

Noisy Channel, N-grams & Smoothing

Lecture #6

Introduction to Natural Language Processing
CMPSCI 585, Spring 2004

University of Massachusetts Amherst



Andrew McCallum

Today's Main Points

- Discussion of PA#1.
 - Your experience, gained insights.
- Discussion of PA#2.
 - Review of some pointers, help.
- Application of the Noisy Channel Model
- Markov Models
 - including Markov property definition
- Smoothing
 - Laplace, Lidstone's, Held-out, Good-Turing.

PA#2 Help Accuracy Evaluation

Result of running classifier on a test set:

```
filename trueclass predclass p(predclass|doc)
filename trueclass predclass p(predclass|doc)
filename trueclass predclass p(predclass|doc)
...
```

	true spam	true ham
pred spam	TP	FP
pred ham	FN	TN

Accuracy = $(TP+TN) / (TP+TN+FP+FN)$
Precision = $TP / (TP+FP)$
Recall = $TP / (TP+FN)$
F1 = harmonic mean of Precision & Recall

PA#2 Help Precision-Recall Curve

Result of running classifier on a test set:

```
filename trueclass predclass p(predclass|doc)
filename trueclass predclass p(predclass|doc)
filename trueclass predclass p(predclass|doc)
...
```

	true spam	true ham
pred spam	TP	FP
pred ham	FN	TN

Accuracy = $(TP+TN) / (TP+TN+FP+FN)$
Precision = $TP / (TP+FP)$
Recall = $TP / (TP+FN)$
F1 = harmonic mean of Precision & Recall

PA#2 Help Accuracy-Coverage Curve

Result of running classifier on a test set:

```
filename trueclass predclass p(predclass|doc)
filename trueclass predclass p(predclass|doc)
filename trueclass predclass p(predclass|doc)
...
```

	true spam	true ham
pred spam	TP	FP
pred ham	FN	TN

Accuracy = $(TP+TN) / (TP+TN+FP+FN)$
Precision = $TP / (TP+FP)$
Recall = $TP / (TP+FN)$
F1 = harmonic mean of Precision & Recall

PA#2 Help Working with log-probabilities

$$p(c|d) \propto p(c) \prod_i p(w_i|c)$$

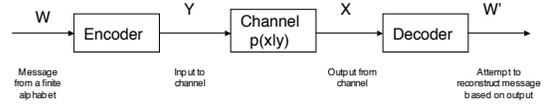
$$\log(p(c|d)) \propto \log(p(c)) + \sum_i \log(p(w_i|c))$$

- Getting back to $p(c|d)$
 - Subtract a constant to make all non-positive
 - $\exp()$

PA#2 Help

The importance of test/train splits

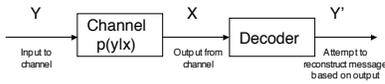
Noisy Channel Model



- Optimize Encoder for throughput and accuracy.
 - compression: remove all redundancy
 - accuracy: adding controlled redundancy
- Capacity: rate at which can transmit information with arbitrarily low probability of error in W'

$$C = \max_{p(Y)} I(X; Y)$$

Noisy Channel in NLP



$$Y' = \arg \max_y p(y|x) = \arg \max_y \frac{p(y)p(x|y)}{p(x)} = \arg \max_y p(y)p(x|y)$$

Application	Input	Output	$p(y)$	$p(x y)$
Machine Translation	L1 word sequences	L2 word sequences	$p(L1)$ in a language model	translation model
Optical Character Recognition (OCR)	actual text	text with OCR errors	prob of language text	model of OCR errors
Part of Speech (POS) tagging	POS tag sequences	English word sequence	prob of POS sequences	probability of word given tag
Speech Recognition	word sequences	acoustic speech signal	prob of word sequences	acoustic model
Document classification	class label	word sequence in document	class prior probability	$p(L1)$ from each class

Probabilistic Language Modeling

- Assigns probability $p(t)$ to a word sequence $t = w_1 w_2 w_3 w_4 w_5 w_6 \dots$
- Chain rule and joint/conditional probabilities for text t :

$$p(t) = p(w_1 \dots w_n) = p(w_1) \dots p(w_n | w_1, \dots, w_{n-1})$$

$$= \prod_{i=1}^n p(w_i | w_1 \dots w_{i-1})$$

where

$$p(w_k | w_1 \dots w_{k-1}) = \frac{p(w_1 \dots w_k)}{p(w_1 \dots w_{k-1})} \sim \frac{C(w_1 \dots w_k)}{C(w_1 \dots w_{k-1})}$$

