

Recap

- Topics covered so far:
 Word segmentation, syntax, discourse,
 - maximum-likelihood, EM, EM, EM.
- Today:
 - Morphology
 - When not to use EM (and how).







- Frequency

Yarowsky & Wicentowski (2000)

- A "minimally-supervised" alignment algorithm.
 - Input: easily available resources.
 - Output: alignments between morphological variants of each word (including irregulars).
- Example of combining many sources of information to achieve excellent results.



Output								
 List of roots and inflected forms, with morphophonological analysis and POS: 								
		STEM				ĺ		
	ROOT	CHANGE	SUFFIX	INFLECTION	POS	1		
	take	$ake \rightarrow ook$	$+\epsilon$	took	VBD			
	take	$e \rightarrow \epsilon$	+ing	taking	VBG			
	take	$\epsilon \rightarrow \epsilon$	+s	takes	VBZ			
	take	$e \rightarrow \epsilon$	+en	taken	VBN			
	skip	$\epsilon \rightarrow p$	+ed	skipped	VBD			
	defy	$y \rightarrow i$	+ed	defied	VBD			
	defy	$y \rightarrow ie$	+s	defies	VBZ			
	defy	$\epsilon \rightarrow \epsilon$	+ing	defying	VBG			
	jugar	$gar \rightarrow eg$	+a	juega	VPI3S			
	jugar	$gar \rightarrow eg$	+an	juegan	VPI3P			
	jugar	$ar \rightarrow \epsilon$	+amos	jugamos	VPI1P			
	tener	$ener \rightarrow ien$	+en	tienen	VPI3P			
	Table 1: Target output (English and Spanish)							
	Table from Y&W (2000)							

Method

- Bootstrapping
 - Create several weak learners using complementary sources of information.
 - Each produces a ranked list of related pairs, with confidence scores.
 - Combine individual decisions to get better decisions.
 - Retrain individual components on output.
 - Repeat to convergence.

Sources of information

- 1. Frequency similarity
 - Inflected forms have similar frequencies:



Sources of information

- 2. Context similarity
 - Inflected forms have similar arguments: read the book, reading a book.
 - To avoid confusion from function words, use regular expressions:
 CW-subj (AUX|NEG)* Verb DET? CW* CW-obj

Sources of information

- 3. Orthographic similarity
 - Inflected forms have similar spellings.
 - Use weighted edit distance to compute similarity, and retrain weights.

Sources of information

- 4. Morphological rule probabilities
 - How often does each rule occur with each stem-final context, suffix, POS?

Context	Change	Suffix	Count	Examples	
ray	$\epsilon \rightarrow \epsilon$	+ed	5	spray, stray,	
ay	$\epsilon \rightarrow \epsilon$	+ed	13	play, spray,	1
oy	$\epsilon \rightarrow \epsilon$	+ed	3	annoy, enjoy,	
ey	$\epsilon \rightarrow \epsilon$	+ed	5	obey, key,	
fy	$y \rightarrow i$	+ed	21	beautify,	
ry	$y \rightarrow i$	+ed	7	carry,	
dy	$y \rightarrow i$	+ed	4	bloody,	
y	$y \rightarrow i$	+ed	43	carry,	1
y	$\epsilon \rightarrow \epsilon$	+ed	21	spray,	
y	$\epsilon \rightarrow \epsilon$	+ing	83	carry, spray,	
e	$e \rightarrow \epsilon$	+ed	728	dance,	
e	$e \rightarrow \epsilon$	+ing	783	dance, take,	
e	$\epsilon \rightarrow \epsilon$	+ing	1	singe	
					 Table from Y&VV (2000)

Parameter re-estimation

- Update alignment pairs by combining individual weighted lists.
- New list of pairs is used to re-estimate

 frequency ratios between inflected forms.
 - weights for edit distance measure.
 - probabilities for morphological rules.

Results

• Evaluated on English verbs:

Combination	# of	All	Highly	Simple	Non-
of Similarity	Iter-	Words	Irregular	Concat.	Concat.
Models	ations	(3888)	(128)	(1877)	(1883)
FS (Frequency Sim)	(Iter 1)	9.8	18.6	8.8	10.1
LS (Levenshtein Sim)	(Iter 1)	31.3	19.6	20.0	34.4
CS (Context Sim)	(Iter 1)	28.0	32.8	30.0	25.8
CS+FS	(Iter 1)	32.5	64.8	32.0	30.7
CS+FS+LS	(Iter 1)	71.6	76.5	71.1	71.9
CS+FS+LS+MS	(Iter 1)	96.5	74.0	97.3	97.4
CS+FS+LS+MS	(Convg)	99.2	80.4	99.9	99.7

Right: knew/know, made/make, brought/bring
Wrong: got/go, slew/slit, went/want

Schone & Jurafsky (2001)

- Combining multiple sources of information in a completely unsupervised way.
- · Input: text corpus.
- Output: "conflation sets" of related words: {abuse, abusive, abusing, abuses, abusers, abused}

Sources of information

• Orthography

 Initially identify possible related word pairs with different prefix/suffix.

- Semantics
 - Use LSA to compute semantic vectors for words.
 - Estimate whether semantic correlation between pairs is significantly greater than chance.

Sources of information

- Syntax
 - Compute local context vectors for words.
 - Estimate whether syntactic correlation between pairs is significantly different from chance.
- Frequency
 - Eliminate pair relations that are too infrequent.
- Transitive closure
 - Pair is related if there is a path between words.

Results, discussion

- Tested on English, German, Dutch.
 - Performs better than Linguistica (Goldsmith, 2001)
 Each source of information improves scores.
- Weaknesses:
 - Lots of free parameters!
 - Heuristic search, no model.

