

Handout 13: Information extraction

1. Tabular data (relational database)

a. Input

Atlas Widgets, Inc. announced on **Friday** that **Constance White** would take the helm on **March 1**. **Ms. White** will replace **Milifred Peirce** as CEO.

b. Relation **office**:

org	office	holder	begin	end
O-49	CEO	P-0155	2008 Mar 1	NA
O-49	CEO	P-0023	NA	2008 Mar 1

c. Relation **person**:

id	first	last
P-0023	Milifred	Peirce
P-0155	Constance	White

d. Relation **organization**:

id	name
O-49	Atlas Widgets, Inc.

2. Queries

a. Which organizations are located in Atlanta?

```
1  >>> ans = [org for (ent, loc) in location_table if loc == 'Atlanta']
```

b. Who is/was CEO of Atlas Widgets?

```
1  >>> orgname = dict(organization_table)
2  >>> pername = dict((id, first + ' ' + last)
3                      for (id, first, last) in person_table)
4  >>> ans = [pername[per]
5             for (org, ofc, per, _, _) in office_table
6             if ofc == 'CEO' and orgname[org].startswith('Atlas Widgets')]
```

3. Components

- Text classification; ad hoc retrieval. Which documents contain useful information?
- General-purpose preprocessing: tokenization, sentence segmentation, part of speech tagging
- (Named) entity extraction (\approx passage retrieval). Boldface phrases in example.
- Relation recognition. Underlined phrases in example.

- e. Coreference. Recognizing that “Constance White” and “Ms. White” are the same person.
- f. Integration: pulling together pieces of an event possibly from multiple sentences.

4. Preprocessing

```

1 >>> nltk.sent_tokenize('Hi there. This is a test.')
2 ['Hi there.', 'This is a test.']
3 >>> [nltk.word_tokenize(s) for s in _]
4 [['Hi', 'there', '.'], ['This', 'is', 'a', 'test', '.']]
5 >>> [nltk.pos_tag(s) for s in _]
6 [[('Hi', 'NNP'), ('there', 'EX'), ('.', '.')], [('This', 'DT'), ('is',
7 'VBZ'), ('a', 'DT'), ('test', 'NN'), ('.', '.')]]

```

5. Partial parsing (chunk parsing)

- a. Example sentence

```

1 >>> s = [('The', 'DT'), ('little', 'JJ'), ('dog', 'NN'), ('likes', 'VBZ'),
2 ...      ('smelly', 'JJ'), ('shoes', 'NN')]

```

- b. Grammar

```

1 >>> g = 'NP: {<DT>?<JJ>*<NN.*>+}'

```

- c. Parser

```

1 >>> p = nltk.RegexpParser(g)

```

- d. Tree

```

1 >>> t = p.parse(s)
2 >>> t
3 Tree('S', [Tree('NP', [('The', 'DT'), ('little', 'JJ'), ('dog',
4 'NN')]), ('likes', 'VBZ'), Tree('NP', [('expensive', 'JJ'), ('shoes',
5 'NNS')])])
6 >>> print(t)
7 (S (NP The/DT little/JJ dog/NN) likes/VBZ (NP expensive/JJ shoes/NNS))

```

6. Refinements

- a. Multiple rules can be specified; applied in sequence

```

1 g = r'''
2     NP: {<DT>?<JJ>*<NN>}
3         {<NNP>+}
4     '''

```

- b. Chinks: things that break up a chunk

7. Longest-match rule

the emergency crews hate most is domestic violence

8. Using a classifier to recognize chunks

B	O	B	I	B	I	I
PRP	VBD	DT	NN	DT	JJ	NN
He	gave	the	dog	a	new	toy

9. Evaluation

a. Test sentences

```
1 >>> from nltk.corpus import conll2000 as conll
2 >>> test = conll.chunked_sents('test.txt', chunk_types='NP')
```

b. A chunk parser

```
1 >>> p = nltk.RegexpParser(r'NP: {<[CDJNP].*>+}')
```

c. Evaluate the parser

```
1 >>> print(p.evaluate(test))
2 ChunkParse score:
3 IOB Accuracy: 87.7%
4 Precision: 70.6%
5 Recall: 67.8%
6 F-Measure: 69.2%
```

d. Why don't we measure chunking accuracy?

e. F-measure: harmonic mean of precision, recall

$$\frac{1}{\frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)}$$

Equals regular mean when $P = R$, but penalizes $P > R$ or $R > P$

10. Relation extraction: quick and dirty

a. Consider adjacent pairs of entities, classify based on the words between

```
1 [PER Ms. White] will replace [PER Milifred Peirce]
```

b. Single entities with preceding or following text: relation to larger event

```
1 [COM Atlas Widgets, Inc.] announced
2 on [DATE March 1]
```

c. Keywords not associated with adjacent entity: property of larger event

```
1 as CEO
```

11. Higher-quality approach

- a.** Parse, classify relationships between entities based on relation in parse tree
- b.** Do full interpretation. But recall usually falls as precision improves.