

Maximum Entropy

Lecture #13

Introduction to Natural Language Processing

CMPSCI 585, Fall 2007

University of Massachusetts Amherst



Andrew McCallum

(Slides from Jason Eisner and Dan Klein)

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;**
 - Smoothing: **$\max p(\theta | \text{data}) = \max p(\theta, \text{data}) = p(\theta)p(\text{data} | \theta)$**

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**

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

really so alternative?

~~An Alternative~~ Tradition

- Old A
- Pos
- Pic
- Ho
-
-
-
- “Learn
- pena
- Total
- Ca

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

uses and
nce. 😊

e

Nuthin' but adding weights

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

Nuthin' but adding weights

- n-grams: ... + $\log p(w_7 \mid w_5, w_6) + \log(w_8 \mid w_6, w_7) + \dots$
- PCFG: $\log p(\text{NP VP} \mid S) + \log p(\text{Papa} \mid \text{NP}) + \log p(\text{VP PP} \mid \text{VP}) \dots$
- HMM tagging: ... + $\log p(t_7 \mid t_5, t_6) + \log p(w_7 \mid t_7) + \dots$
- Noisy channel: $[\log p(\text{source})] + [\log p(\text{data} \mid \text{source})]$
- Naïve Bayes:
 - $\log(\text{Class}) + \log(\text{feature1} \mid \text{Class}) + \log(\text{feature2} \mid \text{Class}) + \dots$
 - 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!

Probabilists Rally Behind their Paradigm

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

- Can estimate our parameters *automatically*
 - e.g., $\log p(t_7 | t_5, t_6)$ (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
 - $p(\text{English text}) * p(\text{English phonemes} | \text{English text}) * p(\text{Jap. phonemes} | \text{English phonemes}) * p(\text{Jap. text} | \text{Jap. phonemes})$
 - $p(\text{semantics}) * p(\text{syntax} | \text{semantics}) * p(\text{morphology} | \text{syntax}) * p(\text{phonology} | \text{morphology}) * p(\text{sounds} | \text{phonology})$

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 $p(C) * p(\text{feat 1} | C) * \dots$
- What assumption does Naïve Bayes make? True here?

spam
.5

ham
.02

.9

.1

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:

spam	ham	
.5	.02	<ul style="list-style-type: none"> • 50% of spam has this – 25x more likely than in ham • Contains a dollar amount under \$100 <i>but here are the emails with both features – only 25x!</i>
.9	.1	<ul style="list-style-type: none"> • 90% of spam has this – 9x more likely than in ham • Mentions money

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

- Naïve Bayes: pick C maximizing $p(C) * p(\text{feat 1} | C) * \dots$
- What assumption does Naïve Bayes make? True here?

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:

subtract "money" score
already included

log prob adjusted

spam	ham		spam	ham	spam	ham
.5	.02	• Contains a dollar amount under \$100	-1	-5.6	-.85	-2.3
.9	.1	• Mentions money	-.15	-3.3	-.15	-3.3

- Naïve Bayes: pick C maximizing $p(C) * p(\text{feat 1} | C) * \dots$
- What assumption does Naïve Bayes make? True here?

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:

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 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: **$\log p(\text{feats} \mid \text{spam}) = 5.77$**
- Oops, then **$p(\text{feats} \mid \text{spam}) = \exp 5.77 = 320.5$**

total: 5.77

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

scale down so everything < 1 and sums to 1!

- $p(\text{feats} \mid \text{spam}) = \text{exp } 5.77 = 320.5$
- $p(m \mid \text{spam}) = (1/Z(\lambda)) \exp \sum_i \lambda_i f_i(m)$ 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"?

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_j that maximize logprob of labeled training data = $\log \prod_j p(c_j) p(m_j | c_j)$
 - where $c_j \in \{\text{ham}, \text{spam}\}$ is classification of message m_j
 - and $p(m_j | c_j)$ is log-linear model from previous slide
 - Convex function – easy to maximize! (why?)
- **But:** $p(m_j | c_j)$ for a given λ requires $Z(\lambda)$: hard!

Attempt to Cancel out Z

- Set weights to maximize $\prod_j p(c_j) p(m_j | c_j)$
 - where $p(m | \text{spam}) = (1/Z(\lambda)) \exp \sum_i \lambda_i f_i(m)$
 - **But** normalizer $Z(\lambda)$ is awful sum over all possible emails
- **So instead:** Maximize $\prod_j p(c_j | m_j)$
 - Doesn't model the emails m_j , only their classifications c_j
 - Makes more sense anyway given our feature set
- $p(\text{spam} | m) = p(\text{spam})p(m|\text{spam}) / (p(\text{spam})p(m|\text{spam})+p(\text{ham})p(m|\text{ham}))$
- 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 ...

