

Linear Models for Statistical Natural Language Processing

Jacob Eisenstein

September 3, 2014

Chapter 1

Introduction

This is a collection of notes that I use for teaching Georgia Tech Computer Science 4650 and 7650, “Natural Language.” The notes focus on what I view as a core subset of the field of natural language processing, unified by the concept of linear models. This includes approaches to document classification, word sense disambiguation, sequence labeling (part-of-speech tagging and named entity recognition), parsing, coreference resolution, relation extraction, discourse analysis, and, to a limited degree, language modeling and machine translation. The theme was inspired by Fernando Pereira’s EMNLP 2008 keynote, “Are linear models right for language.”¹ The notes are heavily influenced by several other good resources (e.g., Manning and Schütze, 1999; Jurafsky and Martin, 2009; Figueiredo et al., 2013; Collins, 2013), but for various reasons I wanted to create something of my own.

¹You can see a version of this talk — not the one I saw — online at vimeo.com/30676245

Chapter 2

Notation

w_n	word token at position n
\mathbf{x}_i	a vector of feature counts for instance i , often word counts
N	number of training instances
V	number of words in vocabulary
$\boldsymbol{\theta}$	a vector of weights
y_i	the label for instance i
\mathbf{y}	vector of labels across all instances
\mathcal{Y}	set of all possible labels
K	number of possible labels $K = \# \mathcal{Y} $
$\mathbf{f}(\mathbf{x}_i, y_i)$	feature vector for instance i with label y_i
$P(A)$	probability function of event A
$p_B(b)$	the marginal probability of random variable B taking value b

Chapter 3

Linear classification and features

Suppose you want to build a spam detector. Spam vs. Ham. How would you do it, using only the text in the email?

One solution is to represent document i as a column vector of word counts: $\mathbf{x}_i = [0 \ 1 \ 1 \ 0 \ 0 \ 2 \ 0 \ 1 \ 13 \ 0 \dots]^\top$, where $x_{i,j}$ is the count of word j in document i . Suppose the size of the vocabulary is V , so that the length of \mathbf{x}_i is also V .

We’ve thrown out grammar, sentence boundaries, paragraphs — everything but the words! But this could still work. If you see the word *free*, is it spam or ham? How about *calls*? How about *Bayesian*? One approach would be to define a “spamminess” score for every word in the dictionary, and then just add them up. This is also a commonly-used approach to sentiment analysis, where each word is scored as one of $\{1, 0, -1\}$, with 1 indicating positive sentiment and -1 indicating negative sentiment.

These scores are called **weights**, written θ , and we’ll spend a lot of time later talking about where they come from. But for now, let’s generalize: suppose we want to build a multi-way classifier to distinguish stories about sports, celebrities, music, and business. Each label is an element y_i in a set of K possible labels \mathcal{Y} . Then for any pair $\langle \mathbf{x}_i, y_i \rangle$, we can define a *feature vector* $\mathbf{f}(\mathbf{x}_i, y_i)$, such that:

$$\mathbf{f}(\mathbf{x}, y = 0) = [\mathbf{x}_i^\top \ \mathbf{0}_{V(K-1)}^\top]^\top \quad (3.1)$$

$$\mathbf{f}(\mathbf{x}, y = 1) = [\mathbf{0}_V^\top \ \mathbf{x}_i^\top \ \mathbf{0}_{V(K-2)}^\top]^\top \quad (3.2)$$

$$\mathbf{f}(\mathbf{x}, y = 2) = [\mathbf{0}_{2V}^\top \ \mathbf{x}_i^\top \ \mathbf{0}_{V(K-3)}^\top]^\top \quad (3.3)$$

$$\dots \quad (3.4)$$

$$\mathbf{f}(\mathbf{x}, K) = [\mathbf{0}_{V(K-1)}^\top \ \mathbf{x}_i^\top]^\top, \quad (3.5)$$

where $\mathbf{0}_{VK}$ is a column vector of VK zeros. Often we’ll add an **offset** feature at

the end of \mathbf{x} , which is always 1; we then have to also add an extra zero to each of the zero vectors. This gives the entire feature vector $\mathbf{f}(\mathbf{x}, y)$ a length of $(V + 1)K$.

