

Lecture 10-12: Language Models

USC VSoE CSCI 544: Applied Natural Language Processing

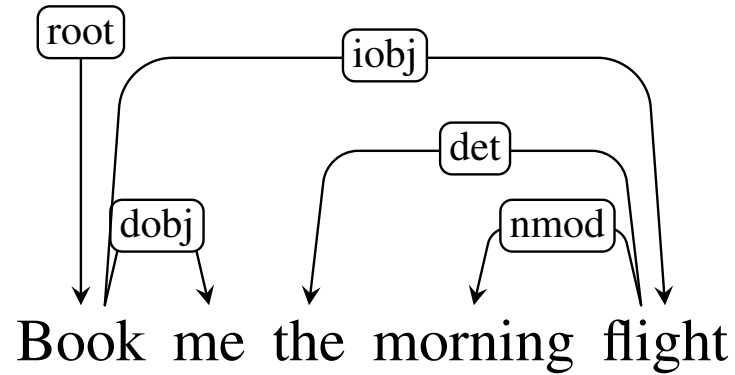
Jonathan May -- 梅約納

September 27-October 4, 2017

Quiz 1

- Which of the following CFG rules are in CNF (assume capital letters signify nonterminals, lowercase letters signify terminals)?
 - $S \rightarrow NP VP$
 - $NP \rightarrow NP PP$
 - $NP \rightarrow DT JJ NN$
 - $NP \rightarrow \text{the boy}$
 - $NP \rightarrow PP$
 - $JJ \rightarrow \text{red}$
 - $NP \rightarrow \text{the NP}$

Quiz 2



Step	Stack	Word List	Action	Relation Added
0	[root]	[book, me, the, morning, flight]	SHIFT	book $\xrightarrow{\text{dobj}}$ me
1	[root, book]	[me, the, morning, flight]	SHIFT	
2	[root, book, me]	[the, morning, flight]	RArc-dobj	
3	[root, book]	[the, morning, flight]	SHIFT	
4	[root, book, the]	[morning, flight]	SHIFT	
5	[root, book, the, morning]	[flight]	SHIFT	
6	[root, book, the, morning, flight]	[]		

- What is the next step?

- shift
- reduce
- LArc
- RArc
- LArc-det

- RArc-det
- RArc-nmod
- LArc-nmod

Please turn your homework ...

- What word comes next?
 - in
 - over
 - into
 - the
 - refrigerator
- What are the probabilities of each of these?
- And why should we care?

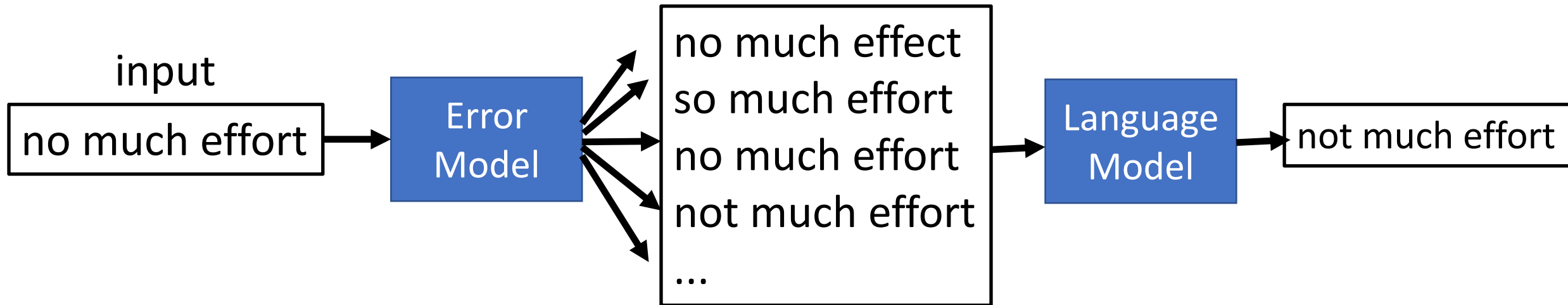
Probability of a sentence; $P(s)$

- How likely is it to occur?
- Colloquially, how likely is any speaker of language X to utter s?
 - $P(\text{the cat slept peacefully}) > P(\text{slept the peacefully cat})$
 - $P(\text{she studies morphosyntax}) > P(\text{she studies more faux syntax})$

