

By the way: what are named/funky entities?

Irene Cramer

IGK Doctoral Colloquium – Saarland University

January 2006

Intro

My research question:

Explore semi-supervised learning methods to train a classifier that annotates classical named and funky entities.

Intro (cont.)

- o Today I'll focus on two aspects of my work:
 - n definitions and
 - n characteristics (◇ features) of named and funky entities respectively
-

Intro (cont.)

I don't want to reinvent the wheel, but

- MUC/CoNLL definitions of named entities are “sort of fuzzy”,
 - naturally, there are no definitions of funky entities at all ...
-

Outline

- o Intro

- o Entities

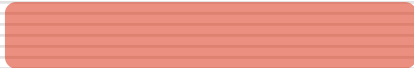
 - n Named Entities (◇ Proper Nouns)

 - n Funky Entities

 - n Construction of an inter-annotator agreement test

- o Feature sets for semi-supervised learning

Entities – named entities

“The Named Entity task consists of three subtasks (entity names, temporal expressions,  number expressions). The expressions to be annotated are "unique identifiers" of entities (organizations, persons, locations), times (dates, times), and quantities (monetary values, percentages).”

MUC-6 guidelines

Entities – named entities (cont.)

In 1989, U.S. Fish and Wildlife Service devoured 175 to 180 million dollars.

In <TIMEX TYPE="DATE">1986</TIMEX>, <ENAMEX TYPE="ORGANIZATION">U.S. Fish and Wildlife Service</ENAMEX> devoured <NUMEX TYPE="MONEY">175</NUMEX> to <NUMEX TYPE="MONEY">180 million dollars</NUMEX>.

Entities – named entities (cont.)

“Named entities are **phrases** that contain the names of persons, organizations, locations, times and quantities.”

CoNLL 2002 guidelines

Entities – named entities (cont.)

Wolff B-PER	in O
, O	the O
currently O	final O
a O	years O
journalist O	of O
in O	the O
Argentina B-LOC	seventies O
, O	in O
played O	Real B-ORG
with O	Madrid I-ORG
Del B-PER	. O
Bosque I-PER	

O: no entity; B-PER: beginning of a person name;
I-PER: within a person name etc.

Entities – named entities (cont.)

Named entities (more or less) equal
proper nouns.

Entities – proper nouns (cont.)

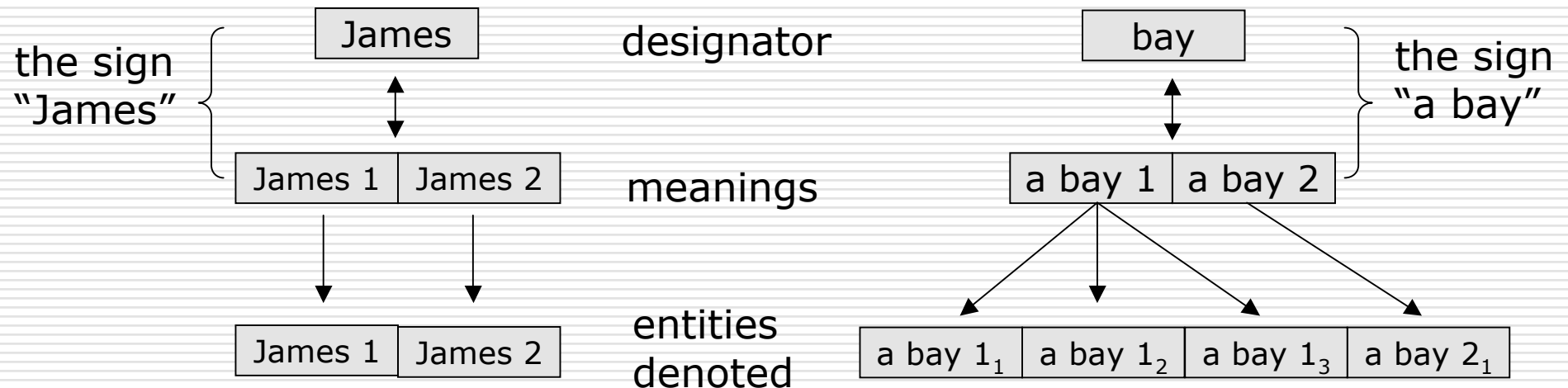
- o Proper nouns:
 - n are linguistic more precisely **lexical signs** (◇ semiotic signs);
 - n concerning linguistical properties, they **resemble substantives** (nouns);
 - n (in most languages proper nouns are **capitalized.**)
-

Entities – proper nouns (cont.)

