

Maximum Entropy
Lecture #13
Introduction to Natural Language Processing
CMPSCI 585, Spring 2004
University of Massachusetts Amherst



Andrew McCallum
(Slides from Jason Eisner)

1

summary of half of the course (statistics)

Probability is Useful

- We love probability distributions!
 - We've learned how to define & **use** $p(\dots)$ functions.
- Pick best output text T from a set of candidates
 - speech recognition; machine translation; OCR; spell correction...
 - maximize $p_1(T)$ for some appropriate distribution p_1
- Pick best annotation T for a fixed input I
 - text categorization; parsing; part-of-speech tagging ...
 - maximize $p(T | I)$; equivalently maximize joint probability $p(I, T)$
 - often define $p(I, T)$ by noisy channel: $p(I, T) = p(T) * p(I | T)$
 - speech recognition & other tasks above are cases of this too:
 - we're maximizing an appropriate $p_1(T)$ defined by $p(T | I)$
- Pick best probability distribution (a meta-problem!)
 - really, pick best parameters θ : train HMM, PCFG, n-grams, clusters ...
 - maximum likelihood; smoothing; EM if unsupervised (incomplete data)
 - Smoothing: $\max p(\theta | \text{data}) = \max p(\theta, \text{data}) = p(\theta)p(\text{data} | \theta)$

2

summary of other half of the course (linguistics)

Probability is Flexible

- We love probability distributions!
 - We've learned how to **define** & use $p(\dots)$ functions.
- We want $p(\dots)$ to define probability of *linguistic* objects
 - Sequences of words, tags, morphemes, phonemes (n-grams, FSMs, FSTs; Viterbi, collocations)
 - Vectors (naïve Bayes; clustering word senses)
 - Trees of (non)terminals (PCFGs; CKY, Earley)
- We've also seen some not-so-probabilistic stuff
 - Syntactic features, morphology. Could be stochasticized?
 - Methods can be quantitative & data-driven but not fully probabilistic: clustering, collocations, ...
- But probabilities have wormed their way into most things
- $p(\dots)$ has to capture our intuitions about the ling. data**

3

really so alternative?

An Alternative Tradition

- Old AI hacking technique:
 - Possible parses (or whatever) have scores.
 - Pick the one with the best score.
 - How do you define the score?
 - Completely ad hoc!
 - Throw anything you want into the stew
 - Add a bonus for this, a penalty for that, etc.
- "Learns" over time – as you adjust bonuses and penalties by hand to improve performance.
- Total kludge, but totally flexible too ...
 - Can throw in **any** intuitions you might have

4

really so alternative?

An Alternative Tradition

- Old AI hacking technique:
 - Possible parses (or whatever) have scores.
 - Pick the one with the best score.
 - How do you define the score?
 - Completely ad hoc!
 - Throw anything you want into the stew
 - Add a bonus for this, a penalty for that, etc.
- "Learns" over time – as you adjust bonuses and penalties by hand to improve performance.
- Total kludge, but totally flexible too ...
 - Can throw in **any** intuitions you might have

Probabilistic Revolution
Not Really a Revolution,
Critics Say

Log-probabilities no more
 than scores in disguise

“We’re just adding stuff up
 like the old corrupt regime
 did,” admits spokesperson

5

Nuthin' but adding weights

- n-grams: ... + $\log p(w_7 | w_5, w_6) + \log p(w_8 | w_6, w_7) + \dots$
- PCFG: $\log p(\text{NP VP} | \text{S}) + \log p(\text{Papa} | \text{NP}) + \log p(\text{VP PP} | \text{VP}) \dots$
- HMM tagging: ... + $\log p(t_7 | t_5, t_6) + \log p(w_7 | t_7) + \dots$
- Noisy channel: $[\log p(\text{source})] + [\log p(\text{data} | \text{source})]$
- Naïve Bayes: $\log p(\text{Class}) + \log p(\text{feature1} | \text{Class}) + \log p(\text{feature2} | \text{Class}) \dots$
- Note: Just as in probability, bigger weights are better.**

6

Nuthin' but adding weights

- n-grams: ... + log p(w7 | w5,w6) + log(w8 | w6, w7) + ...
- PCFG: log p(NP VP | S) + log p(Papa | NP) + log p(VP PP | VP) ...
- HMM tagging: ... + log p(t7 | t5, t6) + log p(w7 | t7) + ...
- Noisy channel: [log p(source)] + [log p(data | source)]
- Naïve Bayes:
 - log(Class) + log(feature1 | Class) + log(feature2 | Class) + ...
 - Can regard any linguistic object as a collection of features (here, doc = a collection of words, but could have non-word features)
 - Weight of the object = total weight of features
 - Our weights have always been conditional log-probs (≤ 0) but that is going to change in a few minutes!