The chain rule leads to a **history-based model**: we predict following things from past things.

n-gram models

the classic example of a statistical model of language

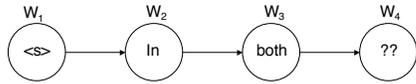
- Each word is predicted according to a conditional distribution based on limited context
- Conditional Probability Table (CPT): $p(X|“both”)$
 - $p(\text{of}|\text{both}) = 0.066$
 - $p(\text{to}|\text{both}) = 0.041$
 - $p(\text{in}|\text{both}) = 0.038$
- a.k.a. Markov (chain) models
 - sequences of random variables in which the future variable is determined by the present variable, but is independent of the way in which the present state arose from its predecessors

The Markov Property

- A stochastic process has the **Markov property** if the conditional probability distribution of future states of the process, given the present state, depends only upon the current state, and conditionally independent of the past states (the *path* of the process) given the present state.
- A process with the Markov property is usually called a **Markov process**, and may be described as **Markovian**.

$$\Pr[X(t+h) = y | X(s) = x(s), s \leq t] = \Pr[X(t+h) = y | X(t) = x(t)], \quad \forall h > 0.$$

1-gram model



First-order Markov model, $P(w_i|w_{i-1})$

- Simplest linear graphical model
- Words are random variables, arrows are direct dependencies between them (CPTs)
- These simple engineering models have been amazingly successful.

n-th order Markov models

- First order Markov assumption = bigram

$$p(w_k|w_1...w_{k-1}) \simeq p(w_k|w_{k-1}) = \frac{p(w_{k-1}w_k)}{p(w_{k-1})}$$

- Similarly, n-th order Markov assumption

- Most commonly, trigram (2nd order)

$$p(w_k|w_1...w_{k-1}) \simeq p(w_k|w_{k-2}w_{k-1}) = \frac{p(w_{k-2}w_{k-1}w_k)}{p(w_{k-2}w_{k-1})}$$

Andrei Andreyevich Markov



1856 - 1922

- Graduate of Saint Petersburg University (1878), where he began a professor in 1886.
- Mathematician, teacher political activist
 - In 1913, when the government celebrated the 300th anniversary of the House of Romanov family, Markov organized a counter-celebration of the 200th anniversary of Bernoulli's discovery of the Law of Large Numbers.
- Markov was also interested in poetry and he made studies of poetic style.

Markov's Model

- Took 20,000 characters from Pushkin's *Eugene Onegin* to see if it could be approximated by a simple chain of characters.

	vowel	consonant
vowel	0.128	0.872
consonant	0.663	0.337

Markov Approximations to English

- Zero-order approximation, $P(c)$
 - XFOML RXXXRJFFUJ ZLPWCFWKRJ
FFJEYVKCQSGHYD QPAAMKBZAACIBZLHJQD
- First-order approximation, $P(c|c)$
 - OCRO HLI RGWR NWIELWIS EU LL
NBNESBYA TH EEI ALHENHTTPA OOBTTVA
- Second-order approximation, $P(c|c,c)$
 - ON IE ANTSOUTINYNS ARE T INCTORE ST BE
S DEAMY ACHIN D ILONASIVE TUCOOWE AT
TEASONARE FUSO TIZIN ANDY TOBE SEACE
CTISBE

[From Shannon's original paper]

Markov Approximations to English (cont.)