- o Naive assumption: proper nouns refer to just “*one thing*” in the world...but
 - n **Jane, Chris**, etc. certainly refer to more than one person...
 - n **the current president of France** certainly refers to *one* person but is not a proper noun all the same
-

Entities – proper nouns (cont.)

- o Naive assumption: proper nouns refer to just “*one thing*” in the world...



adapted from a graphic in Wimmer (1973): Der Eigenname im Deutschen

Entities – proper nouns (cont.)

There is consensus among lexicographers that proper nouns aren't part of dictionaries

◇ encyclopaedic knowledge!

Entities – proper nouns (cont.)

- o Most frequent proper nouns are:
 - n person names (real and fictive ones)
 - n toponyms, hydronyms (cities, countries, rivers, etc.) and
 - n bodies (companies, organisations, etc.).
- o In addition there are names of:
 - n creatures (domestic animals),
 - n products (Big Mac, Cleanex, etc.)
 - n (political, historical ...) events, awards,
 - n ...

MUC and
CoNLL

???

Entities – funky entities

- o (Currently my) funky entities
 - n are those proper names that MUC/CoNLL do not cover;
 - n equal lexical signs (= words, word groups ≠ phrases);
 - n refer to entities in the real or a fictional world.

note: this doesn't violate classical named entity definitions

Entities – funky entities (cont.)

- o develop guidelines: how to annotate funky entities.
- o use these guidelines to conduct an inter-annotator agreement test:

∅ subjects most welcome ∅

annotation of the funky entities syntactically will follow CoNLL guidelines

Entities – inter-annotator agreement test

- o Crucial questions:
 - n Do “normal people” (and linguists) have an intuition about proper nouns?
 - n How do they identify proper nouns?
 - n Is it possible to classify proper nouns according to (more or less) unambiguous classes of entities?
 - n What is the accuracy of that labelling?
-

Entities – inter-annotator agreement test

o Subjects will annotate small corpus:

o 3 test phases:

