

This Table is Different: A WordNet-Based Approach to Identifying References to Document Entities

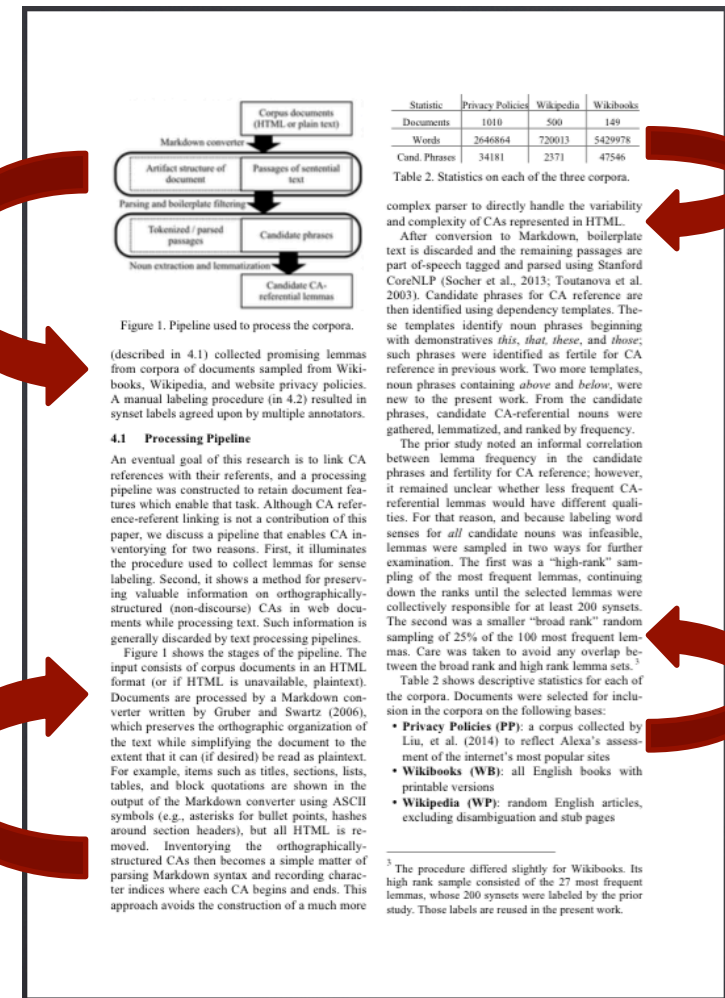
2016-01-28

Shomir Wilson, Alan W Black, and Jon Oberlander

Observations and Motivation

2

- Communication in a document (or anywhere) is not linear.
- References to document entities (DEs) are implicit but important.
- The references can serve as conduits of meaning for:
 - ▣ Labeling parts of a document
 - ▣ Contextual summarization
 - ▣ Document layout generation



