#### Natural Language Processing Phrase-based Machine Translation, etc.



**Christopher Manning** 

Borrows slides Kevin Knight and Dan Klein

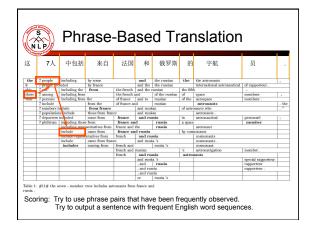


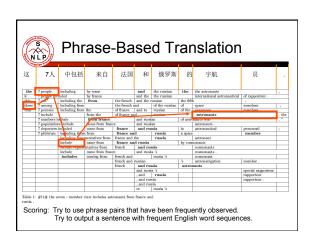
## Feature gains

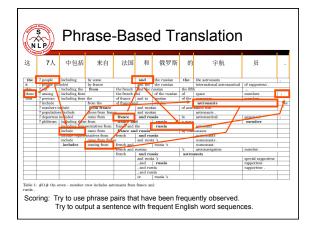
- The core numeric features should get you a decent baseline MT system
- Expect and be pleased by getting small incremental gains from features you devise
- 0.25 BLEU from a feature is good
- 0.5 BLEU from a feature is fantastic

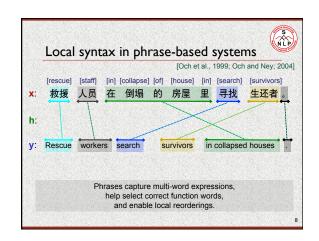
#### Phrase-Based Translation Overview lo haré rápidamente . tries different segmentations, Input: Translations: I'll do it quickly . translates phrase by phrase, quickly I'll do it . and considers reorderings.

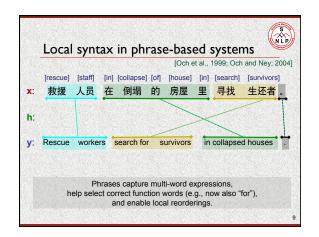
						<u>~</u>	uı	nslation		
这	7人	中包括	来自	法国	和	俄罗斯	的	宇航	员	
the	7 people	including	by some by france		and the	the russian	the	the astronauts	,	3
it	7 people inc						1 10 000	international astronautical	of rapporteur .	
this	7 out	including the	from	the french	and the	of the ressian	the fifth		members	_
these	7 among	including from		of france	and to		of the	space		
that	7 persons 7 include	including from		of france at		russian	or the	aerospace	members .	- 0
			from the	of france at		russian		astronauts		. the
	7 numbers include		from france				onauts who			
			those from france		and russian		astronauts .			
			come from	france	and ru		in	astronautical	personnel	i
	7 philtrum	including those from				russia	a space		member	
					france and the russia			astronaut		
		include			france and russia		by cost	y cosmonauts		
		include representatives from		french and russia				cosmonauts		
		include						cosmonauts .		
		includes	coming from french and					cosmonaut		
				french and re				astronavigation member .		_
				french			astro	stronauts		
					and russ				special rapporteur	
					, and ru-	russia			rapporteur	
									rapporteur.	_
	_				, and rus	sia   vocain 'c				
					or russia 's					

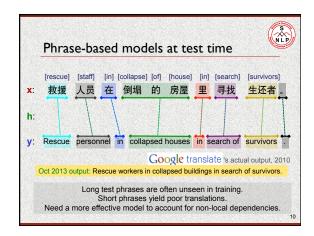


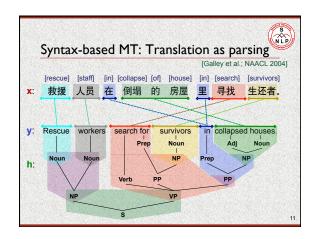


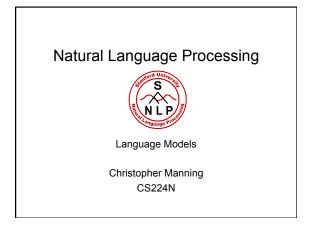














### Language Models

- Traditional grammars (e.g., regular, context free) give a hard ("categorical") model of the sentences in a language
- For NLP, and other applied work, a probabilistic model of a language is much more useful
  - It says what people usually say (next)
  - It enables more fine-grained prediction and
- Called a Language Model ... strange but standard



### Uses of language models

- Speech recognition
- "I saw a van" is a more likely sentence than "eyes awe of an"
- OCR & Handwriting recognition
  - More probable sentences are more likely correct readings
- Machine translation
  - More likely sentences are probably better translations
- (Fluent Text) Generation
  - More likely sentences are probably better NL generations
  - Context sensitive spelling correction
  - "Their are problems wit this sentence."
- Predictive text input systems
  - Please turn your cell phone of\_
  - Google query completion suggestions