7

Probabilists Rally Behind their Paradigm

“.2, .4, .6, .8! We're not gonna take your bait!”

- Can estimate our parameters automatically
 - e.g., log p(t7 | t5, t6) (trigram tag probability)
 - from supervised or unsupervised data (ratio of counts)
- Our results are more meaningful
 - Can use probabilities to place bets, quantify risk
 - e.g., how sure are we that this is the correct parse?
- Our results can be meaningfully combined \Rightarrow modularity!
 - Multiply indep. conditional probs – normalized, unlike scores
 - p(English text) * p(English phonemes | English text) * p(Jap. phonemes | English phonemes) * p(Jap. text | Jap. phonemes)
 - p(semantics) * p(syntax | semantics) * p(morphology | syntax) * p(phonology | morphology) * p(sounds | phonology)

8

Probabilists Regret Being Bound by Principle

- Ad-hoc approach does have one advantage
- Consider e.g. Naïve Bayes for text categorization:
 - Buy this supercalifragilistic Ginsu knife set for only \$39 today ...
- Some useful features:
 - Contains Buy
 - Contains supercalifragilistic
 - Contains a dollar amount under \$100
 - Contains an imperative sentence
 - Reading level = 8th grade
 - Mentions money (use word classes and/or regexp to detect this)
- Naïve Bayes: pick C maximizing p(C) * p(feats 1 | C) * ...
- What assumption does Naïve Bayes make? True here?

spam
.5
.02
.9 .1

ham

9

Probabilists Regret Being Bound by Principle

- Ad-hoc approach does have one advantage
- Consider e.g. Naïve Bayes for text categorization:
 - Buy this supercalifragilistic Ginsu knife set for only \$39 today ...
- Some useful features:
 - Contains Buy
 - Contains supercalifragilistic
 - Contains a dollar amount under \$100
 - Contains an imperative sentence
 - Reading level = 8th grade
 - Mentions money (use word classes and/or regexp to detect this)
- Naïve Bayes: pick C maximizing p(C) * p(feats 1 | C) * ...
- What assumption does Naïve Bayes make? True here?

spam
.5
.02
.9 .1

ham

50% of spam has this – 25x more likely than in ham
Contains a dollar amount under \$100
but here are the emails with both features – only 25x!
90% of spam has this – 9x more likely than in ham

Naive Bayes claims .5*.9=45% of spam has both features – 25*.9=225x more likely than in ham.

10

Probabilists Regret Being Bound by Principle

- But ad-hoc approach does have one advantage
 - Can adjust scores to compensate for feature overlap ...
- Some useful features of this message:
 - Contains a dollar amount under \$100
 - Mentions money
- Naïve Bayes: pick C maximizing p(C) * p(feats 1 | C) * ...
- What assumption does Naïve Bayes make? True here?

spam
.5
.02
.9 .1

ham

log prob adjusted
-1 -5.6 -0.85 -2.3
-1.15 -3.3 -1.15 -3.3
subtract "money" score already included

11

Revolution Corrupted by Bourgeois Values

- Naïve Bayes needs overlapping but independent features
- But not clear how to restructure these features like that:
 - Contains Buy
 - Contains supercalifragilistic
 - Contains a dollar amount under \$100
 - Contains an imperative sentence
 - Reading level = 7th grade
 - Mentions money (use word classes and/or regexp to detect this)
 - ...
- Boy, we'd like to be able to throw all that useful stuff in without worrying about feature overlap/independence.
- Well, maybe we can add up scores and pretend like we got a log probability:

12

Revolution Corrupted by Bourgeois Values

- Naïve Bayes needs overlapping but **independent** features
 - But not clear how to restructure these features like that:
 - +4 Contains Buy
 - +0.2 Contains supercalifragilistic
 - +1 Contains a dollar amount under \$100
 - +2 Contains an imperative sentence
 - 3 Reading level = 7th grade
 - +5 Mentions money (use word classes and/or regex to detect this)
 - ...
- } total: 5.77
- Boy, we'd like to be able to throw all that useful stuff in without worrying about feature overlap/independence.
 - Well, maybe we can add up scores and **pretend** like we got a log probability: **$\log p(\text{feats} \mid \text{spam}) = 5.77$**
 - Oops, then $p(\text{feats} \mid \text{spam}) = \exp 5.77 = 320.5$

13

Renormalize by 1/Z to get a Log-Linear Model

- $p(\text{feats} \mid \text{spam}) = \exp 5.77 = 320.5$ scale down so everything < 1 and sums to 1!
- $p(m \mid \text{spam}) = (1/Z(\lambda)) \exp \sum_i \lambda_i f_i(m)$ where
 - m is the email message
 - λ_i is weight of feature i
 - $f_i(m) \in \{0,1\}$ according to whether m has feature i
 - More generally, allow $f_i(m) = \text{count or strength of feature}$.
 - $1/Z(\lambda)$ is a normalizing factor making $\sum_m p(m \mid \text{spam}) = 1$ (summed over all possible messages m ! hard to find!)
- The weights we add up are basically arbitrary.
- They don't have to mean anything, so long as they give us a good probability.
- Why is it called "log-linear"?

14

Why Bother?

- Gives us probs, not just scores.
 - Can use them to bet, or combine w/ other probs.
- We can now learn weights from data!
 - Choose weights λ_i that maximize logprob of labeled training data = $\log \prod_j p(c_j) p(m_j \mid c_j)$
 - where $c_j \in \{\text{ham}, \text{spam}\}$ is classification of message m_j
 - and $p(m_j \mid c_j)$ is log-linear model from previous slide
 - Convex function – easy to maximize! (why?)
- But:** $p(m_j \mid c_j)$ for a given λ requires $Z(\lambda)$: hard!

15

Attempt to Cancel out Z

- Set weights to maximize $\prod_j p(c_j) p(m_j \mid c_j)$
 - where $p(m \mid \text{spam}) = (1/Z(\lambda)) \exp \sum_i \lambda_i f_i(m)$
 - But** normalizer $Z(\lambda)$ is awful sum over all possible emails
- So instead:** Maximize $\prod_j p(c_j \mid m_j)$
 - Doesn't model the emails m_j , only their classifications c_j
 - Makes more sense anyway given our feature set
- $p(\text{spam} \mid m) = p(\text{spam})p(m \mid \text{spam}) / (p(\text{spam})p(m \mid \text{spam}) + p(\text{ham})p(m \mid \text{ham}))$
- Z appears in both numerator and denominator
- Alas, doesn't cancel out because Z differs for the spam and ham models
- But we can fix this ...

16

So: Modify Setup a Bit

- Instead of having separate models $p(m \mid \text{spam}) * p(\text{spam})$ vs. $p(m \mid \text{ham}) * p(\text{ham})$
- Have just one joint model $p(m, c)$ gives us both $p(m, \text{spam})$ and $p(m, \text{ham})$
- Equivalent to changing feature set to:
 - spam ← weight of this feature is $\log p(\text{spam}) + \text{a constant}$
 - spam and Contains Buy ← old spam model's weight for "contains Buy"
 - spam and Contains supercalifragilistic
 - ...
 - ham ← weight of this feature is $\log p(\text{ham}) + \text{a constant}$
 - ham and Contains Buy ← old ham model's weight for "contains Buy"
 - ham and Contains supercalifragilistic
- No real change, but 2 categories now share single feature set and single value of $Z(\lambda)$