References to Document Entities: Some Examples

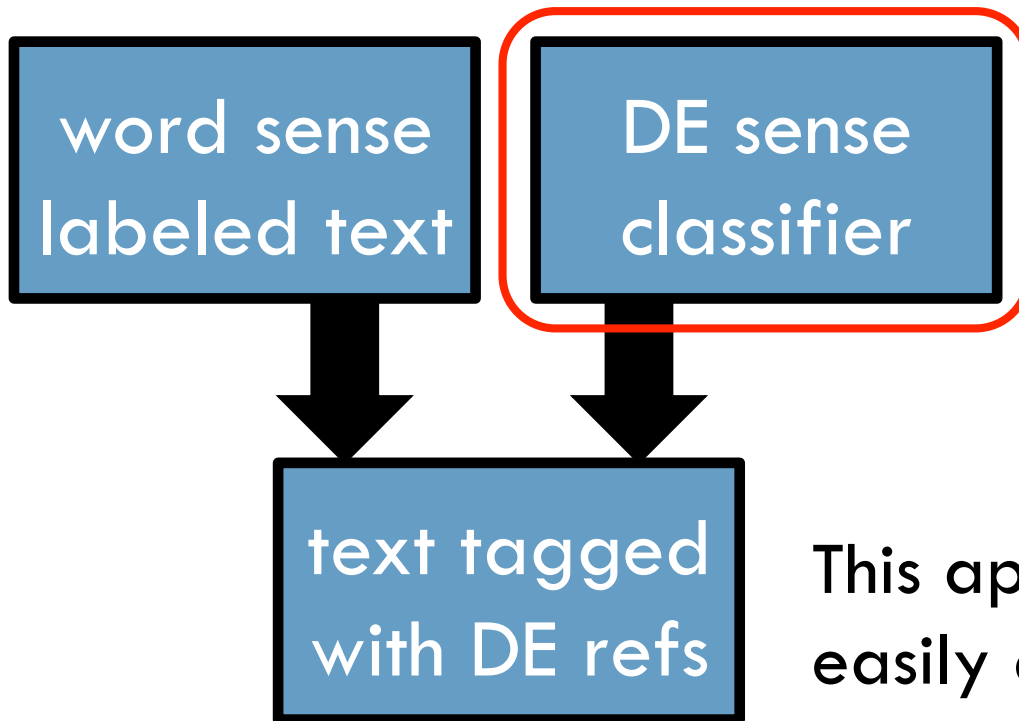
3

Category	Examples
Structural	Many of the resources listed elsewhere in this section have...
	In this chapter , we will show you how to draw...
Illustrative	Consider these sentences : [followed by example sentences]
	[following a source code fragment] ...the first time the computer sees this statement , ‘a’ is zero, so it is less than 10.
Discourse	Utilizing this idea , subunit analogies were invented...
	In this case , you’ve narrowed the topic down to “Badges.”
Non-Artifact Reference	Devices similar to resistors turn this energy into light, motion...
	What type of things does a person in that career field know?

What if we could identify the *word senses* that represent DEs? If one of those senses occurs in a phrase in text, the phrase is a reference to a DE.

A Word-Sense Based Approach

4



We developed a method to automatically label synsets from English WordNet for their capacity to refer to DEs.

This approach makes our results easily adaptable to many domains of text.

However, WSD is not a contribution of this paper.

Example

5

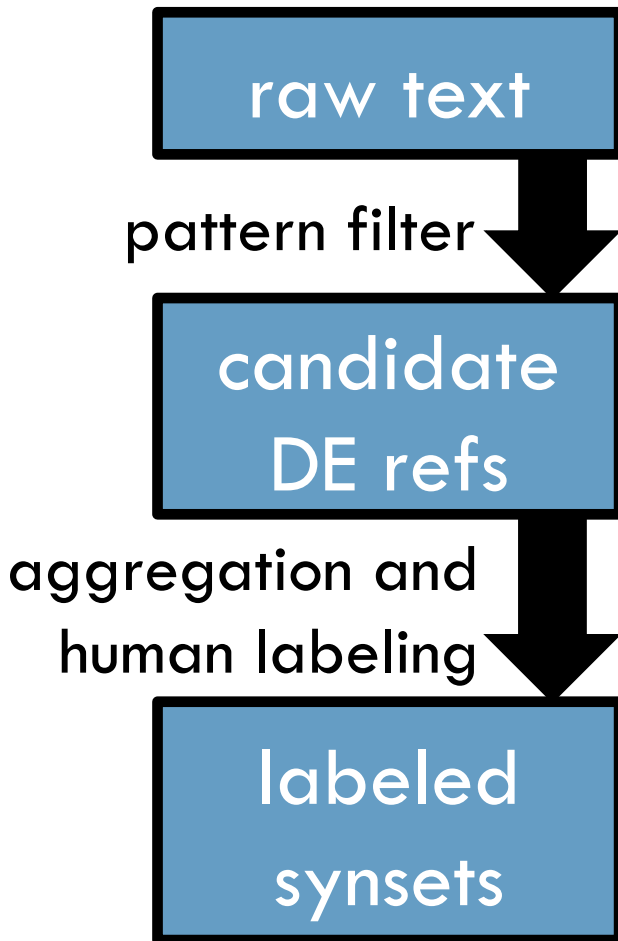
Noun

- ✓ **S: (n) table, [tabular array](#)** (a set of data arranged in rows and columns) "*see table 1*"
- ✗ **S: (n) table** (a piece of furniture having a smooth flat top that is usually supported by one or more vertical legs) "*it was a sturdy table*"
- ✗ **S: (n) table** (a piece of furniture with tableware for a meal laid out on it) "*I reserved a table at my favorite restaurant*"
- ✗ **S: (n) [mesa](#), table** (flat tableland with steep edges) "*the tribe was relatively safe on the mesa but they had to descend into the valley for water*"
- ✗ **S: (n) table** (a company of people assembled at a table for a meal or game) "*he entertained the whole table with his witty remarks*"
- ✗ **S: (n) [board](#), table** (food or meals in general) "*she sets a fine table*"; "*room and board*"

Human annotators read each synset's gloss and used a rubric to label the synset.

