CS388: Natural Language Processing

Lecture 25: Multilinguality and Morphology



when your parser works in 90 different languages







- Project 2 back today/tomorrow
- TACC allocations
- learning"

Jacob Andreas talk Friday 11am GDC 6.302 "Language as a scaffold for



- Other languages present some phenomena not seen in English at all!
- Many algorithms so far have been developed for English
 - Some structures like constituency parsing don't make sense for other languages
 - Neural methods are typically tuned to English-scale resources, may not be the best for other languages where less data is available
- Question:
 - 1) What other phenomena / challenges do we need to solve?

2) How can we leverage existing resources to do better in other languages without just annotating massive data?

Dealing with other languages



- Morphological richness: effects and challenges
- Morphology tasks: analysis, inflection, word segmentation
- Cross-lingual tagging and parsing
- Cross-lingual embeddings and word representations

This Lecture

Morphology



- Study of how words form
- Derivational morphology: create a new *lexeme* from a base estrange (v) => estrangement (n) become (v) => unbecoming (adj)
 - May not be totally regular: enflame => inflammable
- Inflectional morphology: word is inflected based on its context
 - I become / she becomes
 - Mostly applies to verbs and nouns

What is morphology?



Morphological Inflection

In English: I arrive you arrive

we arrive you arrive

In French:



he/she/it arrives they arrive

[X] arrived

singular		plural						
second	third	first	second	th				
tu	il, elle	nous	vous	ils,				
arrives	arrive	arrivons	arrivez	arriven				
/а.віv/	/а.віv/	/a.ʁi.vɔ̃/	/a.ʁi.ve/	/a.ĸiv				
arrivais	arrivait	arrivions	arriviez	arrivaie				
/а.кі.vε/	/а.кі.vɛ/	/a.ʁi.vjɔ̃/	/a.ʁi.vje/	/а.кі.				
arrivas	arriva	arrivâmes	arrivâtes	arrivère				
/а.ві.vа/	/a.ʁi.va/	/a.ʁi.vam/	/a.ʁi.vat/	/а.кі.				
arriveras	arrivera	arriverons	arriverez	arriver				
/а.кі.vка/	/а.кі.vка/	/a.ĸi.vĸɔ̃/	/a.ĸi.vĸe/	/а.кі.				
arriverais	arriverait	arriverions	arriveriez	arrivera				
/a.ĸi.vĸɛ/	/a.ĸi.vĸɛ/	/a.ĸi.və.ĸjɔ̃/	/a.ĸi.və.ĸje/	/а.кі.				





Morphological Inflection

In Spanish:

			singular		plural				
		1st person	2nd person	3rd person	1st person	2nd person	3rd person		
		yo	tú vos	él/ella/ello usted	nosotros nosotras	vosotros vosotras	ellos/ellas ustedes		
	present	llego	llegas ^{tú} llegás ^{vos}	llega	llegamos	llegáis	llegan		
indicative	imperfect	llegaba	llegabas	llegaba	llegábamos	llegabais	llegaban		
	preterite	llegué	llegaste	llegó	llegamos	llegasteis	llegaron		
	future	llegaré	llegarás	llegará	llegaremos	llegaréis	llegarán		
	conditional	llegaría	llegarías	llegaría	llegaríamos	llegaríais	llegarían		





Not just verbs either; gender, number, case complicate things

Declension of Kind									
			singular	plural					
	indef.	def.	noun	def.	noun				
nominative	ein	das	Kind	die	Kinder				
genitive	eines	des	Kindes, Kinds	der	Kinder				
dative	einem	dem	Kind, <mark>Kinde</mark> ¹	den	Kindern				
accusative	ein	das	Kind	die	Kinder				

- Nominative: I/he/she, accusative: me/him/her, genitive: mine/his/hers
- Dative: merged with accusative in English, shows recipient of something I taught the children <=> Ich unterrichte die Kinder
 - I give the children a book <=> Ich gebe den Kindern ein Buch