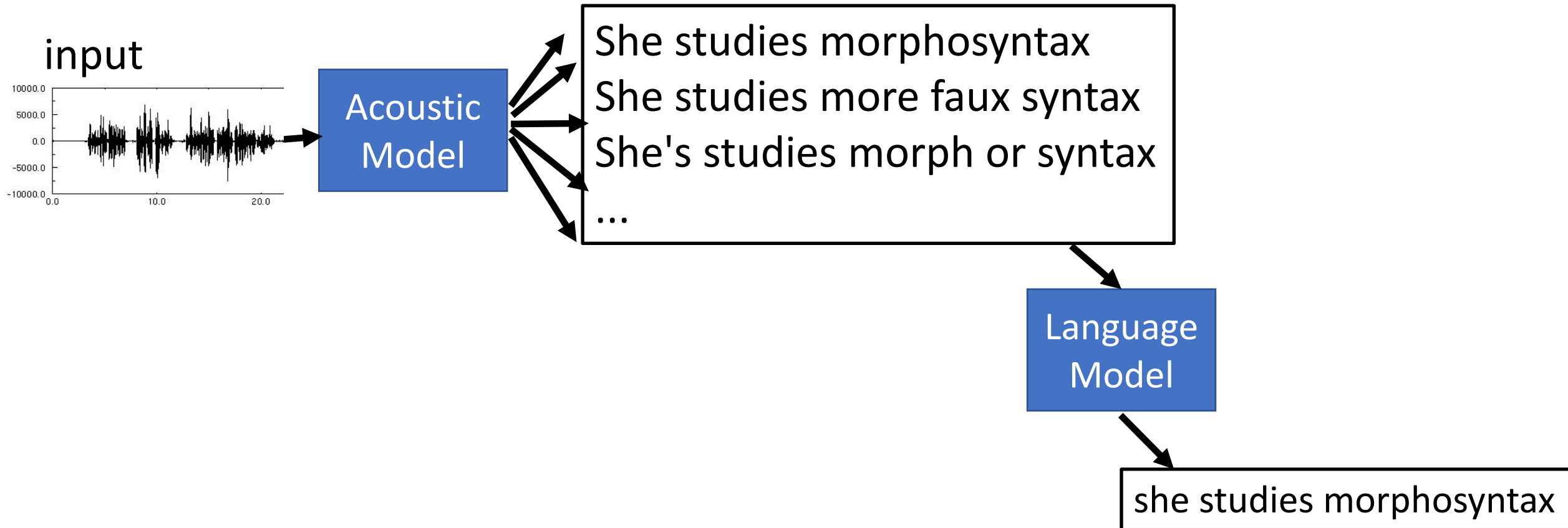
Language Models in NLP

- It's very difficult to know the true $P(s)$ for an arbitrary sequence of words
- But we can define a language model that will give us good approximations
- Language models (LMs) are very useful whenever we are generating output
 - Machine Translation
 - Spelling Correction
 - Summarization
 - Speech Recognition
- There are some very good, easy-to-use toolkits for building and using LMs
 - SRILM: around since the 90s. not advised with >300m tokens
 - KenLM: preferred choice; great scalability in memory and time with even billions of tokens
 - NLTK has some LM training/using support (ok for prototyping)