Now, given a vector of weights, $\boldsymbol{\theta} \in \mathcal{R}^{(V+1)K}$, we can compute the inner product $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$. Then for any document \mathbf{x}_i , we can predict a label \hat{y} as

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y) \quad (3.6)$$

We could just set the weights by hand. If we wanted to distinguish, say, English from Spanish, we could just use English and Spanish dictionaries, and set each weight to 1. For example,

$$\begin{array}{ll} \theta_{\text{english}, \text{bicycle}} = 1 & \theta_{\text{spanish}, \text{bicycle}} = 0 \\ \theta_{\text{english}, \text{bicicleta}} = 0 & \theta_{\text{spanish}, \text{bicicleta}} = 1 \\ \theta_{\text{english}, \text{con}} = 1 & \theta_{\text{spanish}, \text{con}} = 1 \\ \theta_{\text{english}, \text{ordinateur}} = 0 & \theta_{\text{spanish}, \text{ordinateur}} = 0 \end{array}$$

Similarly, if we want to distinguish positive and negative sentiment, we could use positive and negative *sentiment lexicons*, which are defined by expert psychologists (Tausczik and Pennebaker, 2010). You'll try this in Project 1.

But it's usually not easy to set the weights by hand. Instead, we will learn them from data. For example, suppose that an email user has manually labeled thousands of messages as "spam" or "not spam"; or a newspaper may label its own articles as "business" or "fashion." Such **instance labels** are a typical form of labeled data that we will encounter in NLP. In **supervised machine learning**, we use instance labels to automatically set the weights for a classifier. An important tool for this is probability.

3.1 Review of basic probability

This section is inspired/borrowed from Manning and Schütze (1999).

- **Formally:** When we write $P(\cdot)$, this denotes a function $P : \mathcal{F} \rightarrow [0, 1]$ from an **event space** \mathcal{F} to a **probability**. A probability is a real number between zero and one, with zero representing impossibility and one representing certainty.
- The probabilities of disjoint event sets are additive: $A_i \cap A_j = \emptyset \Rightarrow P(A_i \cup A_j) = P(A_i) + P(A_j)$. This is a restatement of the Third Axiom of probability.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- For example, you might ask what is the probability of two heads on three coin flips. There are eight possible series of three flips HHH, HHT, \dots , and each is an equally likely event. Of these events, three meet the criterion, HHT, HTH, THH . So the probability is $\frac{3}{8}$.
- More generally, $P(A_i \cup A_j) = P(A_i) + P(A_j) - P(A_i \cap A_j)$. This can be derived from the third axiom.

$$P(A_i \cup A_j) = P(A_i) + P(A_j - (A_i \cap A_j)) \quad (3.7)$$

$$P(A_j) = P(A_j - (A_i \cap A_j)) + P(A_i \cap A_j) \quad (3.8)$$

$$P(A_j - (A_i \cap A_j)) = P(A_j) - P(A_i \cap A_j) \quad (3.9)$$

$$P(A_i \cup A_j) = P(A_i) + P(A_j) - P(A_i \cap A_j) \quad (3.10)$$

- If the probability $P(A \cap B) = P(A)P(B)$, then the events A and B are *independent*, written $A \perp B$.

Conditional probability and Bayes' Rule

A conditional probability is an expression like $P(A | B)$, where we are interested in the probability of A conditioned on B happening.

- Conditional probability: $P(A | B) = P(A \cap B) / P(B)$
- If $P(A \cap B | C) = P(A | C)P(B | C)$, then the events A and B are **conditionally independent**, written $A \perp B | C$.
- Chain rule: $P(A \cap B) = P(A | B)P(B)$, which is just a rearrangement of terms.
- We can apply the chain rule multiple times:

$$\begin{aligned} P(A \cap B \cap C) &= P(A | B \cap C)P(B \cap C) \\ &= P(A | B \cap C)P(B | C)P(C) \end{aligned}$$

We'll do this a lot later in the course.