Noun Inflection





Irregular Inflection

- Common words are often irregular
 - I am / you are / she is
 - Je suis / tu es / elle est
 - Yo soy / usted está / ella es
- Less common words typically fall into some regular paradigm these are somewhat predictable



Agglutinating Langauges

 Finnish/Turkish/ Hungarian (Finno-Ugric): what a preposition would do in English is instead part of the verb

		active	passive		indicative mood present tense person 1st sing. 2nd sing. 3rd sing. 1st plur.	positive halaan halaat halaa halaamme	negative en halaa et halaa ei halaa emme halaa	perfect person 1st sing. 2nd sing. 3rd sing. 1st plur.	positive olen halannut olet halannut on halannut olemme halanneet
1st		halata			2nd plur. 3rd plur. passive past tense person	halaatte halaavat halataan positive	ette halaa eivät halaa ei halata negative	2nd plur. 3rd plur. passive pluperfect person	olette halanneet ovat halanneet on halattu positive
long	1st ²	halatakseen			1st sing. 2nd sing. 3rd sing. 1st plur. 2nd plur. 3rd plur. passive	halasin halasit halasi halasimme halasitte halasitvat halattiin	en halannut et halannut ei halannut emme halanneet ette halanneet eivät halanneet ei halattu	, 1st sing. 3rd sing. 1st plur. 2nd plur. 3rd plur. passive	olin halannut oli halannut oli halannut olimme halanneet olitte halanneet olivat halanneet oli halatu
Que el	inessive ¹	halatessa	halattaessa		conditional mood present person positive preson and sing. halaisin and sing. halaisit and plur. halaisite and plur. halaisite and plur. halaisite preson present person tat sing. hala and and sing. hala and and sing. halakaamme and and and sing. halakaamme and	positive halaisin halaisi1 halaisi	negative en halaisi et halaisi ei halaisi	perfect person 1st sing. 2nd sing. 3rd sing.	positive olisin halannut olisit halannut olisi halannut
zna	instructive	halaten	_			halaisimme halaisitte halaisivat halattaisiin	emme halaisi ette halaisi eivät halaisi ei halattaisi	1st plur. 2nd plur. 3rd plur. passive perfect	olisimme halanneet olisitte halanneet olisivat halanneet olisi halattu
	inessive	halaamassa	_			negative — älä halaa älköön halatko älkäämme halatko älkää halatko	person 1st sing. 2nd sing. 3rd sing. 1st plur. 2nd plur.	positive 	
	elative	halaamasta			3rd plur. passive potential mood present person 1st sing.	halatkoot halattakoon positive halannen	älkööt halatiko älköön halattako negative en halanne	passive perfect person 1st sing.	olkoot halanneet olkoon halattu positive lienen halannut
Qued	illative	halaamaan	_		2nd sing. Ird sing. Net plur. 2nd plur. 3rd plur.	halannet halannee halannemme halannette halannevat	et halanne ei halanne emme halanne ette halanne eivät halanne	2nd sing. 3rd sing. 1st plur. 2nd plur. 3rd plur. assive	lienet halannut lienee halannut lienemme halanneet lienette halanneet lienevät halanneet lienee halattu
Sra	adessive	halaamalla	_		lominal forms nfinitives st ong 1st ² nd ^{inessive¹}	active halata halatakseen halatessa	passive halattaessa	earticiples resent east gent ^{1, 3}	active halaava halannut halaama
	abessive	halaamatta	—		instructive inessive elative illative adessive abessive	halaten halaamassa halaamasta halaamaan halaamalla halaamatta		 Usually with a possessi Used only with a posses Used only with a posses Does not exist in the ca 	halaamaton e suffix. sive suffix; this is the form for the third- te of intransitive verbs. Do not confuse
	instructive	halaaman	halattaman		th instructive nominative partitive	halaaman halaaminen halaamista halaamaisillaan	halattaman		
/tb	nominative	halaaminen			h		\ +~	. // [
401	partitive	halaamista				dla	ald	•	iug
5th ²		halaamaisillaan							

illative: "into"

Many possible forms — and in newswire data, only a few are observed

adessive: "on"

negative en ole halannut et ole halannut et ole halannut ei ole halannut ei ole halannet ette ole halannet ette ole halannet ette ole halannet ette ole halannet et ollet halannut ei ollut halannut ei ollut halannut ei ollut halannut ei ollut halannut ei ollet halannet eivät olleet halanneet eivät olle halannet ei olla halannut et ollet halannet ei olla halannut et ollet halannet ei olla halannut enme ole halannet ei olla halannut enme ole halannet ei olla halannut et ollet halannet ei ollet ha

passive halattava halattu

erson singular and third-person plural. ith nouns formed with the -ma suffix.

