Lecture 4: Machine Translation and IBM Model I (Learning)

Intro to NLP, CS585, Fall 2014 http://people.cs.umass.edu/~brenocon/inlp2014/ Brendan O'Connor (http://brenocon.com)

- Homework questions?
- Today
 - Quick LM review
 - Machine translation -- learning
- Next week
 - Machine translation -- decoding, broader issues

Interpolation LM [from last time]

Let \hat{P} be the MLE probability distribution Interpolation estimate is:

Mixing weights: $\lambda_0 + \lambda_1 + \lambda_2 = 1$

Hamlet from NLTK,V=4812

http://people.cs.umass.edu/~brenocon/inlp2014/lectures/04-hamlet_ngrams.txt

$$= \lambda_2 \hat{P}(w|\text{blacke as}) + \lambda_1 \hat{P}(w|\text{as}) + \lambda_0 \hat{P}(w)$$

Hamlet from NLTK,V=4812

http://people.cs.umass.edu/~brenocon/inlp2014/lectures/04-hamlet_ngrams.txt

$$= \lambda_2 \frac{\hat{P}(w|\text{blacke as}) + \lambda_1 \hat{P}(w|\text{as}) + \lambda_0 \hat{P}(w)}{\checkmark}$$

3 blacke as

- l blacke as death
- l blacke as hell
- I blacke as his

3 non-zeros

Hamlet from NLTK, V=4812

http://people.cs.umass.edu/~brenocon/inlp2014/lectures/04-hamlet_ngrams.txt

 $= \lambda_2 \hat{P}(w|\text{blacke as}) + \lambda_1 \hat{P}(w|\text{as}) + \lambda_0 \hat{P}(w)$ 205 7 as he as as peace as 'twer as as proper 3 blacke as healthfull as 'twere as pure as hell as quickas a as her as actively siluer blacke as death as againsť as his as reason as how as sharpe as all as hush as sinewes as an blacke as hell as any 18 as i as sinnes as ice as snow as are as if blacke as his as bad 6 as so as before as in as some as swift as by as it as th'art as chast as just as kill 13 as the as non-zeros as kinde as therein checking as leuell as they as as liue as this common as thou as cunning as loue as low as thus as damn as lying as thy as day as death as mad as to as made as deepe as vnuallued as dicers as make as doth as vulcans as many 2 as may as as easie as england watchmen as meditation as waxe as ere as false as mine as we as mortall as well as farre as white as most as fits as much 2 as will as foule as fresh as my as as néedfull as from womans 2 as would as of as gaming as girdle 12 as you as oft as had 2 2 2 as one as your as hamlet as yours as our as patient as hardy

107 non-zeros

Hamlet from NLTK, V=4812

http://people.cs.umass.edu/~brenocon/inlp2014/lectures/04-hamlet_ngrams.txt



Thursday, September 11, 14

Hamlet from NLTK, V=4812

http://people.cs.umass.edu/~brenocon/inlp2014/lectures/04-hamlet_ngrams.txt



Convention in Collins/Knight: translating from **f** into **e**



Convention today (sorry!): translating from **e** into **f**



Convention today (sorry!): translating from **e** into **f**



Convention today (sorry!): translating from **e** into **f**



P(f | e) = P(e | f) P(f) / P(e)

Convention today (sorry!): translating from **e** into **f**



Convention today (sorry!): translating from **e** into **f**



Data to learn statistical MT

- Statistical machine translation is almost entirely based on *parallel corpora*, a.k.a. *bitexts*.
- Training data: human-translated sentence pairs
 (e, f)

How do we translate a word? Look it up in the dictionary

Haus : house, home, shell, household

- Multiple translations
 - Different word senses, different registers, different inflections (?)
 - *house*, *home* are common
 - *shell* is specialized (the Haus of a snail is a shell)

How common is each?

Translation	Count
house	5000
home	2000
shell	100
household	80

MLE

$$\hat{p}_{\mathrm{MLE}}(e \mid \mathrm{Haus}) = \begin{cases} 0.696 & \text{if } e = \mathrm{house} \\ 0.279 & \text{if } e = \mathrm{home} \\ 0.014 & \text{if } e = \mathrm{shell} \\ 0.011 & \text{if } e = \mathrm{household} \\ 0 & \text{otherwise} \end{cases}$$

- Goal: a model $p(\mathbf{e} \mid \mathbf{f}, m)$
- where e and f are complete English and Foreign sentences

 $\mathbf{e} = \langle e_1, e_2, \dots, e_m \rangle$ $\mathbf{f} = \langle f_1, f_2, \dots, f_n \rangle$

- Goal: a model $p(\mathbf{e} \mid \mathbf{f}, m)$
- where e and f are complete English and Foreign sentences
- Lexical translation makes the following **assumptions**:
 - Each word in e_i in e is generated from exactly one word in f
 - Thus, we have an alignment a_i that indicates which word e_i "came from", specifically it came from f_{a_i} .
 - Given the alignments **a**, translation decisions are conditionally independent of each other and depend *only* on the aligned source word f_{a_i} .