- Bayes' rule follows from the Chain rule: $P(A | B) = P(A \cap B) / P(B) = P(B | A)P(A) / P(B)$

Often we want the maximum a posteriori (MAP) estimate

$$\begin{aligned}\hat{B} &= \arg \max_B P(B \mid A) \\ &= \arg \max_B P(A \mid B)P(B)/P(A) \\ &\propto \arg \max_B P(A \mid B)P(B)\end{aligned}$$

- We don't need to normalize the probability because $P(A)$ is the same for all values of B .
- If we do need to compute the conditional $P(A \mid B)$, we can compute $P(A)$ by summing over $P(A \cap B) + P(A \cap \overline{B})$, where $B \cap \overline{B} = \emptyset$ and $B \cup \overline{B} = \Omega$, the entire sample space (such that $P(\Omega) = 1$).
- More generally, if $\bigcup_i B_i = \Omega$ and $\forall_{i,j}, B_i \cap B_j = \emptyset$, then $P(A) = \sum_i P(A \mid B_i)P(B_i)$.

Example Manning and Schütze (1999) have a nice example of Bayes Rule (Bayes Law) in a linguistic setting.

- Suppose one is interested in a rare syntactic construction, perhaps parasitic gaps, which occurs on average once in 100,000 sentences.
 - (An example of a sentence with a parasitic gap is *Which class did you attend __ without registering for __?* -JE)
- Lana Linguist has developed a complicated pattern matcher that attempts to identify sentences with parasitic gaps. Its pretty good, but it's not perfect:
 - If a sentence has a parasitic gap, it will say so with probability 0.95 (this is the **recall** -JE).
 - If it doesn't, it will wrongly say it does with probability 0.005 (this is the **false positive rate**, the additive inverse of **precision** -JE).
- Suppose the test says that a sentence contains a parasitic gap. What is the probability that this is true?
- (This example is usually framed in terms of tests for rare diseases. -JE)

(c) Jacob Eisenstein 2014-2015. Work in progress.

Solution: Let G be the event of a sentence having a parasitic gap, and T be the event of the test being positive.

$$P(G | T) = \frac{P(G | T)P(T)}{P(G | T)P(T) + P(G | \bar{T})P(\bar{T})} \quad (3.11)$$

$$= \frac{0.95 \times 0.00001}{0.95 \times 0.00001 + 0.005 \times 0.99999} \approx 0.002 \quad (3.12)$$

Random variables

A random variable takes on a specific value in \mathbb{R}^n , typically with $n = 1$, but not always. Discrete random variables can take values only in some countable subset of \mathbb{R} .

- Recall the coin flip example. The number of heads, H , can be viewed as a discrete random variable, $H \in 0, 1, 2, 3$.
- The probability mass associated with each number is $\{\frac{1}{8}, \frac{3}{8}, \frac{3}{8}, \frac{1}{8}\}$.
- This set of numbers represents the **probability distribution** over H , written $P(H = h) = p(h)$.
- To indicate that the RV H is distributed as $p(h)$, we write $H \sim p(h)$.
- The function $p(h)$ is called a probability **mass** function (pmf) if h is discrete, and a probability **density** function (pdf) if h is continuous.
- If we have more than one variable, we can write a joint probability $p(a, b) = P(A = a, B = b)$.
- We can write a **marginal** probability $p_A(a) = \sum_b p(a, b)$.
- Random variables are independent iff $p_{A,B}(a, b) = p_A(a)p_B(b)$.
- We can write a conditional probability as $p(a | b) = \frac{p(a,b)}{p_B(b)}$.

Expectations

Sometimes we want the **expectation** of a function, such as $E[g(x)] = \sum_{x \in \mathcal{X}} g(x)p(x)$.