Related work

- Wicentowski (2004)

 Extends morphological rule structure and evaluates on 32 langs (incl. lcelandic, Hindi, Estonian, Klingon).
- Monson (2007), Dasgupta & Ng (2007), etc.
 Other procedural methods for morphology induction.
- Yarowsky (1995)
 Semi-supervised bootstrapping algorithm.
 - Blum & Mitchell (1998)
 - Co-training.
- Abney (2004)
 - Mathematical analysis of Yarowsky algorithm.

Summary

- Advantages of procedural methods for morphology induction:
 - Can combine many sources of information.
 Often achieve good performance.
- Disadvantages:
 - No probabilistic model.
 - What is being optimized?
 - How to combine into larger systems?

Model-based induction

- Thought experiment: use maximumlikelihood estimation to learn morphology.
 - Generative model:
 - generate morphological class c
 - generate stem *t* conditioned on class
 - generate suffix *f* conditioned on class P(c,t,f) = P(c)P(t/c)P(f/c)
- What happens?





Is this always feasible?

Another solution

• Introduce a prior:

$$P(\theta \mid d) \propto P(d \mid \theta) P(\theta)$$

 $\hat{\theta} = \operatorname{argmax} P(d \mid \theta) P(\theta)$

· What sort of prior should we use?

Obligatory Chomsky quote

In careful descriptive work, we almost always find that one of the considerations involved in choosing among alternative analyses is the simplicity of the resulting grammar. If we can set up elements in such a way that very few rules need be given about their distribution, or that these rules are very similar to the rules for other elements, this fact certainly seems to be a valid support for the analysis in question. It seems reasonable, then, to inquire into the possibility of defining linguistic notions in the general theory partly in terms of such properties of grammar as simplicity. (pp. 113-114)

Lis tempting, then, to consider the possibility of devising a notational system which converts considerations of simplicity into considerations of length...More generally, simplicity might be determined as a weighted function of the number of symbols, the weighting devised so as to favor reductions in certain parts of the grammar. (p. 117)

(Chomsky, 1955) [emphasis mine]

Minimum description length

(Rissanen, 1989)

- Derived from information theory:
 Need to encode the corpus so that
 - Codebook (grammar) is short.
 - Encoded corpus is also short.
 - For codebook *h* and corpus *d*, minimize Length(*h*) + Length(encoding_h(*d*))

Corpus: walk walks jump jumps jumped eats								
Codebook 1:		Со	Codebook 2:			Codebook 3:		
walk walks jump jumps 1 jumped 11 eats 111	0 10 110 110 110 110	wal jum eat s ed	k ip	0 10 110 1110 11110		walk jump eat s ed	110 0 1110 10 11110	
Enco Enco Enco	:	010110111011110111110 001110101011101011110110						

Information theory

- If we choose codes for each item optimally, the length of the encoded corpus will be -log₂ P(d/h).
- We want to find

$$\hat{h} = \operatorname{argmin} \left(\operatorname{len}(e_h(d)) + \operatorname{len}(h) \right)$$

$$= \operatorname{argmin} \left(-\log_2 P(d \mid h) + \operatorname{len}(h) \right)$$

$$= \operatorname*{argmax}_{h} P(d \mid h) \cdot 2^{-\operatorname{len}(h)}$$

$$= \operatorname{argmax} P(d \mid h) \cdot P(h)$$

where P(h) increases exponentially with length.

Goldsmith (2001)

Organizes grammar into signatures:

$$g_{I} = \begin{cases} \text{jump} \\ \text{walk} \\ \text{laugh} \\ \text{saving} \end{cases} \mathbf{X} \begin{cases} \text{NULL} \\ \text{ed} \\ \text{s} \\ \text{ing} \end{cases}$$
$$g_{2} = \begin{cases} \text{sav} \\ \text{lik} \\ \text{escap} \end{cases} \mathbf{X} \begin{cases} \text{e} \\ \text{ed} \\ \text{es} \\ \text{ing} \end{cases}$$
$$g_{3} = \begin{cases} \text{cat} \\ \text{dog} \end{cases} \mathbf{X} \begin{cases} \text{NULL} \\ \text{s} \end{cases}$$













- Orthographic rules. {sav,lik} x {e.ed.ing.es} vs. {walk,jump} x {NULL.ed.ing.s}
- Profusion of signatures.
 {NULL,ed,ing,s}, {NULL,ing,s}, {ed,ing,s}, {NULL,ed,er,ing,s}, etc.
- Less effective on agglutinative, other morphology.

Creutz & Lagus (2005; 2007)

- Designed for agglutinative morphology (e.g., Finnish, Turkish)
 - Words split into multiple morphs.
 - Model morphotactics using HMM:
 - Hidden classes: stem, suffix, prefix, none-ofabove
 - Search is iterative:
 - 1. Find initial segmentation (uses older model).
 - 2. Find morphs to split and/or join.
 - 3. Resegment and re-estimate parameters with EM.
 - 4. Repeat 2-3.

Related work

- Other Goldsmith papers
 - Beginnings of: orthographic rules, syntactic context, agglutinative morphology.
- Other MDL
 - Phonology (Ellison, 1994), word segmentation (de Marcken, 1995; Cartwright & Brent, 1996), syntax (Dowman, 2000).

Summary

- Advantages of MDL:
 - Can define (more or less) arbitrary priors.
 - Provides fair comparison between models
 - with different structures and parameters.
- Disadvantages:
 - Uninteresting choices in grammar definition may have major and non-obvious results (see Goldwater & Johnson, 2003).
 - Requires specialized search procedures.Results may be partly due to search, not MDL.