Lexical Translation $\mathbf{e} = \langle e_1, e_2, \dots e_m \rangle$ $\mathbf{f} = \langle f_1, f_2, \dots f_n \rangle$ $\mathbf{a} = \langle a_1, a_2, \dots a_m \rangle$ each $a_i \in \{0, 1, \dots, n\}$

Chain rule

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in \{0, 1, \dots, n\}^m} p(\mathbf{a} \mid \mathbf{f}, m) \times p(\mathbf{e} \mid \mathbf{a}, \mathbf{f}, m)$$

Modeling assumptions

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in \{0, 1, \dots, n\}^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

[Alignment] x [Translation | Alignment]

m

 $p(e_i \mid f_{a_i})$

 $p(e_i \mid f_{a_i})$ p(house | Haus)



 $\mathbf{e} = \langle e_1, e_2, \dots e_m \rangle \qquad \mathbf{f} = \langle f_1, f_2, \dots f_n \rangle$ $\mathbf{a} = \langle a_1, a_2, \dots a_m \rangle \quad \text{each } a_i \in \{0, 1, \dots, n\}$

Modeling assumptions

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in \{0, 1, \dots, n\}^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^m p(e_i \mid f_{a_i})$$

[Alignment] x [Translation | Alignment]

m

Alignment

$p(\mathbf{a} \mid \mathbf{f}, m)$

Most of the action for the first 10 years of MT was here. Words weren't the problem, word *order* was hard.

Alignment

 Alignments can be visualized in by drawing links between two sentences, and they are represented as vectors of positions:



$$\mathbf{a} = (1, 2, 3, 4)^{+}$$

Reordering

Words may be reordered during translation.



$$\mathbf{a} = (3, 4, 2, 1)^{\top}$$

Word Dropping

A source word may not be translated at all



$$\mathbf{a} = (2, 3, 4)^{\top}$$

Word Insertion

Words may be inserted during translation
 English just does not have an equivalent

But it must be explained - we typically assume every source sentence contains a NULL token



One-to-many Translation

 A source word may translate into more than one target word



$$\mathbf{a} = (1, 2, 3, 4, 4)^{\top}$$

Many-to-one Translation

 More than one source word may not translate as a unit in lexical translation



Many-to-one Translation

 More than one source word may not translate as a unit in lexical translation



IBM Model I

- Simplest possible lexical translation model
- Additional assumptions
 - The *m* alignment decisions are independent
 - The alignment distribution for each a_i is uniform over all source words and NULL

for each $i \in [1, 2, ..., m]$ $a_i \sim \text{Uniform}(0, 1, 2, ..., n)$ $e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$



$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m}$$

IBM Model I

for each $i \in [1, 2, \dots, m]$ $a_i \sim \text{Uniform}(0, 1, 2, \dots, n)$ $e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{1+n}$$
IBM Model I

for each $i \in [1, 2, ..., m]$ $a_i \sim \text{Uniform}(0, 1, 2, ..., n)$ $e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{1+n} p(e_i \mid f_{a_i})$$

IBM Model I

for each $i \in [1, 2, ..., m]$ $a_i \sim \text{Uniform}(0, 1, 2, ..., n)$ $e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{1+n} p(e_i \mid f_{a_i})$$

IBM Model I

for each
$$i \in [1, 2, ..., m]$$

 $a_i \sim \text{Uniform}(0, 1, 2, ..., n)$
 $e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$

$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(e_i, a_i \mid \mathbf{f}, m) = \frac{1}{1+n} p(e_i \mid f_{a_i})$$
$$p(\mathbf{e}, \mathbf{a} \mid \mathbf{f}, m) = \prod_{i=1}^{m} p(e_i, a_i \mid \mathbf{f}, m)$$