#### Uses of language models

- Text classification
- A topic is a language ("the language of finance")
- Gender/style detection
- Information retrieval: Language models for IR
  - Treat either or both the query or each document as a
  - One of the most pursued research approaches recently (at UMass/CMU)
- Certain aspects of grammar checking
- E.g., preposition choice
- Text compression
- Way better than gzip/bzip2 for human language text



#### Probabilistic Language Models

- Idea is to build models which assign scores to sentences
  - P(I saw a van) >> P(eyes awe of an)
  - Not really grammaticality
    - P(artichokes intimidate zippers) ≈ 0
- Formally, a probability distribution over sentences of a language ... sums to 1 over whole language
- One option: empirical distribution over corpus sentences?
  - Problem: doesn't generalize (at all)
  - Whereas languages are infinite



## Probabilistic Language Models

Three major components of generalization

- Decomposition: sentences generated in small
- Discounting: save some probability mass for the possibility of unseen events
- **Backoff** contexts that words are generated from to equivalence classes of contexts which generalize better

After that, there are a lot of details

 But the details are very important in getting good performance in many NLP systems



#### Decomposition:

### N-Gram Language Models

No loss of generality to break sentence probability down with the chain rule

$$P(w_1 w_2 ... w_n) = \prod_{i=1}^{n} P(w_i \mid w_1 w_2 ... w_{i-1})$$

- Too many histories!
  - P(??? | No loss of generality to break sentence) ?
    P(??? | the water is so transparent that) ?
- N-gram solution: assume each word depends only on a short linear history (a Markov assumption) = equivalence classing

$$P(w_1w_2...w_n) = \prod_{i} P(w_i \mid w_{i-k}...w_{i-1})$$
$$= \prod_{i} P(w_i \mid w_{i-1}) \text{ for bigram}$$



## Character-level



- Claude Shannon (1951): the entropy of English

$$H(X) = E_P \log \frac{1}{P(X)}$$
$$= -\sum_{x \in \mathcal{X}} P(x) \log P(x)$$

Cross entropy

$$H(S \mid M) = \frac{-\log_2 P_M(S)}{\mid S \mid} = \frac{-\sum\limits_{i=1,\dots N} \log_2 P_M(w_i \mid w_{1,\dots,i-1})}{N} \underbrace{\sum\limits_{j} \log_2 P_M(w_j \mid w_{j-1})}^{\text{e.g.}}$$



### Word level

... denied the \_\_\_\_



The Palestinian security chief in Gaza denied the report Judge Kathleen Kennedy-Powell denied the motion to strike Pineau-Valencienne has denied the charges

The FDA denied the group's request

the show's writer and co-star, denied the characters had real-life The district attorney's office had denied the KCBS-TV report

Coleman denied the charge Defense attorney Al Kitching denied the allegations

Local officials have consistently denied the existence of armed

Kraft has categorically denied the remarks

Goddard has denied the charges

congressional employees are denied the legal protections

who denied the accusation of the woman



### Discounting/Smoothing

- We often want to make estimates from sparse statistics
  - P(w | denied the) 3 allegations

  - 2 reports 1 claims
  - 1 request
  - 7 total
- Smoothing flattens spiky distributions so they generalize better
  - P(w | denied the) 2.5 allegations 1.5 reports 0.5 claims

  - 0.5 request

  - Very important all over NLP, but easy to do badly! Illustration with trigrams (h = previous word, could be anything).





## Discounting/Backoff/Interpolation

- $P(w_i|h)$  is just a multinomial
  - but we need to estimate it well
  - We want to know how often a word follow some history h
  - There's some true distribution  $P(w \mid h)$
  - We saw some small sample of N words from P(w | h)
  - We want to reconstruct a useful approximation of  $P(w \mid h)$  Counts of events we didn't see are always too low

  - Counts of events we did see are in aggregate too high
- Discounting: providing mass for what we haven't seen
- Backoff: Increasing N by decreasing the amount of
- Interpolation between backed-off distributions: how to allocate that mass amongst unseen events



#### Evaluation

- What we want to know is:
  - Will our language model prefer natural sentences?
  - Does it assign higher probability to "real" or "frequently observed" sentences than "rarely observed" sentences?
- We train parameters of our model on a training set.
- To evaluate how well our model works, we look at the model's performance on some different data
- This is what happens in the real world; we want to know how our model performs on data we haven't seen So a test set. A dataset which is different from our
  - training set Preferably totally unseen/unused!
- So we can do this, we do model development with a separate development test (devtest) set



# Language models

- Language models are a cool technology
- You can have them for not only a language like "English" but for particular languages/topics

  Papers about language modeling

  - "Spam emails"
  - Seventeenth century novels
- Because they flexibly model higher order context, they can be very powerful models
  - And work very well

Look at the videos and J&M chapter 4!