17

Now we can cancel out Z

- Now $p(m, c) = (1/Z(\lambda)) \exp \sum_i \lambda_i f_i(m, c)$ where $c \in \{\text{ham}, \text{spam}\}$
- Old:** choose weights λ_i that maximize prob of labeled training data = $\prod_j p(m_j, c_j)$
 - New:** choose weights λ_i that maximize prob of labels given messages = $\prod_j p(c_j \mid m_j)$
 - Now Z cancels out of conditional probability!
 - $p(\text{spam} \mid m) = p(m, \text{spam}) / (p(m, \text{spam}) + p(m, \text{ham}))$
 - $= \exp \sum_i \lambda_i f_i(m, \text{spam}) / (\exp \sum_i \lambda_i f_i(m, \text{spam}) + \exp \sum_i \lambda_i f_i(m, \text{ham}))$
 - Easy to compute now ...
 - $\prod_j p(c_j \mid m_j)$ is still convex, so easy to maximize too

18

Maximum Entropy

- Suppose there are 10 classes, A through J.
- I don't give you any other information.
- **Question:** Given message m: what is your guess for $p(C | m)$?
- Suppose I tell you that 55% of all messages are in class A.
- **Question:** Now what is your guess for $p(C | m)$?
- Suppose I **also** tell you that 10% of all messages contain `Buy` and 80% of these are in class A or C.
- **Question:** Now what is your guess for $p(C | m)$, if m contains `Buy`?
- **OUCH!**

19

Maximum Entropy

	A	B	C	D	E	F	G	H	I	J
Buy	.051	.0025	.029	.0025	.0025	.0025	.0025	.0025	.0025	.0025
Other	.499	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446

- Column A sums to 0.55 ("55% of all messages are in class A")

20

Maximum Entropy

	A	B	C	D	E	F	G	H	I	J
Buy	.051	.0025	.029	.0025	.0025	.0025	.0025	.0025	.0025	.0025
Other	.499	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446

- Column A sums to 0.55
- Row `Buy` sums to 0.1 ("10% of all messages contain `Buy`")

21

Maximum Entropy

	A	B	C	D	E	F	G	H	I	J
Buy	.051	.0025	.029	.0025	.0025	.0025	.0025	.0025	.0025	.0025
Other	.499	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446

- Column A sums to 0.55
- Row `Buy` sums to 0.1
- (`Buy`, A) and (`Buy`, C) cells sum to 0.08 ("80% of the 10%")

- Given these constraints, fill in cells "as equally as possible": maximize the entropy (related to cross-entropy, perplexity)

Entropy = $-.051 \log .051 - .0025 \log .0025 - .029 \log .029 - \dots$
 Largest if probabilities are evenly distributed

22

Maximum Entropy

	A	B	C	D	E	F	G	H	I	J
Buy	.051	.0025	.029	.0025	.0025	.0025	.0025	.0025	.0025	.0025
Other	.499	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446

- Column A sums to 0.55
- Row `Buy` sums to 0.1
- (`Buy`, A) and (`Buy`, C) cells sum to 0.08 ("80% of the 10%")
- Given these constraints, fill in cells "as equally as possible": maximize the entropy
- Now $p(\text{Buy}, C) = .029$ and $p(C | \text{Buy}) = .29$
- We got a compromise: $p(C | \text{Buy}) < p(A | \text{Buy}) < .55$

23

Generalizing to More Features

		<\$100									
Other		A	B	C	D	E	F	G	H	...	
Buy		.051	.0025	.029	.0025	.0025	.0025	.0025	.0025	.0025	
Other		.499	.0446	.0446	.0446	.0446	.0446	.0446	.0446	.0446	

24

What we just did

- For each feature (“contains Buy”), see what fraction of training data has it
- Many distributions $p(c,m)$ would predict these fractions (including the unsmoothed one where all mass goes to feature combos we’ve actually seen)
- Of these, pick distribution that has max entropy
- **Amazing Theorem:** This distribution has the form $p(m,c) = (1/Z(\lambda)) \exp \sum_i \lambda_i f_i(m,c)$
 - So it is log-linear. In fact it is the same log-linear distribution that maximizes $\prod_j p(m_j, c_j)$ as before!
- Gives another motivation for our log-linear approach.