Thursday, September 11, 14

$$p(e_i, a_i \mid \mathbf{f}, m) = \frac{1}{1+n} p(e_i \mid f_{a_i})$$
$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(e_i, a_i \mid \mathbf{f}, m) = \frac{1}{1+n} p(e_i \mid f_{a_i})$$
$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(a, b, c, d) = p(a)p(b)p(c)p(d)$$

$$p(e_i, a_i \mid \mathbf{f}, m) = \frac{1}{1+n} p(e_i \mid f_{a_i})$$
$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(e_i, a_i \mid \mathbf{f}, m) = \frac{1}{1+n} p(e_i \mid f_{a_i})$$
$$p(e_i \mid \mathbf{f}, m) = \sum_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(\mathbf{e} \mid \mathbf{f}, m) = \prod_{i=1}^{m} p(e_i \mid \mathbf{f}, m)$$

$$p(e_i, a_i \mid \mathbf{f}, m) = \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(e_i \mid \mathbf{f}, m) = \sum_{\substack{a_i = 0 \\ m}}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(\mathbf{e} \mid \mathbf{f}, m) = \prod_{\substack{i=1}}^m p(e_i \mid \mathbf{f}, m)$$

$$p(e_i, a_i \mid \mathbf{f}, m) = \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(e_i \mid \mathbf{f}, m) = \sum_{\substack{a_i=0 \ m}}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$

$$p(\mathbf{e} \mid \mathbf{f}, m) = \prod_{\substack{i=1 \ m}}^m p(e_i \mid \mathbf{f}, m)$$

$$= \prod_{\substack{i=1 \ a_i=0}}^m \sum_{\substack{a_i=0}}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$





Start with a foreign sentence and a target length.

01234NULLdasHausistklein

1 2 3 4

















01234NULLdasHausistklein

1 2 3 4

















What is this good for?

- I. Evaluate translation quality
- 2. Find the best alignment
 - The IBM models are still used for this, though not as translation models

 $\mathbf{a}^* = \arg \max_{\mathbf{a} \in [0,1,\dots,n]^m} p(\mathbf{a} \mid \mathbf{e}, \mathbf{f})$ $= \arg \max_{\mathbf{a} \in [0,1,\dots,n]^m} \frac{p(\mathbf{e}, \mathbf{a} \mid \mathbf{f})}{\sum_{\mathbf{a}'} p(\mathbf{e}, \mathbf{a}' \mid \mathbf{f})}$ $= \arg \max_{\mathbf{a} \in [0,1,\dots,n]^m} p(\mathbf{e}, \mathbf{a} \mid \mathbf{f})$ $= p(\mathbf{a} \mid \mathbf{f}) p(\mathbf{e} \mid \mathbf{a}, \mathbf{f})$

$$a_i^* = \arg \max_{a_i=0}^n \frac{1}{1+n} p(e_i \mid f_{a_i})$$
$$= \arg \max_{a_i=0}^n p(e_i \mid f_{a_i})$$

01234NULLdasHausistklein

the home is little 1 2 3 4






















































Learning Lexical Translation Models

- How do we learn the parameters $p(e \mid f)$
- "Chicken and egg" problem
 - If we had the alignments, we could estimate the parameters (MLE)
 - If we had parameters, we could find the most likely alignments



How to learn?

- Training data: tons of (e,f) pairs. Want to learn theta: lexical translation probabilities.
- If knew the alignments **a**, MLE would be easy!

$$\max_{\theta} \sum_{(\mathbf{e}, \mathbf{f})} \log p(\mathbf{e} \mid \mathbf{a}, \mathbf{f}; \theta)$$

 But **a** is a latent variable. The marginal loglikelihood needs to sum-out the alignments. No closed form.

$$\max_{\theta} \sum_{(\mathbf{e}, \mathbf{f})} \log p(\mathbf{e} \mid \mathbf{f}; \theta)$$
$$= \sum_{(\mathbf{e}, \mathbf{f})} \log \left(\sum_{\mathbf{a}} p(\mathbf{a} \mid \mathbf{f}; \theta) \times p(\mathbf{e} \mid \mathbf{a}, \mathbf{f}; \theta) \right)$$

EM Algorithm

- pick some random (or uniform) parameters
- Repeat until you get bored (~ 5 iterations for lexical translation models)

• using your current parameters, compute "expected" alignments for every target word token in the training data $p(a_i \mid \mathbf{e}, \mathbf{f})$ (on board)

- keep track of the expected number of times f translates into e throughout the whole corpus
- keep track of the expected number of times that f is used as the source of any translation
- use these expected counts as if they were "real" counts in the standard MLE equation

... la maison ... la maison blue ... la fleur ...

- .. the house ... the blue house ... the flower ...
- Initial step: all alignments equally likely
- Model learns that, e.g., la is often aligned with the



- After one iteration
- Alignments, e.g., between la and the are more likely



- After another iteration
- It becomes apparent that alignments, e.g., between fleur and flower are more likely (pigeon hole principle)



- Convergence
- Inherent hidden structure revealed by EM



Thursday, September 11, 14



Now what?

- Decoding: find the best translation for a sentence
- Alignments: find the best alignment
- Evaluation?
- Other approaches?