Use of Language Models: Spelling Correction



Use of Language Models: Speech Recognition



Use of Language Models: Machine Translation

input

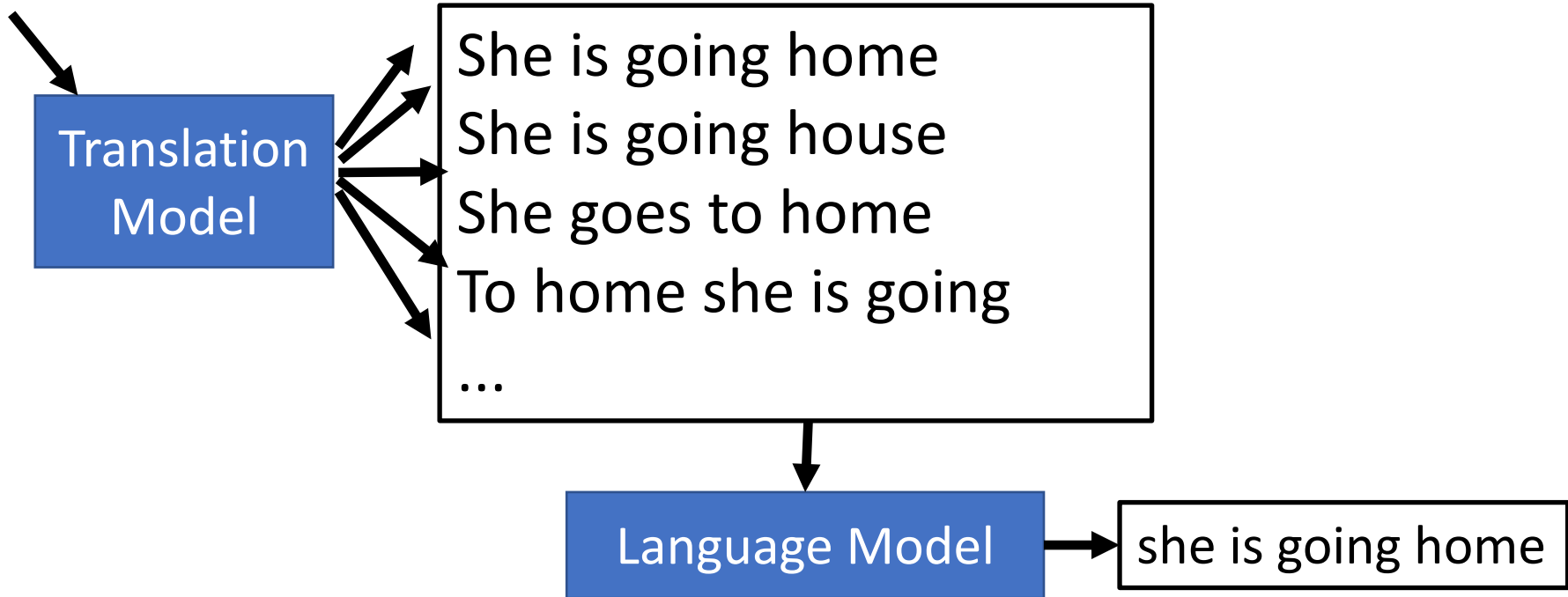
ella se va a casa

Translation
Model

She is going home
She is going house
She goes to home
To home she is going
...

Language Model

she is going home



LMs for Prediction

- LMs can be used to predict what a human will do next, rather than correct a possibly faulty model
- Example: predictive text correction/completion on your phone
 - Keyboard is tiny, so it's easy to touch a spot slightly off from the letter you intend
 - Correct these errors as you go and also provide possible completions

i n e f f i c i e n t

- In this case, LM may be defined over sequences of *characters* instead of (or in addition to) sequences of words

Noisy Channel Model



But How To Estimate These Probabilities?

- We want to know the probability of word sequence $\mathbf{w} = w_1, w_2, \dots, w_n$ occurring in English
- Assume we have some training data: large corpus of general English text
- We use this data to estimate the probability of \mathbf{w} (even if we never see it in the corpus)

Bit of Notation: Random Variables

- Random Variable = variable that represents all the possible events in some partition of Ω
- So if X is a coin flip I can say $P(X=\text{heads})$ to mean probability the flip comes up heads or just $P(X)$ to mean the probability table for the events (heads, tails)
- We can treat Random Variables like events
 - Flip coin A. Independently, flip coin B
 - $P(A=\text{heads} \mid B=\text{tails}) = P(A=\text{heads})$; $A=\text{heads}$ and $B=\text{tails}$ are independent events
 - for all x in {heads, tails}, for all y in {heads, tails}, $P(A=x \mid B=y) = P(A=x)$; A and B are independent random variables
- We've been using Random Variables all along, actually, but have been a bit sloppy
- (I probably should have discussed these back in Lecture 4)

Probability of a word sequence

- $P(\mathbf{w}) = P(w_1, w_2, w_3, w_4, \dots, w_n)$
- e.g. $P(\mathbf{w} = \text{the cat slept quietly}) = P(w_1=\text{the}, w_2=\text{cat}, w_3=\text{slept}, w_4=\text{quietly})$
- We'll often abuse notation when talking about specific events and context is clear, e.g. $P(\text{the cat slept quietly})$

Maximum Likelihood Estimation?