"



- Many languages used all over the world have much richer morphology than English (Chinese is the main exception)
 - CoNLL 2006 / 2007: dependency parsing + morphological analyses for ~15 mostly Indo-European languages
 - SPMRL shared tasks (2013-2014): Syntactic Parsing of Morphologically-Rich Languages
- Word piece / byte-pair encoding models for MT are pretty good at handling these if there's enough data

Morphologically-Rich Languages





MORGAN & CLAYPOOL PUBLISHERS

Linguistic Fundamentals for Natural Language Processing

100 Essentials from Morphology and Syntax

Emily M. Bender

SYNTHESIS LECTURES ON HUMAN LANGUAGE TECHNOLOGIES

Graeme Hirst, Series Editor

Morphologically-Rich Languages

Great resources for challenging your assumptions about language and for understanding multilingual models!

Morphological Analysis/Inflection

Morphological Analysis: Hungarian



But the government does not recommend reducing taxes. Ám a kormány egyetlen adó csökkentését sem javasolja.





Morphological Analysis

- Given a word, need to predict what its morphological features are
- Basic approach:
 - Lexicon: tells you what possibilities are
 - Analyzer: statistical model that disambiguates
- Models are largely CRF-like: score morphological features in context
- Lots of work on Arabic inflection (high amounts of ambiguity)

Predicting Inflection



Other direction: given base form + features, inflect the word Hard for unknown words — need models that generalize



			winden								
			windend								
	gewunden										
			haben								
indic	ative		subjunctive								
de	wir winden		ich winde	wir winden							
est	ihr windet	i	du windest	ihr windet							
let	sie winden		er winde	sie winden							
nd	wir wanden		ich wände	wir wänden							
lest	ihr wandet	ii	du wändest	ihr wändet							
nd	sie wanden		er wände	sie wänden							
du)	windet (ihr)										
				[

Durrett and DeNero (2013)







Other direction: given base form + features, inflect the word

*i*1

- Hard for unknown words need models that generalize
- Take a bunch of existing verbs from Wiktionary, extract these change rules using character alignments
- Train a CRF with character ngram context features to learn where to apply them

Predicting Inflection



to wind (de)



Durrett and DeNero (2013)





Morphological Reinflection





- inflection based on source side

Machine translation where phrase table is defined in terms of lemmas "Translate-and-inflect": translate into uninflected words and predict

Chahuneau et al. (2013)



Word Segmentation



- analyses?
- common pieces and split them off
- How do we do this?

Morpheme Segmentation

Can we do something unsupervised rather than these complicated

unbecoming => un+becom+ing — we should be able to recognize these

Creutz and Lagus (2002)





- $Cost(Source text) = \sum_{i=1}^{n} -\log p(m_i)$ Simple probabilistic model morph tokens
- $p(m_i) = count(token)/count(all tokens)$
- Train with EM: E-step involves estimating best segmentation with Viterbi, M-step: collect token counts
- allowed expected need needed all+owe+d expe+cted n+e+ed ne+ed+ed EO
- MO: ed has count 3 all+ow+ed expect+ed ne+ed ne+ed+ed
- Some heuristics: reject rare morphemes, one-letter morphemes
- Doesn't handle stem changes: becoming => becom + ing

Morpheme Segmentation

Creutz and Lagus (2002)







Chinese Word Segmentation

- Some languages including Chinese don't have easy whitespace tokenization
- LSTMs over character embeddings / character bigram embeddings to predict word boundaries
- Having the right segmentation can help machine translation

冬天 (winter), 能 (can) 穿 (wear) 多少 (amount) 穿 (wear) 多少 (amount); 夏天 (summer), 能 (can) 穿 (wear) 多 (more) 少 (little) 穿 (wear) 多 (more) 少 (little)。 Without the word "夏天 (summer)" or "冬天 (winter)", it is difficult to segment the phrase "能 穿多少穿多少".