So: Modify Setup a Bit

- Instead of having separate models
 $p(m|\text{spam}) * p(\text{spam})$ vs. $p(m|\text{ham}) * p(\text{ham})$
- 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 $\log p(\text{ling}) + \text{a constant}$
 - `ham and Contains Buy` ← old ling 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 $Z(\lambda)$

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 = $\prod_j p(c_j | m_j)$
- Now Z cancels out of conditional probability!
 - $p(\text{spam} | m) = p(m, \text{spam}) / (p(m, \text{spam}) + p(m, \text{ham}))$
 $= \exp \sum_i \lambda_i f_i(m, \text{spam}) / (\exp \sum_i \lambda_i f_i(m, \text{spam}) + \exp \sum_i \lambda_i f_i(m, \text{ham}))$
 - Easy to compute now ...
 - $\prod_j p(c_j | m_j)$ is still convex, so easy to maximize too

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!**

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”)

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`”)

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

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: $p(C | \text{Buy}) < p(A | \text{Buy}) < .55$

Generalizing to More Features

A 3D grid representing a feature space. The top face is labeled with $< \$100$. The left face is labeled with Other. The front face is a table with columns A, B, C, D, E, F, G, H, and ... and rows Buy and Other. The cells for Buy A and Buy C are highlighted in yellow.

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

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 $p(m,c) = (1/Z(\lambda)) \exp \sum_j \lambda_j f_j(m,c)$
 - 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.

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---Occum's Razor)
- This yields an optimization problem

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

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

Log-linear form derivation

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

$$L = - \sum_{\mathbf{x}} p(\mathbf{x}) \log p(\mathbf{x}) - \sum_i \theta_i \left(\sum_{\mathbf{x}} p(\mathbf{x}) f_i(\mathbf{x}) - \alpha_i \right) - \mu \left(\sum_{\mathbf{x}} p(\mathbf{x}) - 1 \right)$$

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

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

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

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

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

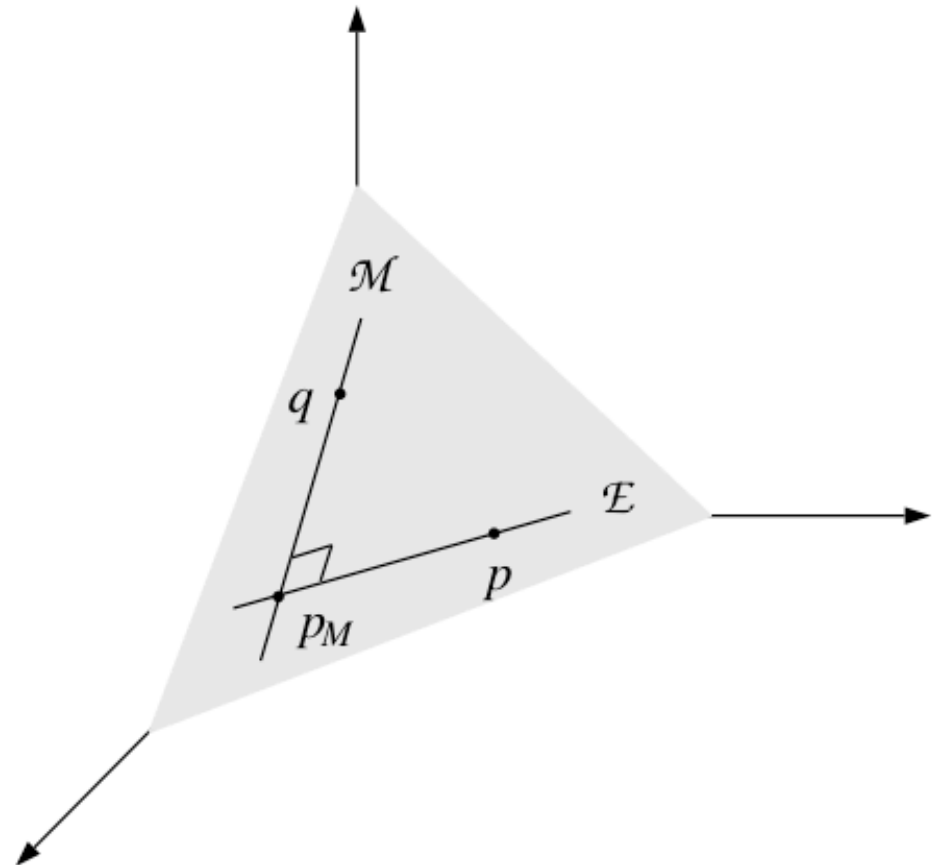
MaxEnt = Max Likelihood

Define two submanifolds on the probability simplex $p(\mathbf{x})$.