- Third-order approximation, $P(w|w,w,w)$
 - IN NO IST LAT WHEY CRATICT FROURE BIRS
GROCID PONDENOME OF DEMONSTURES
OF THE REPTABIN IS REGOACTON OF CRE
- Markov Random Field with 1000 "features"
 - WAS REASER IN THERE TO WILL WAS BY
HOMES THING BE RELOVERATED THER
WHICH CONSISTS AT FORES ANDITING WITH
PROVERAL THE CHESTRAING FOR HAVE TO
INTRALLY OF QUT DIVERAL THIS OFFECT
INATEVER THIFER CONTRANDED STATER

[Della Pietra, Della Pietra & Lafferty, 1997]

Character n-grams for language identification

Word-based Approximations

- First-order approximation
 - representing and speedily is an good apt or come can different natural here he the a in came the to of to expert gray come to furnishes the line message had be
- Second-order approximation
 - the head and in frontal attack on an English writer that the character of this point is therefore another method for the letters that the time of who ever told the problem for an unexpected

Shannon's comment (1948): "It would be interesting if further approximations could be constructed, but the labor involved becomes enormous at the next stage."

n-gram models

- Core language model for the engineering task of better predicting the next word:
 - Speech recognition
 - OCR
 - Context-sensitive spelling correction
- It has only recently that improvements have been made for these tasks [Alshawi '96, Wu '97]
- But linguistically, they are appallingly simple and naïve.

Why might n-gram models not work?

- Relationships (say between subject and verb) can be arbitrarily distant and convoluted, as linguists love to point out:
 - The **man** I was watching through the window without pausing to look at what was happening down the street, and quite oblivious to the situation that was about to befall him confidently **strode** into the center of the road.

Why do they work?

- That kind of thing doesn't happen much
- Collins (1997)
 - 74% of dependencies (in the Penn Treebank, WSJ) are with an adjacent word (95% with one less than 5 words away), once one treats simple NPs as units

Evaluation of language models

- Best evaluation of probability model is task-based!
- As substitute for evaluating one component, standardly use corpus per-word cross entropy:

$$H(X, p) = -\frac{1}{n} \sum_{i=1}^n \log_2 p(w_i | w_1, \dots, w_{i-1})$$

- Or perplexity
 - units = average number of choices, scaled for uniform distr.
 - high = unpredictable

$$PP(X, p) = 2^{H(X, p)} = \left[\prod_{i=1}^n p(w_i | w_1, \dots, w_{i-1}) \right]^{-1/n}$$

Parameter Estimation

Maximum Likelihood Estimate

- Relative frequency
- Makes training data as probable as possible
- Overfits

$$p(w_2|w_1) = \frac{C(w_1, w_2)}{C(w_1)}$$

Limitations of the Maximum Likelihood Estimator

- Problem: often infinitely surprised when unseen word appears, $P(\text{unseen}) = 0$
 - Problem: this happens commonly
 - Probabilities of zero-count words are too low
 - Probabilities of nonzero-count words are too high
 - Estimates for high count words are fairly accurate
 - Estimates for low count words are unstable
 - We need “smoothing”

Sparsity

- How often does an every day word like “kick” occur in a million words of text?
 - “kick”: about 10 [depends vastly on genre, of course]
 - “wrist”: about 5
- Normally we want to know about something bigger than a single word, like how often you “kick a ball”, or how often the dative alternation “he kicked the baby a toy” occurs.
- How often can we expect that to occur in 1 million words?
- Almost never.
- “There’s no data like more data”
 - Must be of the right domain

Severity of the sparse data problem

count	2-grams	3-grams
1	8,045,024	53,737,350
2	2,065,469	9,229,958
3	970,434	3,654,791
>4	3,413,290	8,728,789
>0	14,494,217	75,349,888
possible	6.8×10^{10}	1.7×10^{16}

Vocab size 260,741 words, 365M words training

The Zero Problem

- Necessarily some zeros
 - trigram model: 1.7×10^{16} parameters
 - but only 2.6×10^8 words of training data
- How should we distribute some probability mass over all possibilities in the model
 - optimal situation: even the least frequent trigram would occur several times, in order to distinguish its probability versus other trigrams
 - optimal situation cannot happen, unfortunately (how much data would we need?)
- Two kinds of zeros: $p(w|h)=0$, or even $p(h)=0$

Laplace smoothing

$$p(w_2|w_1) = \frac{C(w_1, w_2) + 1}{C(w_1) + V}$$

- V is the vocabulary size (assume fixed, closed vocabulary)
- This is the Bayesian *maximum a posteriori* estimator you get by assuming a uniform prior on multinomials (a Dirichlet prior)

