CSC 594 Topics in AI – Applied Natural Language Processing

Fall 2009/2010

7. Named Entity Recognition

Named Entity Recognition

- Named Entities (NEs) are proper names in texts, i.e. the names of persons, organizations, locations, times and quantities
- NE Recognition (NER) is a sub-task of Information Extraction (IE)
- NER is to process a text and identify named entities – e.g. "U.N. official Ekeus heads for Baghdad."
 - [ORG U.N.] official [PER Ekeus] heads for [LOC Baghdad] .

2

• NER is also an important task for texts in **specific domains** such as biomedical texts

Source: J. Choi, CSE842, MSU; Marti Hearst, i256, at UC Berkeley

Difficulties with NER

- · Names are too numerous to include in dictionaries
- · Variations
- e.g. John Smith, Mr Smith, John
- Changing constantly
 - new names invent unknown words
- Ambiguities
 - Same name refers to different entities,
 - e.g.
 - JFK the former president
 JFK his son
 - JFK airport in NY

Source: J. Choi, CSE842, MSU









Detail	s	on	Т	rai	n	ina/	Τ	est	Sets	5
	-									

	Articles	Sentences	Tokens	English data	LOC	AUSC	ORG	PER
Training set	946	14,987	203,621	Training set	7140	3438	6321	6600
Development set	216	3,466	51,362	Development set	1837	922	1341	1842
Test set	231	3,684	-46,435	Test set	1668	702	1661	1617
German data	Articles	Sentences	Tokens	German data	LOC	MISC	ORG	PER
Training set	-553	12,705	206,931	Training set	4363	2288	2427	2773
Development set	201	3.068	51,444	Development set	1181	1010	1241	1401
Test set	155	3.160	51,943	Test set	1035	670	773	1195
Table 1: Number o rach data file. Reuter	farticles, s s News	wire +	Etokens in	Table 2: Number	of name	sd entitie	s per da	ta file
Table 1: Number o cach data file. Reuter	farticles, s	wire +	I tokens in Europea	Table 2: Number	of name	d entitie	s per da	ta file







D -			D		- 0		
English test Florian	Precision 88.99%	ON, Recall 88.54%	Кеса _{F_{β=1} 88.76±0.7}	German test	Precision 83.87%	Recall	S $F_{\beta=1}$ 72.41±1.3
Chieu	88.12%	88.51%	88.31 ± 0.7	Klein	80.38%	65.04%	71.90 ± 1.2
Klein	85.93%	86.21%	86.07 ± 0.8	Zhang	82.00%	63.03%	71.27 ± 1.5
Zhang	86.13%	84.88%	85.50 ± 0.9	Mayfield	75.97%	64.82%	69.96 ± 1.4
Carreras (b)	84.05%	85.96%	$85.00 {\pm} 0.8$	Carreras (b)	75.47%	63.82%	69.15 ± 1.3
Curran	84.29%	85.50%	84.89 ± 0.9	Bender	74.82%	63.82%	68.88 ± 1.3
Mayfield	84.45%	84.90%	84.67 ± 1.0	Curran	75.61%	62.46%	68.41 ± 1.4
Carreras (a)	85.81%	82.84%	84.30 ± 0.9	McCallum	75.97%	61.72%	68.11 ± 1.4
McCallum	84.52%	83.55%	84.04 ± 0.9	Munro	69.37%	66.21%	67.75 ± 1.4
Bender	84.68%	83.18%	83.92 ± 1.0	Carreras (a)	77.83%	58.02%	66.48 ± 1.5
Munro	80.87%	84.21%	82.50 ± 1.0	Wu	75.20%	59.35%	66.34 ± 1.3
Wu	82.02%	81.39%	81.70 ± 0.9	Chieu	76.83%	57.34%	65.67 ± 1.4
Whitelaw	81.60%	78.05%	79.78 ± 1.0	Hendrickx	71.15%	56.55%	63.02 ± 1.4
Hendrickx	76.33%	80.17%	78.20 ± 1.0	De Meulder	63.93%	51.86%	57.27 ± 1.6
De Meulder	75.84%	78.13%	76.97 ± 1.2	Whitelaw	71.05%	44.11%	54.43 ± 1.4
Hammerton	69.09%	53.26%	60.15 ± 1.3	Hammerton	63.49%	38.25%	47.74 ± 1.5
Baseline	71.91%	50.90%	59.61 ± 1.2	Baseline	31.86%	28.89%	30.30 ± 1.3
		-	* Not signif	cantly different			
Source: Marti Hea	arst, i256, at	UC Berke	ley				1



State of the Art Performance

- Named entity recognition
 - Person, Location, Organization, ...
 - F1 in high 80's or low- to mid-90's
- However, performance depends on the entity types
 [Wikipedia] At least two <u>hierarchies</u> of named entity types have been
 proposed in the literature. <u>BBN</u> categories [1], proposed in 2002, is
 used for <u>Question Answering</u> and consists of 29 types and 64
 subtypes. Sekine's extended hierarchy [2], proposed in 2002, is made
 of 200 subtypes.
- Also, various domains use different entity types (e.g. concepts in biomedical texts)