• separating nouns and pre-modifying adjectives: 高血压 (high blood pressure) \rightarrow 高(high) 血压(blood pressure)

• separating compound nouns: 内政部 (Department of Internal Affairs) \rightarrow 内政(Internal Affairs) 部(Department).



Cross-Lingual Tagging and Parsing



- Labeling POS datasets is expensive
- Can we transfer annotation from high-resource languages (English, etc.) to *low-resource* languages?



Cross-Lingual Tagging



- Multilingual POS induction
- Generative model of two languages simultaneously, joint alignment + tag learning
- Complex generative model, requires Gibbs sampling for inference

Unsupervised Tagging



Snyder et al. (2008)





Tagging by Annotation Projection

Rather than doing unsupervised learning, can we use supervised learning in combination with alignments?



- Tag with English tagger, project across bitext, train French tagger?
- Can do something smarter

Das and Petrov (2011)



Cross-Lingual Tagging

	Model	Danish	Dutch	German	Greek	Italian	Portuguese	Spanish	Swedish
	EM-HMM	68.7	57.0	75.9	65.8	63.7	62.9	71.5	68.4
baselines	Feature-HMM	69.1	65.1	81.3	71.8	68.1	78.4	80.2	70.1
	Projection	73.6	77.0	83.2	79.3	79.7	82.6	80.1	74.7
our annroach	No LP	79.0	78.8	82.4	76.3	84.8	87.0	82.8	79.4
our approach	With LP	83.2	79.5	82.8	82.5	86.8	87.9	84.2	80.5
oracles	TB Dictionary	93.1	94.7	93.5	96.6	96.4	94.0	95.8	85.5
oracies	Supervised	96.9	94.9	98.2	97.8	95.8	97.2	<i>96</i> .8	94.8

- from learned tags to gold tags
- on that
- LP: additionally model that words in similar contexts should have similar tags

EM-HMM/feature HMM: unsupervised methods with a greedy mapping

Projection: project tags across bitext to make pseudogold corpus, train

Das and Petrov (2011)







- apply it to another language



Cross-Lingual Parsing

Now that we can POS tag other languages, can we parse them too?

Direct transfer: train a parser over POS sequences in one language, then

McDonald et al. (2011)





	best	-source	avg-source	gold	I-POS	pred-POS		
	source	gold-POS	gold-POS	multi-dir.	multi-proj.	multi-dir.	multi-pro	
da	it	48.6	46.3	48.9	49.5	46.2	47.5	
de	nl	55.8	48.9	56.7	56.6	51.7	52.0	
el	en	63.9	51.7	60.1	65.1	58.5	63.0	
es	it	68.4	53.2	64.2	64.5	55.6	56.5	
it	pt	69 .1	58.5	64.1	65.0	56.8	58.9	
nl	el	62.1	49.9	55.8	65.7	54.3	64.4	
pt	it	74.8	61.6	74.0	75.6	67.7	70.3	
SV	pt	66.8	54.8	65.3	68.0	58.3	62.1	
avg		63.7	51.6	61.1	63.8	56.1	59.3	

- target language
- Multi-proj: more complex annotation projection approach

Cross-Lingual Parsing

Multi-dir: transfer a parser trained on several source treebanks to the

McDonald et al. (2011)





Cross-Lingual Word Representations



Multilingual Embeddings

Input: corpora in many languages. Output: embeddings where similar words in different languages have similar embeddings

I have an apple 47 24 18 427

J' ai des oranges 47 24 89 1981

MultiCluster: use bilingual dictionaries to form clusters of words that are translations of one another, replace corpora with cluster IDs, train "monolingual" embeddings over all these corpora