From Raw Text to Labeled DE Senses

6



For classifier training data, we labeled synsets (from English WordNet) of nouns associated with “candidate phrases” (i.e., likely DE references).

The candidates came from:

- Wikibooks textbooks
- Wikipedia articles
- Website privacy policies

Most Frequent Lemmas in Candidate Phrases

7

Privacy Policies		Wikibooks		Wikipedia	
Lemma	Freq.	Lemma	Freq.	Lemma	Freq.
policy	5945	case	790	page	535
information	3862	license	687	article	168
site	2151	book	686	time	67
website	1233	page	574	year	27
statement	859	example	515	period	21
party	852	section	486	list	18
company	720	way	385	case	15
cookie	638	type	363	section	15
service	585	point	344	issue	15
page	462	equation	337	game	15

A Machine Learning Problem

8

We wanted to use **supervised learning** to automatically assign labels to synsets, potentially including ones that our classifier had never seen before.

The **instances** are synsets. Our **training data** consists of synsets we labeled by hand.

The **features** are properties of synsets.

The **label** that we wish to predict for each instance is DE-referential capacity (positive or negative).

Features

9

Name (Type)	Description
ss_rank (numeric)	Rank of synset for its namesake lemma (e.g., 2 for <i>section.n.02</i>)
ss_depth (numeric)	Length of shortest hypernym chain from the instance-synset to the noun root synset
hyper_synset (binary)	Presence of <i>synset</i> in the shortest hypernym chain from the instance-synset to the root noun synset
gloss-self_word (binary)	Presence of <i>word</i> in the instance-synset's definition
gloss-hypo_word (binary)	Presence of <i>word</i> in the definitions of the instance-synset's hyponyms

Preliminary experiments led to the selection of a logistic regression classifier.

Automatic Labeling: Evaluation on High Rank Sets

10

		LOOCV	Cross-Corpus Training		
			PP	WB	WP
Evaluation Set	PP	.53/.89/.67	-	.55/.86/.67	.94/.43/.59
				.41/.77/.53	.91/.33/.49
	WB	.68/.77/.72	.90/.60/.72	-	.96/.36/.52
					.86/.49/.62
	WP	.44/.79/.56	.80/.43/.56	.57/.86/.69	-

precision/recall/F-score

Shaded boxes: results with overlapping synsets included

Automatic Labeling: Evaluation on High Rank Sets

11

		LOOCV	Cross-Corpus Training		
			PP	WB	WP
Evaluation Set	PP	.53/.89/.67	-	.55/.86/.67	.94/.43/.59
				.41/.77/.53	.91/.33/.49
	WB	.68/.77/.72	.90/.60/.72	-	.96/.36/.52
					.86/.49/.62
	WP	.44/.79/.56	.80/.43/.56	.57/.86/.69	-

- Performance similar to some discourse labeling tasks
- F-scores vary widely; inter-domain labeling harder
- Not shown here: training on two and testing on one