- Recall Maximum Likelihood Estimations (MLE) for our HMM POS tagger
 - AKA "Count and divide"
- So get a corpus of N sentences
 - $P_{MLE}(\mathbf{w} = \text{the cat slept quietly}) = C(\text{the cat slept quietly})/N$
- But consider these sentences:
 - the long-winded peripatetic beast munched contentedly on mushrooms
 - parsimonius caught the of about for syntax
- Neither is in a corpus (I just made them up), so $P_{MLE}=0$ for both
 - But one is meaningful and grammatical and the other isn't!

Sparse Data and MLE

- If something doesn't occur, MLE thinks it can't occur
- No matter how much data you get, you won't have enough observations to model all events well with MLE
- We need to make some assumptions so that we can provide a reasonable probability for grammatical sentences, even if we haven't seen them

Independence (Markov) Assumption

- Recall, $P(w_1, w_2, \dots, w_n) = P(w_n | w_1, w_2, \dots, w_{n-1})P(w_{n-1} | w_1, w_2, \dots, w_{n-2})\dots$
 - $\prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$
- Still too sparse (nothing changed)
 - if we want $P(\text{I spent three years before the mast})$
 - we still need $P(\text{mast} | \text{I spent three years before the})$
- Remember definition of independence; A and B are independent if $P(A) = P(A|B)$

Independence (Markov) Assumption

- Make an n-gram independence assumption: probability of a word only depends on a fixed number of previous words (history)
 - **trigram model:** $P(w_i | w_1, \dots, w_{i-1}) \approx P(w_i | w_{i-2}, w_{i-1})$
 - **bigram model:** $P(w_i | w_1, \dots, w_{i-1}) \approx P(w_i | w_{i-1})$
 - **unigram model:** $P(w_i | w_1, \dots, w_{i-1}) \approx P(w_i)$
- I.e. a trigram model says
 - $P(\text{mast} \mid \text{I spent three years before the}) \approx P(\text{mast} \mid \text{before the})$
- It also assumes all these are equal:
 - $P(\text{mast} \mid \text{I spent three years before the})$
 - $P(\text{mast} \mid \text{I went home before the})$
 - $P(\text{mast} \mid \text{I saw the sail before the})$
because all are estimated as $P(\text{mast} \mid \text{before the})$
- Not always a good assumption! But it does reduce the sparse data problem

Estimating Trigram Conditional Probabilities

- $P_{MLE}(\text{mast} \mid \text{before the}) = \text{Count}(\text{before the mast}) / \text{Count}(\text{before the})$
- In general, for any trigram, we have
 - $P_{MLE}(w_i \mid w_{i-2}, w_{i-1}) = \frac{\text{Count}(w_{i-2}, w_{i-1}, w_i)}{\text{Count}(w_{i-2}, w_{i-1})}$

Example from *Moby Dick* corpus

- $C(\text{before, the}) = 25$; $C(\text{before, the, mast}) = 4$
- $C(\text{before, the, mast}) / C(\text{before, the}) = 0.16$
- mast is the most common word to come after "before the" (wind is second most common)
- $P_{\text{MLE}}(\text{mast}) = 56/110927 = .0005$ and $P_{\text{MLE}}(\text{mast}|\text{the}) = .003$
- Seeing "before the" vastly increases the probability of seeing "mast" next

Trigram model summary

- To estimate $P(\mathbf{w})$, use chain rule and make an independence assumption
 - $P(w_1, \dots, w_n) = \prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$
 - $\approx P(w_1)P(w_2 | w_1) \prod_{i=3}^n P(w_i | w_{i-2}, w_{i-1})$
- Then estimate each trigram prob from data (here, using MLE)
 - $P_{MLE}(w_i | w_{i-2}, w_{i-1}) = \frac{\text{Count}(w_{i-2}, w_{i-1}, w_i)}{\text{Count}(w_{i-2}, w_{i-1})}$

Midterm Info

Some of you will
NOT BE IN SAL 101
YOU WILL BE
IN MHP (Mudd) 101



Who is Where?