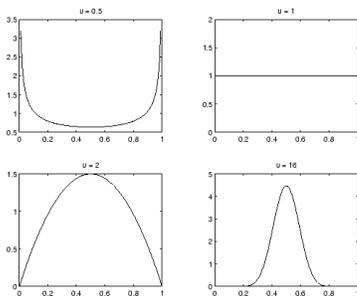
Dirichlet Distribution

- “Multinomial” is a die: a distribution over a finite alphabet of outcomes
- “Dirichlet” is a dice generator: a distribution over multinomials!

$$p(q) \sim \text{Dirichlet}(\alpha_1, \alpha_2, \dots, \alpha_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k q_k^{\alpha_k - 1}$$

- It is a conjugate prior for multinomials

Dirichlet Examples



For $a_1 = s_2 = a$

Laplace Smoothing

- Problem: gives too much probability mass to unseens
- Not good for large vocabulary, comparatively little data (NLP!)
- e.g. 10,000 word vocab, 1,000,000 words of training data, but “comes across” occurs 10 times. Of those, 8 times next word is “as”
 - $P_{MLE}(\text{as}|\text{comes across}) = 0.8$
 - $P_{Laplace}(\text{as}|\text{comes across}) = (8+1)/(10+10000) = 0.0009$
- Quick fix: Lidstone’s law (Mitchell’s 1997 “*m*-estimate”):

$$p(w_2|w_1) = \frac{C(w_1, w_2) + \lambda}{C(w_1) + \lambda V} \quad \text{for } \lambda < 1, \text{ e.g. } 1/2 \text{ or } 0.05$$

How much mass to allocate to unseens?

- For Laplace smoothing, in $P(.|\text{comes across})$, 10,000/10,010 of the prob mass is given to unseen events.
- How do we know that this is too much?

Absolute discounting

- Idea is that we want to discount counts of seen things a little, and reallocate this probability mass to unseens
- By subtracting a fixed count, probability estimates for commonly seen things are scarcely affected, while probabilities of rare things are greatly affected
- If the discount is around $d=0.75$, then seeing something once is not so different than not having seen it at all

$$p(w_2|w_1) = \frac{C(w_1, w_2) - \delta}{C(w_1)}, \text{ if } C(w_1, w_2) > 0$$

$$p(w_2|w_1) = \frac{\frac{(V - N_0)\delta}{N_0}}{C(w_1)}, \text{ otherwise}$$

Held Out Estimator

- How do you know how likely you are to see a new word type in the future (in a certain context)?
 - Examine some further text and find out (empirical held-out estimators = validation)
 - Divide data into two pots: training data, validation data
 - N = number of (non-unique) bigrams in training
 - N_r = number of unique bigrams with freq r in training
 - T_r = number of times that all bigrams appearing r times in the training data appeared in the validation data.

$$p_{ho}(w_1w_2) = \frac{T_r}{N_r N} \text{ where } r = C(w_1w_2)$$

T_r/N_r is an "improved estimate for r ".

Pots of data for estimating and testing models

- Major error: testing on your training data
- **Overfitting**: expect future events to be too much like the events on which it was trained, rather than allowing sufficiently for other possibilities.
- **Training data**
- **Validation data**
- **Testing data**
 - Development test data
 - Final test data
- Don't report results on just one test set, but average of many, and report variance... use a test of statistical significance

Held-out Estimator with Cross-validation

- Reshuffle the training data several times into pots of training and validation data
- Calculate $p_{ho}(w_1w_2)$ for each split, then average them.

$$p_{ho}(w_1w_2) = \frac{T_r^{01} + T_r^{10}}{(N_r^0 N_r^1) N} \text{ where } r = C(w_1w_2)$$

- Extreme case of cross-validation: **leave-one-out** cross validation
 - Train on $N-1$ of the words, validate on 1 word
 - How much probability mass would be reserved for the unseen words in this case?