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

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





Exponential Model Likelihood

- Maximum Likelihood (Conditional) Models :
 - Given a model form, choose values of parameters to maximize the (conditional) likelihood of the data.
- Exponential model form, for a data set (C,D):

$$\log P(C | D, \lambda) = \sum_{(c,d) \in (C,D)} \log P(c | d, \lambda) = \sum_{(c,d) \in (C,D)} \log \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}$$



Building a Maxent Model

- Define features (indicator functions) over data points.
 - Features represent sets of data points which are distinctive enough to deserve model parameters.
 - Usually features are added incrementally to “target” errors.
- For any given feature weights, we want to be able to calculate:
 - Data (conditional) likelihood
 - Derivative of the likelihood wrt each feature weight
 - Use expectations of each feature according to the model
- Find the optimum feature weights (next part).



The Likelihood Value

- The (log) conditional likelihood is a function of the iid data (C,D) and the parameters λ :

$$\log P(C | D, \lambda) = \log \prod_{(c,d) \in (C,D)} P(c | d, \lambda) = \sum_{(c,d) \in (C,D)} \log P(c | d, \lambda)$$

- If there aren't many values of c , it's easy to calculate:

$$\log P(C | D, \lambda) = \sum_{(c,d) \in (C,D)} \log \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}$$

- We can separate this into two components:

$$\log P(C | D, \lambda) = \sum_{(c,d) \in (C,D)} \log \exp \sum_i \lambda_i f_i(c, d) - \sum_{(c,d) \in (C,D)} \log \sum_{c'} \exp \sum_i \lambda_i f_i(c', d)$$

$$\log P(C | D, \lambda) = N(\lambda) - M(\lambda)$$

- The derivative is the difference between the derivatives of each component



The Derivative I: Numerator

$$\begin{aligned}\frac{\partial N(\lambda)}{\partial \lambda_i} &= \frac{\partial \sum_{(c,d) \in (C,D)} \log \exp \sum_i \lambda_{ci} f_i(c,d)}{\partial \lambda_i} = \frac{\partial \sum_{(c,d) \in (C,D)} \sum_i \lambda_i f_i(c,d)}{\partial \lambda_i} \\ &= \sum_{(c,d) \in (C,D)} \frac{\partial \sum_i \lambda_i f_i(c,d)}{\partial \lambda_i} \\ &= \sum_{(c,d) \in (C,D)} f_i(c,d)\end{aligned}$$

Derivative of the numerator is: the empirical count(f_i, c)



The Derivative II: Denominator

$$\begin{aligned}
 \frac{\partial M(\lambda)}{\partial \lambda_i} &= \frac{\partial \sum_{(c,d) \in (C,D)} \log \sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}{\partial \lambda_i} \\
 &= \sum_{(c,d) \in (C,D)} \frac{1}{\sum_{c''} \exp \sum_i \lambda_i f_i(c'', d)} \frac{\partial \sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}{\partial \lambda_i} \\
 &= \sum_{(c,d) \in (C,D)} \frac{1}{\sum_{c''} \exp \sum_i \lambda_i f_i(c'', d)} \sum_{c'} \frac{\exp \sum_i \lambda_i f_i(c', d)}{1} \frac{\partial \sum_i \lambda_i f_i(c', d)}{\partial \lambda_i} \\
 &= \sum_{(c,d) \in (C,D)} \sum_{c'} \frac{\exp \sum_i \lambda_i f_i(c', d)}{\sum_{c''} \exp \sum_i \lambda_i f_i(c'', d)} \frac{\partial \sum_i \lambda_i f_i(c', d)}{\partial \lambda_i} \\
 &= \sum_{(c,d) \in (C,D)} \sum_{c'} P(c' | d, \lambda) f_i(c', d) = \text{predicted count}(f_i, \lambda)
 \end{aligned}$$



The Derivative III

$$\frac{\partial \log P(C | D, \lambda)}{\partial \lambda_i} = \text{actual count}(f_i, C) - \text{predicted count}(f_i, \lambda)$$

- The optimum parameters are the ones for which each feature's **predicted expectation** equals its **empirical expectation**. The optimum distribution is:
 - Always unique (but parameters may not be unique)
 - Always exists (if features counts are from actual data).
- Features can have high model expectations (predicted counts) either because they have large weights or because they occur with other features which have large weights.