11

12

Sequence Labeling
Inputs: x = (x₁, ..., x_n)
Labels: y = (y₁, ..., y_n)
Dypical goal: Given x, predict y
Partof-speech tagging
Named Entity Recognition (NER)
Target class/label is the entity type, in IOB notation

Methods for Sequence Labeling

- Typically the following methods are used for NER: a) Hidden Markov Model (HMM)
 - b) Maximum Entropy Classifier (MaxEnt)
 - c) Maximum Entropy Markov Model (MEMM)
 - d) Conditional Random Fields (CRF)
- These are all classifiers (i.e., supervised learning) which model sequences (rather than individual random variables)

13

Instance/word Features for NER

Characteristics of the token & text in a surrounding window.

Lowland items	Explanation
Lexical hems	The token to be labeled
Stemmed lexical items	Stemmed version of the target token
Shape	The orthographic pattern of the target word
Character affixes	Character-level affixes of the target and surrounding words
Part of speech	Part of speech of the word
Syntactic chunk labels	Base-phrase chunk label
Gazetteer or name list	Presence of the word in one or more named entity lists
Predictive token(s)	Presence of predictive words in surrounding text
Bag of words/Bag of N-grams	Words and/or N-grams occurring in the surrounding contex
Shape/orthographic	features
Shape/orthographic	features
Shape/orthographic	features Example
Shape/orthographic	features Example commings
Shape/orthographic Shape Lower Capitalized	features Example cummings Washington
Shape/orthographic Shape Lower Capitalized All caps Mired ease	features Example commings Washington IRA alter
Shape/orthographic Shape Lower Capitalized All caps Mixed case Capitalized character with	features Example cummings Washington IRA eBay U
Shape/orthographic Shape Lower Capitalized All caps Miscel case Capitalized character with Ends in dirit	features Example commings Washington BA ePA iperiod II 99
Shape/orthographic Shape Lower Capitalized All caps Mixed case Capitalized character with Ends in digit Contains hyphen	features Example cummings Washington IRA eBay H. A9 H-P







iven a sequence	of observat	ions:			
Yesterday Ped	ro Domingos	spoke th	is exam	ole senten	ce.
nd a trained HMM			per loci bac	son name ation name kground]
				0	
Yesterday Per	dro Domingo	s spoke ti	ils exam	ple senter	ice.
Yesterday Per	dro Domingo	S spoke ti) O	Die senter	10





However...

18

- These arbitrary features are not independent.
 - Multiple levels of granularity (chars, words, phrases)
 - Multiple dependent modalities (words, formatting, layout)
 - Past & future
- Possible solutions:
 - 💥 Model dependencies
 - X Ignore dependencies
 - ✓ Conditional Sequence Models

Source: Andrew McCallum, UMass Amherst



- Conditional model p(y|x).
 - Do not waste effort modeling p(x), since x is given at test time anyway.
 - Allows more complicated input features, since we do not need to model dependencies between them.

19

20

- Feature functions f(x,y):
 - f1(x,y) = { word is Boston & y=Location }
 - f2(x,y) = { first letter capitalized & y=Name } - f3(x,y) = { x is an HTML link & y=Location}

Source: Andrew McCallum, UMass Amherst

Maximum Entropy Models

- Same as multinomial logistic regression (thus an exponential model for n-way classification)
- · Features are constraints on the model.
- Of all possible models, we choose that which maximizes • the entropy of all models that satisfy these constraints.
- Choosing one with less entropy means we are adding information that is not justified by the empirical evidence.

Source: Jason Balbridge, U of Texas at Austin



Source: Jason Balbridge, U of Texas at Austin

01131310	nt as	ssigr	nment	M	laximum	entr	ору	assigr
(x,y)	0	1			p(x,y)	0	1	
а	.5	.1			а	.3	.2	
b total	.1	.3	1.0	-	b total	.3	.2	1.0











Source: Andrew McCallum, UMass Amherst

CRFs

- CRFs are widely used and applied in NLP research. CRFs give state-the-art results in many domains and tasks.
 - Part-of-speech tagging
 - Named Entity Recognition
 - Gene prediction

 - Chinese word segmentationExtracting information from research papers.
 - etc.

Source: Andrew McCallum, UMass Amherst

• For a standard linear-chain structure:

 $P(\mathbf{y} \mid \mathbf{x}) = \prod_{t} \Psi_{k}(y_{t}, y_{t-1}, \mathbf{x})$ $\Psi_{k}(y_{t}, y_{t-1}, \mathbf{x}) = \exp\left(\sum_{k} \lambda_{k} f(y_{t}, y_{t-1}, \mathbf{x})\right)$

Source: Andrew McCallum, UMass Amherst

29

28