#### **Open Domain Information Extraction**

Günter Neumann, DFKI, 2012

### Improving TextRunner

- Wu and Weld (2010) Open Information Extraction using Wikipedia, ACL 2010
- Fader et al. (2011) Identifying Relations for Open Information Extraction, EMNLP 2011

## Wu and Weld (2010) Open Information Extraction using Wikipedia

- Novel form of self-supervised learning for open extractors: heuristic matches between Wikipedia infobox attribute values and corresponding sentences to construct training data
- Like TextRunner: avoid lexicalized features and handles unbounded set of semantic relations
- System WOE (Web-based Open Information Extraction):
  - WOE<sup>pos</sup> : POS tag features → as fast as TextRunner but better P & R
  - WOE<sup>parse</sup> : dependency-parse features → even higher P & R, but 30times slower

#### **Open Extractor - Problem Definition**

- An open information extractor is:
  - a function from a document d to a set of triples {<arg1, rel, arg2>}, where
  - arg<sub>i</sub> is a noun phrase and rel is a textual fragment indicating an implicit, semantic relation between arg<sub>1</sub> and arg<sub>2</sub>
- Restriction: Triples represent facts stated *explicitly* in the text, no inference of implicit facts
- Assumption: all relational instances are stated within a single sentence

#### Self-Supervision using Wikipedia

- Goal: Learn an open extractor without direct supervision !
- Input: Wikipedia (as source for sentences) and DBpedia (as source for cleaned infoboxes)
- Output: unlexicalized and relation-independent open extractor
- Objection: Extractor should generalize beyond Wikipedia, e.g., that can handle the general Web

### Wikipedia-based Open IE

- Key idea: automatic construction of training examples by heuristically matching Wikipedia infobox values and corresponding text.
- Used to generate generalized relation-independent extractors.

#### From infoboxes to a training set

Clearfield County, Pennsylvania		
	Statistics	Clearfield County was created in 1804 from parts of
Founded	March 26, 1804	Huntingdon and Lycoming Counties but was
Seat	Clearfield	administered as part of Centre County until 1812.
Area - Total - Land - Water	2,988 <mark>km²</mark> (1,154 mi²) sq mi ( km²) 17 km² (6 mi²), 0.56%	Its county seat is Clearfield. 2,972 km <sup>2</sup> (1,147 mi <sup>2</sup> ) of it is land and
Population - (2000) - Density	83,382 28/km²	17 km <sup>2</sup> (7 mi <sup>2</sup> ) of it (0.56%) is water. As of 2005, the population density was
		28.2/km <sup>2</sup> .

#### Architecture of WOE

- Preprocessor converts raw Wikipedia text into a sequence of sentences
- The matcher constructs training data from attributevalue pairs of infoboxes and matching sentences
- Learner acquires the open extractors using either parser features or POS features.



Figure 1: Architecture of WOE.

#### Preprocessor

- OpenNLP for sentence splitting
- NLP annotation:
  - OpenNLP for POS tagging and NP chunks
  - Stanford Parser for dependency parse; hyperlinked anchor text handled as a single token (using underscore)
- Compiling synonymies (to increase recall of the matcher)
  - Wikipedia articles contain different mentions of same entities (across pages and between infobox and Wikipedia page)
  - Wikipedia redirection pages and backward links are used to construct automatically synonym sets.



## Matcher - Constructing Training Data from Infoboxes

- Given a Wikipedia page with an infobox assumption here: such a page describes an entity - iterate through all its attributes looking for a unique sentence that contain references to both the subject of the article and the attribute value (or its synonym).
- "Stanford University" article: match <established, 1891> with "The university was founded in 1891 by …"
   → <arg<sub>1</sub>=Stanford University, rel=???, arg<sub>2</sub>=1891>
- Ordered heuristics for matching the subject, e.g., full match, synonym match, partial match (prefix/suffix of the entity's name), patterns of "the <type>" (type identification using simple patterns from Wikipedia), article's most frequent pronoun.
- Use **DBpedia's** cleaned infobox data (1,027,744 articles) as basis for attributevalue pairs → leads to **301,962 labeled sentences**

#### Extraction with Parser Features

- Construction of relation text:
  - corePath = shortest dependency path between arg1 and arg2 (expandPath = adding all adverbial and adjectival modifiers, "neg" and "auxpass" labels of the root node)
  - select all tokens of expanded Path as value for rel ("was not born" in our example)
- Generalization of patterns:
  - ignoring corePaths that don't start with subject like dependencies, s.a. nsubj, nsubj-pass → 259,046 corePaths
  - generalized-corePaths: substitute lexical words by their POS; map all noun tags to N, verb tags to V, prep tags to prep etc.
  - yields a database **DB**<sub>p</sub> of 29,005 distinct patterns
  - each pattern receives its number of matching sentences as frequency p (311 patterns have  $f_p \ge 100$ , and 3,519 have  $f_p \ge 5$ )

Dan was not born in Berkeley. The Stanford Parser dependencies are: nsubjpass(born-4, Dan-1) auxpass(born-4, was-2) neg(born-4, not-3) prep\_in(born-4, Berkeley-6)

Dan nsubjpass born prep\_in Berkeley was auxpass not

Dan nsubjpass born prep\_in Berkeley

"N nsubjpass V prep N".