- If your last* name is from Aggarwal to Patel, you are in SAL 101
- If it is from Pathapi to Zhu, you are in MHP 101
- * "Last" is a very western term. I have an alphabetization based on what USC has your surname listed as.
- To avoid confusion, please check the course website, which has a very prominent link to a roster with firstname, last name, partially masked email address, and assignment
- If you are still unsure of where to go, contact us on piazza ASAP

Midterm Logistics

- Length: 1 hour, 40 minutes (you will need time to set up and hand in exam)
- Date: Friday, October 6, 8:00 AM (Please arrive promptly!)
- Please Bring:
 - Pencils/pens/erasers as needed
 - one 8.5x11 inch (or A4) sheet of paper with notes on both sides (optional)
 - NO OTHER NOTES
 - NO ELECTRONIC RESOURCES
 - NO BOOKS
- We will provide extra paper for scratch work
- Sit only at seats with exams on them. Fill up all available space.

What's On The Exam?

- Fair game
 - Anything on the slides
 - Anything in the required reading
 - Anything in the homeworks
- But
 - We're not trying to trick you
 - We're not trying to make this impossible
 - If you understand the lectures well, you should be ok

What's On The Exam

- Major Topics
 - Levels of Linguistic Knowledge (L1)
 - Corpora, Regex, Basic text processing (L2)
 - Morphology, Finite State Automata and Transducers (L3)
 - Probability Theory (L4)
 - Naive Bayes, Features, Perceptron, Logistic Regression (L5)
 - POS Tagging and HMM tagger, Viterbi Decoding (L6)
 - Constituency Syntax Trees, Context-Free Grammars, CKY, CNF, Smoothing, Interpolating, Beam Decoding (L7-8)
 - Dependency Syntax Trees, Arc-Standard and Arc-Eager Dependency Parsing (L9-10)
 - Ngram Language Models, Smoothing, Backoff, Interpolation, alternate language models (L10-11) (Probably no neural; depends on how far we get)
- It won't all be on there because there isn't enough time
 - But there is plenty of room on the final

Practical details (I)

- Trigram model assumes two-word history
- But consider these sentences:

w_1	w_2	w_3	w_4
he	saw	the	yellow
feeds	the	cats	daily

- What's wrong?
 - a sentence shouldn't end with 'yellow'
 - a sentence shouldn't begin with 'feeds'
- Does the model capture these problems?

Beginning / end of sequence

- To capture behavior at beginning/end of sequences, we can augment the input:

w_{-1}	w_0	w_1	w_2	w_3	w_4	w_5
<s>	<s>	he	saw	the	yellow	</s>
<s>	<s>	feeds	the	cats	daily	</s>

- That is, assume $w_{-1}=w_0=<s>$ and $w_{n+1}=</s>$ so:
 - $P(\mathbf{w}) = \prod_{i=1}^{n+1} P(w_i | w_{i-2}, w_{i-1})$
- Now $P(</s> | \text{the, yellow})$ is low, indicating this is not a good sentence
- $P(\text{feeds} | <s>, <s>)$ should also be low

Beginning/end of sequence

- Alternatively, we could model all sentences as one (very long) sequence, including punctuation
 - two cats live in sam 's barn . sam feeds the cats daily . yesterday , he saw the yellow cat catch a mouse . [...]
- Now, trigram probabilities like $P(. \mid \text{cats daily})$ and $P(, \mid . \text{ yesterday})$ tell us about behavior at sentence edges
- Here, all tokens are lowercased. What are the pros/cons of not doing that?

Practical details (II)

- Word probabilities are typically very small.
- Multiplying lots of small probabilities quickly gets so tiny we can't represent the numbers accurately, even with double precision floating point.
- So in practice, we typically use log probabilities (usually base-e)
 - Since probabilities range from 0 to 1, log probs range from $-\infty$ to 0
 - Instead of multiplying probabilities, we add log probs
 - Often, negative log probs are used instead; these are often called "costs"; lower cost = higher prob
- Recall: we saw this with bigram HMM for POS tagging

Interim Summary: N-gram probabilities

- "Probability of a sentence": how likely is it to occur in natural language?
- We can never know the true probability, but we may be able to estimate it from corpus data.
- N-gram models are one way to do this:
 - To alleviate sparse data, assume word probs depend only on short history
 - Tradeoff: longer histories may capture more, but are also sparser
 - So far, we estimated N-gram probabilities using MLE