Automatic labeling: Evaluation on Broad Rank Sets

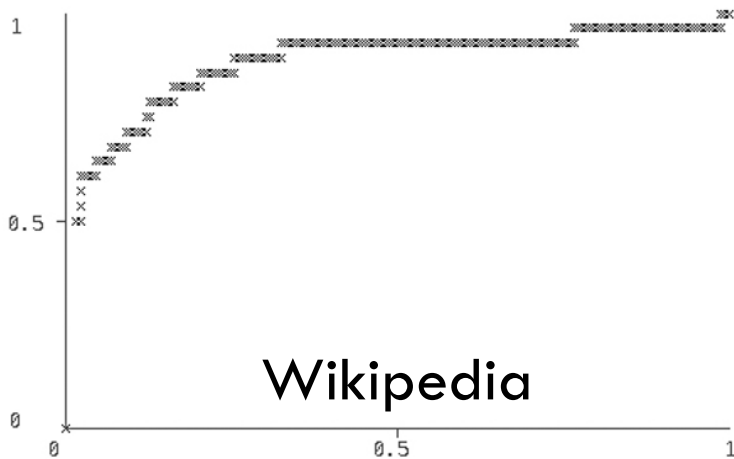
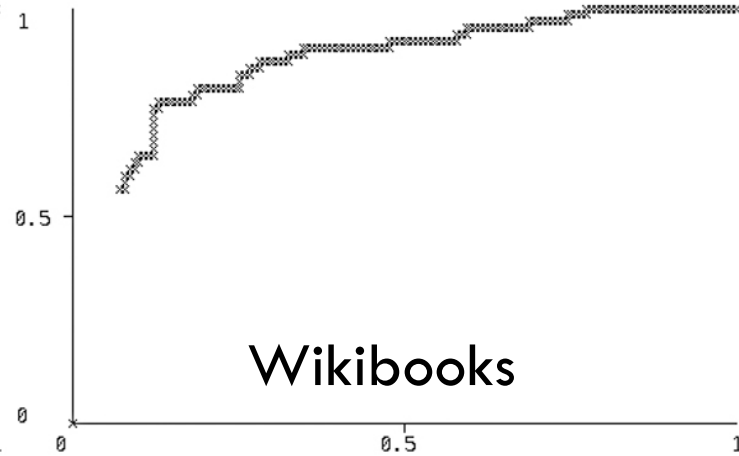
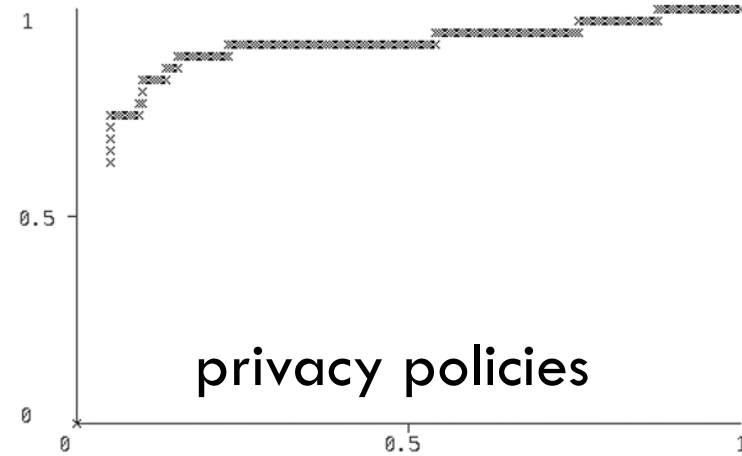
12

		Same Corpus (High Rank)	Cross-Corpus Training		
			PP	WB	WP
Eval. Set	PP	.33/.57/.42	-	.36/.71/.48	.55/.86/.67
	WB	.61/.69/.65	.60/.56/.58	-	.34/.61/.44
	WP	.34/.61/.44	.34/.72/.46	.43/.67/.52	-

- There were few positive instances in the testing data: take these results with a grain of salt.
- Performance was generally lower, suggesting different DE characteristics for the broad rank sets.

ROC Curves for LOOCV

13



Horizontal axis:
false positive rate
Vertical axis:
true positive rate

Work in Progress: Referent Resolution

14

localhost:8080/deixis/annotate_cit.jsp

Advances in domain independent linear text segmentation

Choi

0003083

Abstract

A-0 **DXC-4** [Deixis] [CA] [entri] [DEIXIS; YES; entire paper] [This paper] describes a method for linear text segmentation which is twice as accurate and over seven times as fast as the state-of-the-art (Reynar 1998) . A-1 Inter - sentence similarity is replaced by rank in the local context . A-2 Boundary locations are discovered by divisive clustering .