#### WOEparse - A Simple Pattern Classifier

 Lookup the generalized-corePath of a test triple in DB<sub>p</sub>, and compute normalized logarithmic frequency as the probability:

 $w(p) = \frac{max(log(f_p) - log(f_{min}), 0)}{log(f_{max}) - log(f_{min})}$ 

- f<sub>max</sub> = max. freq. of patterns in DB<sub>p</sub> , f<sub>min</sub> = a controlling threshold
- Example: "Dan was not born in Berkeley" let  $f_{max} = 54,274$ ,  $f_{min} = 1$ , and  $f_p = 31,767$ , where p = "N nsubjpass" V prep N".

• then w(p) = 0.95

• Performance: 79% to 90% higher F-measure than TextRunner ! But: TextRunner is still 30X faster !

#### WOE<sup>pos</sup> - Extraction with Shallow Patterns

- Like TextRunner learn a CRF extractor WOE<sup>pos</sup> based on shallow features. (CRF -conditional random fields - classifier from the <u>Mallet ML toolkit</u>.)
- However, use **generated labeled examples**, where TextRunner learns from a small set of hand-written labeled sentences.
- Use same matching sentence set behind DB<sub>p</sub> :
  - positive examples: use matching arg<sub>1</sub> and arg<sub>2</sub> as ends of text sequences, rel used from expandedPath
  - negative examples: random NP pairs in other sentences which are not covered via DB<sub>p</sub>
- Features: POS-tags, regular expressions, combined features from 6-sized window
- Performance: 15% to 34% higher F-measure than TextRunner ! And: as fast as TextRunner !!



Thank you for protecting Wikipedia. (We're not done yet.)

#### Amazon Mechanical Turk

From Wikipedia, the free encyclopedia

The Amazon Mechanical Turk (MTurk) is a crowdsourcing Internet marketplace that enables computer programmers (known as Requesters) to co-ordinate the use of human intelligence to perform tasks that computers are unable to do yet. It is one of the suites of Amazon Web Services. The Requesters are able to post tasks known as HITs (Human Intelligence Tasks), such as choosing the best among several photographs of a store-front, writing product descriptions, or identifying performers on music CDs. *Workers* (called *Providers* in Mechanical Turk's Terms of Service) can then browse among existing tasks and complete them for a monetary payment set by the Requester. To place HITs, the requesting programs use an open Application Programming Interface, or the more limited MTurk Requester site.<sup>[1]</sup>

- 300 randomly selected sentences form three corpora:

  - - WSJ from Penn Treebank
    - Wikipedia

Experiments

- general Web
- Each sentence was examined by two people to label all reasonable triples. These triples are mixed with pseudo-negative ones.
- Using <u>Amazon Mechanical Turk</u> for verification: each triple was examined by 5 workers, and was labeled positive when more than 3 workers marked it as positive.

### **Overall Performance Analysis**



Figure 2: WOE<sup>pos</sup>performs better than TextRunner, especially on precision. WOE<sup>parse</sup> dramatically improves performance, especially on recall.

- WOE<sup>pos</sup> is better than TextRunner, especially on precision
- WOE<sup>parse</sup> performed best, especially in recall; parser helps to handle complicated and long-distance relations in difficult sentences; further experiments showed that parser errors have negligible effect !
- Extraction errors (numbers based on WSJ corpus) : 1) incorrect args from NP chunking (18,6%), 2) erroneous parses (11.9%), 3) inaccurate meaning (27,1%), 4) a pattern inapplicable for the test sentence (42,4%)

#### Further Performances







Figure 3: WOE<sup>pos</sup> achieves an F-measure, which is between 15% and 34% better than TextRunner's. WOE<sup>parse</sup> achieves an improvement between 79% and 90% over TextRunner. The error bar indicates one standard deviation.

Figure 4: WOE<sup>parse</sup>'s F-measure decreases more slowly with sentence length than WOE<sup>pos</sup> and TextRunner, due to its better handling of difficult sentences using parser features.

Figure 5: Textrnner and WOE<sup>pos</sup>'s running time seems to grow linearly with sentence length, while WOE<sup>parse</sup>'s time grows quadratically.

- WOE<sup>parse</sup> achieves best results (Figure 3), due to deeper parsing, however also has highest costs (Figure 5; 0.679 sec/sent vs. 0.022 sec/sent)
- WOE<sup>parse</sup> decreases more slowly with sentence length, also due to deeper parsing (figure 4)

# Self-supervision with Wikipedia Results in Better Training Data



Figure 6: Matching sentences with Wikipedia infoboxes results in better training data than the handwritten rules used by TextRunner.

- Generating training examples from Wikipedia differently: 1) tr = pos/neg examples selected using TextRunners patterns, 2) w = WOE's heuristic infobox matchers, 3) r = randomly
- CRF<sub>+h1-h2</sub>, where  $h_i = \{tr, w, r\}$ ; using same feature set as TextRunner

### Design Criteria for WOS<sup>parse</sup>



Figure 7: Filtering prepositional phrase attachments ( $\overline{PPa}$ ) shows a strong boost to precision, and we see a smaller boost from enforcing a lexical ordering of relation arguments (1 $\prec$ 2).

- two interesting design choices:
- require  $arg_1$  to appear before  $arg_2$  (1 < 2)
- allow corePath to contain PP-attachment (PPa)

Freitag, 27. Januar 2012

#### Conclusion

- Novel self-supervised method using automatically generated training examples from Wikipedia
- CRF extractor trained with shallow features (WOE<sup>pos</sup>) and dependency-parse features (WOE<sup>parse</sup>)
- Much better P&R compared to TextRunner
- Two sources of WOE's strong performance:
  1) the Wikipedia heuristic is responsible for the bulk of WOE's improved accuracy

2) dependency-parse features are highly informative when performing unlexicalized extraction.