import wikipedia
5
6   p = wikipedia.page("Donald Trump")
7   text = p.content
8
9   sentences = nltk.sent_tokenize(text)
10  tokenized_sentences = [nltk.word_tokenize(sentence) for sentence in sentences]
11  tagged_sentences = [nltk.pos_tag(sentence) for sentence in tokenized_sentences]
12
13  BORN = re.compile(r'.*\bborn\b.*')
14
15  for i, sent in enumerate(tagged_sentences):
16      sent = nltk.ne_chunk(sent)
17      rels = extract_rels('PERSON', 'GPE', sent, corpus='ace', pattern=BORN, window=5)
18      for rel in rels:
19          print('{0:<5}{1}'.format(i, rtuple(rel)))


1   28   [PER: 'Fred/NNP'] 'was/VBD born/VBN in/IN the/DT' [GPE: 'Bronx/NNP']
2   53   [PER: 'Donald/NNP Jr./NNP'] '(/( born/JJ 1977/CD )/) ,/,' [GPE: 'Ivanka/NNP']
3   310  [PER: 'Obama/NNP'] 'was/VBD born/VBN in/IN the/DT' [GPE: 'United/NNP States/NNPS']
4   311  [PER: 'Obama/NNP'] 'was/VBD born/VBN in/IN the/DT' [GPE: 'U.S./NNP']
```

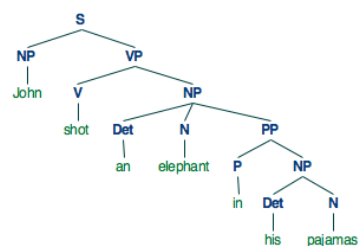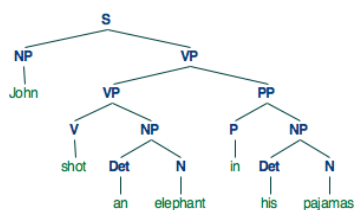## Relation Extraction based on Linguistic Analysis

- Identify the subject, predicate and object of the sentence
  - *the relation is the predicate*
  - *the subject is the source entity*
  - *the object is the target entity*
- Perform dependency parsing of the sentence

```
                        S
          _____|_____
         |                          VP
         |            _____|_____
         |           |                            NP
         |           |              _____|_____
         |           |             |                      PP
         |           |             |              _____|_____
         |           |             |             |             NP
         |           |             |             |        _____|_____
         NP          V     Det     N      P     Det            N
         |           |      |      |      |      |             |
        John       shot    an   elephant  in    his         pajamas
```

## Parsing

- Allows to analyze a sentence structure
  - *Dealing with the ambiguity of the natural language*

```python
import nltk
grammar = nltk.CFG.fromstring("""
  S -> NP VP
  PP -> P NP
  NP -> Det N | Det N PP | 'John'
  VP -> V NP | VP PP
  Det -> 'an' | 'his'
  N -> 'elephant' | 'pajamas'
  V -> 'shot'
  P -> 'in'
""")
sentence = ['John', 'shot', 'an', 'elephant', 'in', 'his', 'pajamas']
parser = nltk.ChartParser(grammar)
for tree in parser.parse(sentence):
    print(tree)
    tree.draw()
```

## Identification of a subject in parsed English tree

```
                           S
               _____|
               |                VP
               |          _____|____
               |          |            NP
               |          |      _____|__
               |          |      |          PP
               |          |      |        __|___
               |          |      |        |      NP
               |          |      |        |    __|___
              NP    V   Det    N    P   Det      N
               |     |     |     |     |    |        |
             John  shot   an  elephant  in  his   pajamas
```
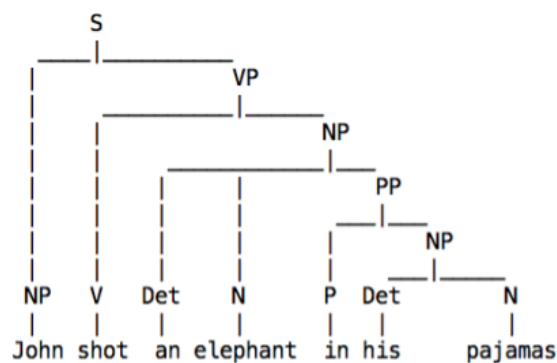
- Identification of a subject for English
  - S *indicates the sentence*
  - NP *is noun phrase*
  - VP *is verb phrase*
  - *subject is NP that is the child of S and the sibling of VP*

## Relation extraction from parsed trees

```
                           S
               _____|
               |                VP
               |          _____|____
               |          |            NP
               |          |      _____|__
               |          |      |          PP
               |          |      |        __|___
               |          |      |        |      NP
               |          |      |        |    __|___
              NP    V   Det    N    P   Det      N
               |     |     |     |     |    |        |
             John  shot   an  elephant  in  his   pajamas
```

- The subject is the NP
  - John
- The relation is the verb (V) in the verb phrase (VP)
  - shot
- The object is the noun (N) in the noun phrase (NP) in the verb phrase (VP)
  - elephant