Introduction

S-0 Even **ATC-0** [moderately long documents] typically address **ATC-1** [several topics] or **ATC-2** [different aspects] of **ATC-3** [the same topic] . S-1 **ATC-4** [The aim] of **ATC-5** [linear text segmentation] is to discover **ATC-6** [the topic boundaries] . S-2 **ATC-7** [The uses] of **ATC-8** **DXC-1** [NP] [ATC-5] [Non-CA] [REF: ATC-5; NO] [this procedure] include **ATC-9** [information] **ATC-10** [retrieval] **ATC-11** [Hearst and Plaunt 1993] , **ATC-12** [Hearst 1994] , **ATC-13** [Yaari 1997] , **ATC-14** [Reynar 1999] , **ATC-15** [summarization] **ATC-16** [Reynar 1998] , **ATC-17** [text understanding] , **ATC-18** [anaphora resolution] **ATC-19** [Kozima 1993] , **ATC-20** [language modelling] **ATC-21** [Morris and Hirst 1991] , **ATC-22** [Beeferman et al. 1997b] and **ATC-23** [improving document navigation] for **ATC-24** [the visually disabled] **ATC-25** [Choi 2000] .

S-3 **DXC-2** [Deixis] [CA] [entri] [DEIXIS; YES; entire paper] [This paper] focuses on domain independent methods for segmenting written text . S-4 We propose that our method propose that order) , or a new algorithm that builds on previous work by Reynar (Reynar 1998) , (Reynar of a ranking scheme and the cosine similarity measure (van Rijsbergen1979) in f arity values of short text segments is statistically insignificant . S-7 Thus , **ATC-2** [rank] , for **ATC-29** [clustering] .

Background

S-8 **ATC-30** [Existing work] falls into **ATC-31** [one] of **ATC-32** [two categories] , **ATC-33** [lexical co source methods] **ATC-35** [Yaari 1997] . S-9 **ATC-36** [The former stem] from **ATC-37** [the work of Ha Hasan 1976] . S-10 **ATC-39** [They] proposed that **ATC-40** [text segments] with **ATC-41** [similar voca **ATC-43** [a coherent topic segment] . S-11 **ATC-44** [Implementations] of **ATC-45** **DXC-3** [Youmans 1991.] **ATC-46** [Reynar 1994] , **ATC-47** [Ponte and Croft 1997] , **ATC-48** [context vectors 1997] , **ATC-51** [Kaufmann 1999] , **ATC-52** [Eichmann et al. 1999] , **ATC-53** [entity repetition] **ATC-**

In progress: annotating data to support bootstrapping and machine learning

Future Work:

Detecting Structure in Online Discussions

15

[Previous](#) | [Next](#) --- Slide 3 of 40

[Back to Lecture Thumbnails](#)



kayvonf 5 months ago

Question: In 15-213's web proxy assignment you gained experience writing concurrent programs using pthreads. Think about your motivation for programming with threads in that assignment. How was it different from the motivation to create multi-threaded programs in this class? (e.g., consider Assignment 1, Program 1)

Hint: What is the difference between *concurrent* execution and *parallel* execution?



rofer 5 months ago

In 15-213 our goal was to handle several concurrent events at once. In this class our goal is going to be to speed up a single task by performing many calculations in parallel.



martin31hao 5 months ago

I think concurrency is more about the operating system giving user the illusion that different programs can run simultaneously on single CPU core, while parallelism gives the idea of speeding up a single task by breaking it into independent pieces on different CPUs (to ensure correctness of the program).



mingf 5 months ago

In my understanding, concurrent execution is that multiple tasks execute interleavedly. It is an implementation detail whether they actually execute on different cores or on a single core. Whereas parallel execution means that multiple tasks execute in different independent cores.



ESINNG 5 months ago

To some extent, I agree with mingf. Concurrent execution means that you start several tasks and they will run interleavedly. If the machine it runs on has several cores or only a single core but it supports hyper-threading, and the OS supports running multiple tasks at the same time, it can be executed in parallel. Besides, when you run the tasks in multiple machines or heterogeneous machines at the same time it's also executed in parallel. So concurrent just guarantee that all tasks will have some time slot to run, and some methods will be used to

More Future Work

16

Flexible retrieval of document entities:

- Automatic document layouts
- Semantic web applications
- Entity linking for rhetoric analysis

Thank You

17

The dataset for this paper (i.e., the set of synset labels) is available on my website.

Shomir Wilson

<http://www.cs.cmu.edu/~shomir>