n entity boundaries

n entity labelling

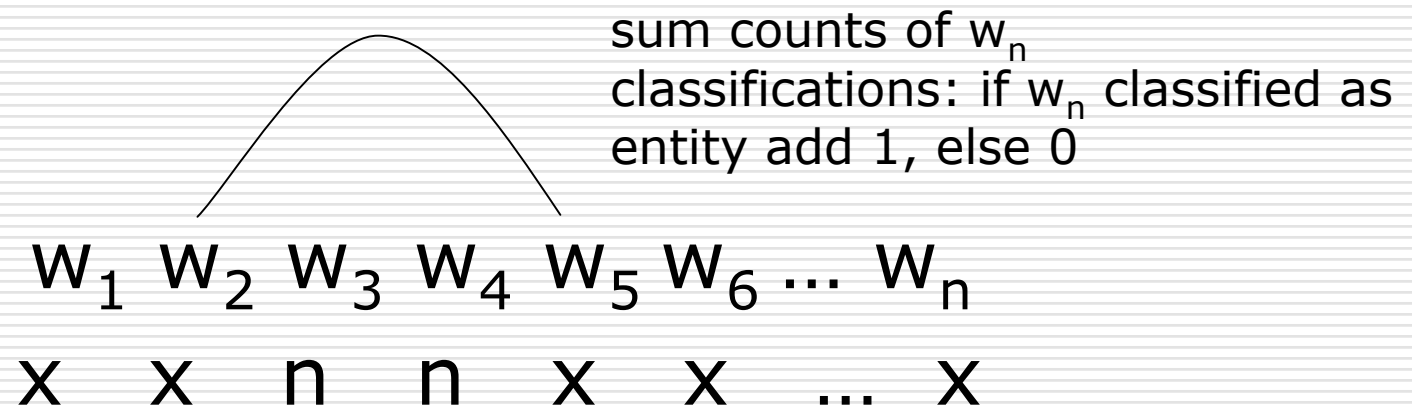
n both

} derive final
annotation guidelines



Entities – inter-annotator agreement test

- o entity boundaries



- n explore “borderland” of the entities
 - n derive new features?
-

Entities – inter-annotator agreement test

- o entity labelling
 - n guidelines include:
 - o specification of the classes (taxonomy),
 - o “boundary description”
 - o examples
 - n result analysis \diamond standard deviation / accuracy

	per	org	...	award
entity 1	2%	87%
entity 2	89%	5%
...
entity n	...	14%	...	78%

Entities – funky entities (cont.)

- After the inter-annotator agreement test: start to annotate a test corpus with these guidelines

◇ gold standard for my experiments!

Outline

- o Intro
 - o Entities
 - n Named Entities (◇ Proper Nouns)
 - n Funky Entities
 - n Construction of an inter-annotator agreement test
 - o Feature sets for semi-supervised learning**
-

Feature sets

- o Why extensively analyze features...
 - o ...isn't that contradicting my semi-supervised approach?
 - à to keep the approach as flexible as possible, I consider all features
 - à however, at the moment, I am glad about every bit of information I could possibly get
-

Feature sets

feature	class	complexity
char n-grams	n-gram level	?
suffix & prefix	word level	+
special char, cap.	word level	+
POS	synt. level	++
POS of context	synt. level	++
words in context	sem. level	+
trigger words in context	sem. level	++
one meaning per context	text level	+
word in gazetteer or lexicon?	look up	++
partial matches	look up	+

Feature sets: Correlation Entity – cap/alpha

				α	
0	α	0		5	4
		0	.		
		N	9		
	α	0	0	0	*
		0	.		
		N	9	5	3
	α	0	*	*	0
		0			.
		N			5

* 0

Feature sets: Correlation Entity – POS

0	0	0	0	0	0	0	0
					.		0
					*		0
							.
							3

* 0

0	0	0	0	0	0	0	0
					.		0
					0		0
					0		.
					8		3

* 0

e

Feature sets: Correlation Entity – POS in window

கு	பு	கு		0	0	0
		N		.1		9
	கு	கு	கு	0		*
		N		0	7 7	8 2
	கு	கு		*	5	0
		N		8	8	.6

..
கு

கு	பு	கு		0	0	0
		N		.		6
	கு	கு		0		4
		N		5	9	4
	கு	கு			*	0
		N			4	.9

..
கு

Feature sets: Correlation Entity – POS in window

த	த	கு	0	.	0	0
	த	கு	1		*	4
	த	கு	2		*	8
		கு	6		0	7

* 0

கு	த	கு	0	.	0	0
	த	கு	3		*	0
	த	கு	4		*	0
		கு	0		0	5

* 0

Feature sets: Correlation Entity – POS in window

0	0	0	0	0
*	*	*	*	*
9	7	7	7	2
0	0	0	0	.

0
*
*
0

0

0	0	0	0	*
0	0	0	0	0
*	*	*	*	*
0	0	0	0	0
4	4	4	4	8
0	0	0	0	0
6	6	6	6	0

0
*
0

Feature sets:

Correlation Entity- Feature sets

feature	correl. coeff.	sig.
suffix & prefix	--	--
character combi	--	--
special char, cap.	-0.295/-0.064	highly
POS	0.089/-0.41	highly
POS of context	$\leq 0.1 $	sometimes
words in context	--	--
trigger words in context	--	--
one meaning per context	--	--
word in gazetteer or lexicon?	--	--
partial matches	--	--

Feature sets

- test char n-grams \diamond mutual information?
 - derive new feature sets from inter-annotator agreement test \diamond what do humans do to recognize entities?
-

Thank you!