Summary

- We have a function to optimize:

$$\log P(C | D, \lambda) = \sum_{(c,d) \in (C,D)} \log \frac{\exp \sum_i \lambda_i f_i(c,d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c',d)}$$

- We know the function's derivatives:

$$\partial \log P(C | D, \lambda) / \partial \lambda_i = \text{actual count}(f_i, C) - \text{predicted count}(f_i, \lambda)$$

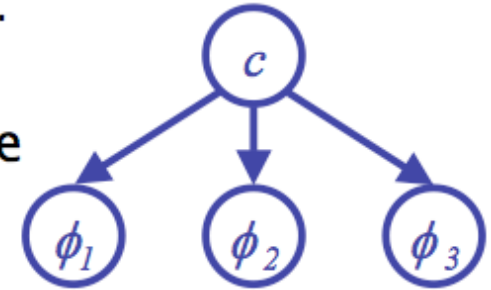
- Perfect situation for general optimization (Part II)

By gradient ascent or conjugate gradient.



Comparison to Naïve-Bayes

- Naïve-Bayes is another tool for classification:
 - We have a bunch of random variables (data features) which we would like to use to predict another variable (the class):
 - The Naïve-Bayes likelihood over classes is:



$$P(c | d, \lambda) = \frac{P(c) \prod_i P(\phi_i | c)}{\sum_{c'} P(c') \prod_i P(\phi_i | c')} \Rightarrow \frac{\exp\left[\log P(c) + \sum_i \log P(\phi_i | c)\right]}{\sum_{c'} \exp\left[\log P(c') + \sum_i \log P(\phi_i | c')\right]}$$

Naïve-Bayes is just an exponential model.

$$\Rightarrow \frac{\exp\left[\sum_i \lambda_{ic} f_{ic}(d, c)\right]}{\sum_{c'} \exp\left[\sum_i \lambda_{ic'} f_{ic'}(d, c')\right]}$$



Comparison to Naïve-Bayes

- The primary differences between Naïve-Bayes and maxent models are:

Naïve-Bayes

Trained to maximize joint likelihood of data and classes.

Features assumed to supply independent evidence.

Feature weights can be set independently.

Features must be of the conjunctive $\Phi(d) \wedge c = c_i$ form.

Maxent

Trained to maximize the conditional likelihood of classes.

Features weights take feature dependence into account.

Feature weights must be mutually estimated.

Features need not be of the conjunctive form (but usually are).

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 $p(\text{spam} \mid \text{feature})$ at 1.
- These behaviors overfit the training data.
- Will probably do poorly on test data.

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



Smoothing: Priors (MAP)

- What if we had a prior expectation that parameter values wouldn't be very large?
- We could then balance evidence suggesting large parameters (or infinite) against our prior.
- The evidence would never totally defeat the prior, and parameters would be smoothed (and kept finite!).
- We can do this explicitly by changing the optimization objective to maximum posterior likelihood:

$$\log P(C, \lambda | D) = \log P(\lambda) + \log P(C | D, \lambda)$$

Posterior

Prior

Evidence

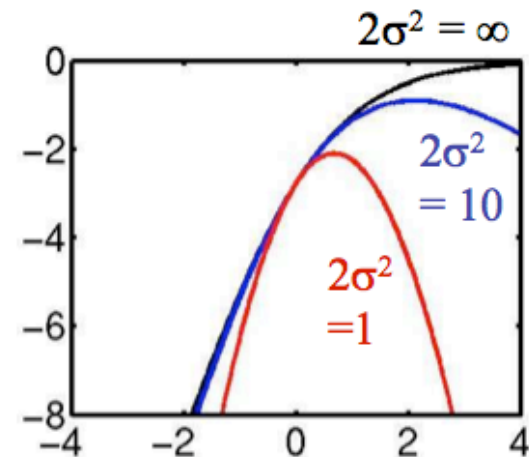


Smoothing: Priors

- Gaussian, or quadratic, priors:
 - Intuition: parameters shouldn't be large.
 - Formalization: prior expectation that each parameter will be distributed according to a gaussian with mean μ and variance σ^2 .

$$P(\lambda_i) = \frac{1}{\sigma_i \sqrt{2\pi}} \exp\left(-\frac{(\lambda_i - \mu_i)^2}{2\sigma_i^2}\right)$$

- Penalizes parameters for drifting to far from their mean prior value (usually $\mu=0$).
- $2\sigma^2=1$ works surprisingly well.



They don't even capitalize my name anymore!



Recipe for a Conditional MaxEnt Classifier

1. Gather *constraints* from training data:

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

2. Initialize all parameters to zero.
3. Classify training data with current parameters. Calculate *expectations*.

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

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.