Interim Summary: Language Models

- Language Models tell us $P(\mathbf{w}) = P(w_1, \dots, w_n)$: How likely is this sequence of words to occur?
 - Roughly: Is this sequence of words a "good" one in my language?
- LMs are used as a component in applications such as speech recognition, machine translation, and predictive text completion
- To reduce sparse data, N-gram LMs assume words depend only on a fixed-length history, even though we know this isn't true
- Next:
 - How to evaluate a language model
 - Weaknesses of MLE and how to address them (more sparsity)

Quiz 3

- Which of These Statements are True?
 - Naive Bayes is a probabilistic model
 - Perceptron is a generative model
 - Logistic regression has a closed-form solution
 - Logistic regression is a discriminative model

Two Types of Evaluation in NLP

- **Extrinsic:** measure performance on a downstream application
 - For LM, plug it into a machine translation/ASR/etc system
 - The most reliable and useful evaluation: We don't use LMs absent other technology
 - But can be time-consuming
 - And of course we still need an evaluation measure for the downstream system
- **Intrinsic:** design a measure that is inherent to the current task
 - much quicker/easier during development cycle
 - not always easy to figure out what the right measure is. Ideally, it's one that correlates with extrinsic measures
 - Extra-hard for LMs

Intrinsically Evaluating a Language Model

- For parsing, tagging, sentiment, etc. it was fairly clear how to evaluate: Hold a set of labeled data out and see how often your model gets it right
- For LM, it's not quite so clear
 - Given a corpus of sentences and non-sentences, see how often the LM thinks you have a sentence?
 - Not a very realistic evaluation of how an LM is used
 - Often we are deciding between not-that-grammatical outputs
- Ideally we want a regression evaluation
 - Given a sentence, how close is the model probability to the true probability
 - But we don't know the true probability of a sentence!

Idea: Model should give high probability to an unseen corpus

- Assume that you have a proper probability model, i.e. for all sentences S in the language L , $\sum_{S \in L} P(S) = 1$
- Then take a held-out test corpus T consisting of sentences in the language you care about
- $\prod_{t \in T} P(t)$ should be as high as possible; model should think each sentence is a good one
- Let's be explicit about evaluating each word in each sentence
 - $\prod_{t \in T} \prod_{w \in t} P(w)$
- Collapse all these words into one big 'sentence' N :
 - $\prod_{w \in N} P(w)$

Resolving Some Problems

- $\prod_{w \in N} P(w)$ is going to result in underflow. Ok, let's use logs again!
- Also we tend to like positive sums.
 - $-\sum_{w \in N} \log(P(w))$
- This can be tough to compare against corpora of different length (or sentences of different length), so normalize by the number of words:
 - $\frac{-\sum_{w \in N} \log(P(w))}{|N|}$ is called the cross-entropy of the data according to the model
- When comparing models, differences between these numbers tend to be pretty small, so we exponentiate
 - $2^{\frac{-\sum_{w \in N} \log(P(w))}{|N|}}$ is called the perplexity of the data

Example

- Three word sentence with probabilities $\frac{1}{4}$, $\frac{1}{2}$, $\frac{1}{4}$
 - $\frac{1}{4} * \frac{1}{2} * \frac{1}{4} = .03125$
 - cross-entropy: $-(\log(1/4) + \log(1/2) + \log(1/4))/3 = 5/3$; $2^{5/3} \approx 3.17$
- Six word sentence with probabilities $\frac{1}{4}$, $\frac{1}{2}$, $\frac{1}{4}$, $\frac{1}{4}$, $\frac{1}{2}$, $\frac{1}{4}$
 - $\frac{1}{4} * \frac{1}{2} * \frac{1}{4} * \frac{1}{4} * \frac{1}{2} * \frac{1}{4} = .00097$
 - cross-entropy: $-(\log(1/4) + \log(1/2) + \log(1/4) + \log(1/4) + \log(1/2) + \log(1/4))/6 = 10/6$; $2^{10/6} \approx 3.17$
- If you overfit your training corpus so that $P(\text{train}) = 1$, then Perplexity on train is 0
- But Perplexity on test (which doesn't overlap with train) will be infinite

Intrinsic Evaluation Big Picture

- Lower Perplexity is better
- Roughly = number of bits needed to communicate information about a word
 - The terms 'cross-entropy' and 'perplexity' come out of information theory; it's in the reading if you're interested but we won't dwell on it
- In principle you could compare on different test sets
- In practice, domains shift. To know which of two LMs is better, train on common training sets, test on common test sets