Works okay but not all that well



Ammar et al. (2019)





Multilingual Sentence Embeddings



- Form BPE vocabulary over all corpora (50k merges); will include characters from every script
- Take a bunch of bitexts and train this as an MT model (one side is always English/Spanish for them, but 93 langs total), use W as sentence embeddings Artetxe et al. (2019)





Multilingual Sentence Embeddings

		FN							EN -	$\rightarrow XX$						
			fr	es	de	el	bg	ru	tr	ar	vi	th	zh	hi	SW	U
Zero-Shot Transfer, one NLI system for all languages:																
Conneau et al.	X-BiLSTM	73.7	67.7	68.7	67.7	68.9	67.9	65.4	64.2	64.8	66.4	64.1	65.8	64.1	55.7	58
(2018b)	X-CBOW	64.5	60.3	60.7	61.0	60.5	60.4	57.8	58.7	57.5	58.8	56.9	58.8	56.3	50.4	52
BERT uncased*	Transformer	<u>81.4</u>	_	<u>74.3</u>	70.5	_	_	_	_	62.1	—	_	63.8	_	_	58
Proposed method	BiLSTM	73.9	71.9	72.9	72.6	72.8	74.2	72.1	69.7	71.4	72.0	69.2	71.4	65.5	62.2	61

Train a system for NLI (entailment/neutral/contradiction of a sentence pair) on English and evaluate on other languages

Artetxe et al. (2019)







- Take top 104 Wikipedias, train BERT on all of them simultaneously
- What does this look like?

Beethoven may have proposed unsuccessfully to Therese Malfatti, the supposed dedicatee of "Für Elise"; his status as a commoner may again have interfered with those plans.

- 当人们在马尔法蒂身后发现这部小曲的手稿时,便误认为上 面写的是"Für Elise"(即《给爱丽丝》)[51]。
- Кита́й (официально Кита́йская Наро́дная Респу́блика, сокращённо — КНР; кит. трад. 中華人民共和國, упр. 中华人民

Multilingual BERT

Devlin et al. (2019)





Fine-tuning \setminus Eval	EN	DE	NL	ES
EN	90.70	69.74	77.36	73.59
DE	73.83	82.00	76.25	70.03
NL	65.46	65.68	89.86	72.10
ES	65.38	59.40	64.39	87.18

Table 1: NER F1 results on the CoNLL data.

Can transfer BERT directly across languages with some success

...but this evaluation is on languages that all share an alphabet

Multilingual BERT: Results

Fine-tuning \setminus Eval	EN	DE	ES	IT
EN	96.82	89.40	85.91	91.60
DE	83.99	93.99	86.32	88.39
ES	81.64	88.87	96.71	93.71
IT	86.79	87.82	91.28	98.11

Table 2: POS accuracy on a subset of UD languages.

Pires et al. (2019)





	HI	UR	
HI	97.1	85.9	
UR	91.1	93.8	

Table 4: POS accuracy on the UD test set for languages with different scripts. Row=fine-tuning, column=eval.

Urdu (Arabic script) => Hindi (Devanagari). Transfers well despite different alphabets!

Japanese => English: different script and very different syntax

Multilingual BERT: Results

	EN	BG	JA
EN	96.8	87.1	49.4
BG	82.2	98.9	51.6
JA	57.4	67.2	96.5

Pires et al. (2019)





Multilingual BERT

mBERT doesn't require word piece overlap between things to do well (but going from 0 overlap to some overlap helps a lot)



Figure 1: Zero-shot NER F1 score versus entity word piece overlap among 16 languages. While performance





- Universal dependencies: treebanks (+ tags) for 70+ languages
- Many languages are still small, so projection techniques may still help
- More corpora in other languages, less and less reliance on structured tools like parsers, and pretraining on unlabeled data means that performance on other languages is better than ever
- BERT has pretrained multilingual models that seem to work pretty well



- challenges
- Problems: how to analyze rich morphology, how to generate with it
- Can leverage resources for English using bitexts
- Next time: wrapup + discussion of ethics

Many languages have richer morphology than English and pose distinct