25

Log-linear form derivation

- Say we are given some **constraints** in the form of feature expectations:

$$\sum_x p(x) f_i(x) = \alpha_i$$

- In general, there may be many distributions $p(x)$ that satisfy the constraints. Which one to pick?
- The one with maximum entropy (making fewest possible additional assumptions---Occur’s Razor)
- This yields an optimization problem

$$\max H(p(x)) = - \sum_x p(x) \log p(x)$$

$$\text{Subject to } \sum_x p(x) f_i(x) = \alpha_i, \forall i \text{ and } \sum_x p(x) = 1$$

26

Log-linear form derivation

- To solve the maxent problem, we use Lagrange multipliers:

$$L = - \sum_x p(x) \log p(x) - \sum_i \theta_i \left(\sum_x p(x) f_i(x) - \alpha_i \right) - \mu \left(\sum_x p(x) - 1 \right)$$

$$\frac{\partial L}{\partial p(x)} = 1 + \log p(x) - \sum_i \theta_i f_i(x) - \mu$$

$$p^*(x) = e^{\mu-1} \exp \left\{ \sum_i \theta_i f_i(x) \right\}$$

$$Z(\theta) = e^{1-\mu} = \sum_x \exp \left\{ \sum_i \theta_i f_i(x) \right\}$$

$$p(x|\theta) = \frac{1}{Z(\theta)} \exp \left\{ \sum_i \theta_i f_i(x) \right\}$$

- So feature constraints + maxent implies exponential family.
- Problem is convex, so solution is unique.

27

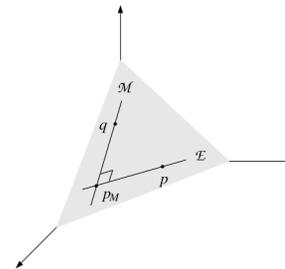
MaxEnt = Max Likelihood

Define two submanifolds on the probability simplex $p(x)$.

The first is \mathcal{E} , the set of all exponential family distributions based on a particular set of features $f_i(x)$.

The second is \mathcal{M} , the set of all distributions that satisfy the feature expectation constraints.

They intersect at a single distribution p_M , the maxent, maximum likelihood



28

$$\begin{aligned} \ell(\theta; \mathcal{D}) &= \sum_x n(x) \log p(x|\theta) \\ &= \sum_x n(x) \left(\sum_i \theta_i f_i(x) - \log Z(\theta) \right) \\ &= \sum_x n(x) \sum_i \theta_i f_i(x) - N \log Z(\theta) \\ \frac{\partial \ell}{\partial \theta_i} &= \sum_x n(x) f_i(x) - N \frac{\partial}{\partial \theta_i} \log Z(\theta) \\ &= \sum_x n(x) f_i(x) - N \sum_x p(x|\theta) f_i(x) \\ \Rightarrow \sum_x p(x|\theta) f_i(x) &= \sum_x \frac{n(x)}{N} f_i(x) = \sum_x \bar{p}(x) f_i(x) \end{aligned}$$

Derivative of log partition function is the expectation of the feature.
At ML estimate, model expectations match empirical feature counts.

29

Recipe for a Conditional MaxEnt Classifier

1. Gather **constraints** from training data:

$$\alpha_{iy} = \bar{E}[f_{iy}] = \sum_{x_j, y_j \in D} f_{iy}(x_j, y_j)$$

2. Initialize all parameters to zero.

3. Classify training data with current parameters. Calculate **expectations**.

$$E_{\Theta}[f_{iy}] = \sum_{x_j \in D} \sum_{y_j'} p_{\Theta}(y_j' | x_j) f_{iy}(x_j, y_j')$$

4. Gradient is $\tilde{E}[f_{iy}] - E_{\Theta}[f_{iy}]$
5. Take a step in the direction of the gradient
6. Until convergence, return to step 3.

30

Overfitting

- If we have too many features, we can choose weights to model the training data perfectly.
- If we have a feature that only appears in spam training, not ling training, it will get weight ∞ to maximize $p(\text{spam} | \text{feature})$ at 1.
- These behaviors overfit the training data.
- Will probably do poorly on test data.

31

Solutions to Overfitting

- Throw out rare features.
 - Require every feature to occur > 4 times, and > 0 times with ling, and > 0 times with spam.
- Only keep 1000 features.
 - Add one at a time, always greedily picking the one that most improves performance on held-out data.
- Smooth the observed feature counts.
- Smooth the weights by using a prior.
 - $\max p(\lambda | \text{data}) = \max p(\lambda, \text{data}) = p(\lambda)p(\text{data} | \lambda)$
 - decree $p(\lambda)$ to be high when most weights close to 0

32