Good-Turing Smoothing

- Derivation reflects leave-one-out estimation
 - For each word token in data, call it the validation set; remaining data is training set.
 - The validation-set word has count r in training set.
 - See how often any word has r counts in training set. (How many different words have count r ?)
 - This will happen every time word left out has $r+1$ counts in original data
 - So total count mass of r count words is assigned from mass of $r+1$ count words = $N_{r+1} \times (r+1)$
 - Apply to low counts; not needed (harmful!) for high count words

Good-Turing smoothing

- All words with same count get same probability (as before)
- Count mass of words with $r+1$ occurrences is assigned to words with r occurrences.
- r^* is corrected frequency estimates for a word occurring r times

$$r^* = (r+1) \frac{E[N_{r+1}]}{E[N_r]}$$

$$p_{gt}(w_1w_2) = \frac{r^*}{N}$$

Estimated frequencies in AP newswire (Church & Gale 1991)

$r=f_{MLE}$	$f_{empirical}$	f_{Lap}	f_{del}	f_{GT}
0	0.000027	0.000137	0.000037	0.000027
1	0.4480	.000274	0.396	0.446
2	1.25	0.000411	1.24	1.26
3	2.24	0.000548	2.23	2.24
4	3.23	0.000685	3.22	3.24
5	4.21	0.000822	4.22	4.22
6	5.23	0.000959	5.20	5.19
7	6.21	0.00109	6.21	6.21
8	7.21	0.00123	7.18	7.24

Differentiating based on history

- So far the methods considered have all used nothing but the raw frequency of an n-gram.
 - “the large” $P(\text{large}|\text{the})$
 - “the mauve” $P(\text{mauve}|\text{the})$
 - if $C(\text{“the large”}) = C(\text{“the mauve”})$ (e.g. = 0) then $P(\text{large}|\text{the}) = P(\text{mauve}|\text{the})$
 - This doesn’t seem right
- Also use frequency of its (n-1)-gram
 - $p(\text{large})$
 - $p(\text{mauve})$

Linear Interpolation

- Estimate probability of an n-gram from a weighted average of low-order...high order n-grams.

$$P_l(w_n|w_{n-2}, w_{n-1}) = \lambda_0 \frac{1}{V} + \lambda_1 P_1(w_n) + \lambda_2 P_2(w_n|w_{n-1}) + \lambda_3 P_3(w_n|w_{n-1}, w_{n-2})$$

- where λ 's sum to 1
- set λ 's by hand or from held-out data
- they can be functions of (equivalence-classed histories)
- Also known as “shrinkage” [Stein 1957]

- Works surprisingly well!

Assigning Probability to the “Unseen” Event

- At test time, see word, u , that wasn’t seen at training time
- $P(u) = ?$

- Replace all singleton word tokens in training data with special token, $\langle \text{UNK} \rangle$

Smoothing: Rest of the story

- Other methods:
 - backoff (Katz 1987): Try four-gram, if zero-count, try tri-gram, if zero-count, try bi-gram,...
 - Kneser and Ney (1995): Backoff n-gram counts not proportional to frequency of n-gram in training data but to expectation of how often it should occur in novel trigram (since one only uses backoff estimate when trigram not found)
 - Witten-Bell discounting
 - Smoothed maximum entropy models
 - See (Chen and Goodman 1998) for a survey.

- Progress in the field is often dominated, not by the need to create fancier more complex models,
- but by the need to do a good job of estimating parameters for the simpler models we already have.

Statistical Language Modeling



Noam Chomsky



Fred Jelinek

But it must be recognized that the notion of “probability of a sentence” is an entirely useless one, under any known interpretation of the term.
(1969)

Anytime a linguist leaves the group, the [speech] recognition rate goes up. (while at IBM speech group, 1968).

- Progress in the field is often dominated, not by the need to create fancier more complex models,
- but by the need to do a good job of estimating parameters for the simpler models we already have.
- Real benefit comes from targeted enhancements, and sharp toolset of excellent estimation techniques

Distinctiveness of NLP as an ML problem

- Most structure is hidden
- Relational, constraint satisfaction nature
- Long pipelines, with cascading errors
- Large and strange, sparse discrete distributions
- Large scale
- Feature-driven; performance driven

Large, sparse, discrete distributions

- Both features and assigned classes regularly involve multinomial distributions over huge numbers of values (often in the tens of thousands).
- The distributions are very uneven, and have fat tails.
- Enormous problems with data sparseness: much work on smoothing distributions...
- We normally have inadequate (labeled) data to estimate probabilities.
- Unknown/ unseen things are usually a central problem
- (Generally dealing with discrete distributions, though)