Sparse data, again

- Suppose we build a trigram model from Moby Dick and evaluate the sentence "I spent three years before the mast"
- "I spent three" never occurs in training, so $P_{MLE}(\text{three} | \text{I spent}) = 0$
- so cross-entropy is infinite
- This is basically right; our model says "I spent three" should never occur so when it does our model is infinitely surprised!
- Even with a unigram model we run into words we never saw, so we need better ways to estimate probabilities from sparse data

Add-1 (Laplace) and Add- α (Lidstone) Smoothing Again

- Pretend we saw everything 1 (α) more times than we did before
- $P_{+1}(w_i | w_{i-2}, w_{i-1}) = (C(w_{i-2}, w_{i-1}, w_i) + 1) / (C(w_{i-2}, w_{i-1}) + |V|)$ where $|V|$ is the size of the vocabulary
- $P_{+\alpha}(w_i | w_{i-2}, w_{i-1}) = (C(w_{i-2}, w_{i-1}, w_i) + \alpha) / (C(w_{i-2}, w_{i-1}) + \alpha |V|)$

Dealing with unknown vocabulary

- Can we also add a new 'OOV' token as was done in HMM emission table?
- It gets kind of complicated...
- $P_{+\alpha}(w_i | w_{i-2}, w_{i-1}) = (C(w_{i-2}, w_{i-1}, w_i) + \alpha) / (C(w_{i-2}, w_{i-1}) + \alpha |V| + 1)$
- $P_{+\alpha}(w_i = \text{OOV} | w_{i-2}, w_{i-1}) = \alpha / (\alpha |V| + 1)$
- But then we also have to deal with, e.g., $P_{+\alpha}(w_i | w_{i-2}, w_{i-1} = \text{OOV})$
- Better solution: replace low-count words in corpus with "OOV"
- Intuition: 1-count is basically the same as 0-count

Remaining Problem

- In a training corpus, suppose we see *Scottish beer* but neither of
 - *Scottish beer drinkers*
 - *Scottish beer eaters*
- If we build a smoothed trigram model (with any kind of smoothing), which example has higher probability?
 - Both the same! Unknown events are treated equally by smoothing!

Remaining Problem

- Previous smoothing methods assign equal probability to unseen events
- Better: use information from lower-order N-grams (shorter histories)
 - beer drinkers
 - beer eaters
- Two ways: backoff and interpolation

Backoff

- Idea: Trust the highest order language model that contains your N-gram
- $P_{BO}(z|x\ y) =$
 $(1 - \alpha_{xy})P(z|x\ y)$ if $\text{count}(x\ y) > 0$
 $\alpha_{xy} P_{BO}(z | y)$ else
- where α_{xy} is an interpolation parameter

Simple Interpolation

- Idea: Trust different amounts of context differently
- $P_{SI}(z|x y) =$
 $\lambda_3 P(z|x y) +$
 $\lambda_2 P(z|y) +$
 $\lambda_1 P(z) +$
 λ_0
- where $\lambda_0 + \lambda_1 + \lambda_2 + \lambda_3 = 1$, all ≥ 0
- We did something similar in lexicalized constituency parsing

Better Interpolation

- Idea: As in backoff, particular contexts matter
- $P_{SI}(z|x y) =$
 $\lambda_{xy} P(z|x y) +$
 $\lambda_{y/xy} P(z|y) +$
 $\lambda_{1/xy} P(z) +$
 $\lambda_{0/xy}$
- where, for each xy , $\lambda_{0/xy} + \lambda_{1/xy} + \lambda_{y/xy} + \lambda_{xy} = 1$, all ≥ 0
- Best not to actually have a different set for each unique context; can group by context count

State-of-the-art Smoothing

- There is lots and lots of work done on smoothing and lots of variants
- See Chen and Goodman (optional reading); it's actually quite comprehensive, though mathy
- Best today is Modified Kneser-Ney
 - replace MLE with estimates based on count of unique histories
 - 4 interpolation lambdas based on ngram counts
- For very large data, Google's Stupid Backoff
 - Really fast to calculate
 - Doesn't give proper perplexities!
 - Works well in practice
- These are available in SRILM (K-N) and KenLM (both)