Expectations are easiest to think about in terms of probability distributions over discrete events:

- If it is sunny, Marcia will eat three ice creams.
- If it is rainy, she will eat only one ice cream.
- There's a 80% chance it will be sunny.
- The expected number of ice creams she will eat is $0.8 \times 3 + 0.2 \times 1 = 2.6$.

If the random variable X is continuous, the sum becomes an integral:

$$E[g(x)] = \int_{\mathcal{X}} g(x)p(x)dx \quad (3.13)$$

For example, a fast food restaurant in Quebec gives a 1% discount on french fries for every degree below zero. Assuming they used a thermometer with infinite precision, the expected price would be an integral over all possible temperatures.

3.2 Naïve Bayes

Back to classification! A Naïve Bayes classifier chooses the weights θ to maximize the *joint* probability of a labeled dataset, $p(\mathbf{x}_{1:N}, \mathbf{y}_{1:N})$, where $\langle \mathbf{x}_i, y_i \rangle$ is a labeled instance.

We first need to define the probability $p(\mathbf{x}, y)$. We'll do that through a "generative model," which describes a hypothesized stochastic process that has generated the observed data.¹

- For each document i ,
 - draw the label $y_i \sim \text{Categorical}(\mu)$
 - draw the vector of counts $\mathbf{x}_i \sim \text{Multinomial}(\phi_{y_i})$,

¹We'll see a lot of different generative models in this course. They are a helpful tool because they clearly and explicitly define the assumptions that underly the form of the probability distribution.

The first thing this generative model tells us is that we can treat each document independently: the probability of the whole dataset is equal to the product of the probabilities of each individual document. The observed word counts and document labels are independent and identically distributed (IID).

$$p(\mathbf{x}, \mathbf{y}; \mu, \phi) = \prod_i p(\mathbf{x}_i, y_i; \mu, \phi) \quad (3.14)$$

This means that the words in each document are **conditionally independent** given the parameters μ and ϕ .

When we write $y_i \sim \text{Categorical}(\mu)$, that means y_i is a stochastic draw from a categorical distribution with **parameter** μ . A categorical distribution is just like a weighted die: $p_{\text{cat}}(y; \mu) = \mu_y$, where μ_y is the probability of the outcome $Y = y$. We require $\sum_y \mu_y = 1$ and $\forall_y, \mu_y \geq 0$.

A multinomial distribution is only slightly more complex:

$$p_{\text{mult}}(\mathbf{x}; \phi) = \frac{(\sum_j x_j)!}{\prod_j x_j!} \prod_j \phi_j^{x_j} \quad (3.15)$$

We again require that $\sum_j \phi_j = 1$ and $\forall_j, \phi_j \geq 0$. The first part of the equation doesn't depend on ϕ , and can usually be ignored. Can you see why we need the first part at all?²

We can write $p(\mathbf{x}_i | y_i; \phi)$ to indicate the conditional probability of word counts \mathbf{x}_i given label y_i , with parameter ϕ , which is equal to $p_{\text{mult}}(\mathbf{x}_i; \phi_{y_i})$.

By specifying the multinomial distribution, we are working with *multinomial naïve Bayes* (MNB). Why “naïve”? Because the multinomial distribution treats each word token independently: the probability mass function factorizes across the counts.³ We'll see this more clearly later, when we show how MNB is an example of linear classification.

Another version of Naïve Bayes

Consider a slight modification to the generative story of NB:

²Technically, a multinomial distribution requires a second parameter, the total number of counts (the number of words in the document). Even more technically, that number should be treated as a random variable, and drawn from some other distribution. But none of that matters for classification.

³You can plug in any probability distribution to the generative story and it will still be naïve Bayes, as long as you are making the “naïve” assumption that your features are generated independently.

- For each document i
 - Draw the label $y_i \sim \text{Categorical}(\mu)$
 - For each word $n \leq D_i$
 - * Draw the word $w_{i,n} \sim \text{Categorical}(\phi_{y_i})$

This is not quite the same model as multinomial Naive Bayes (MNB): it's a product of categorical distributions over words, instead of a multinomial distribution over word counts. This means we would generate the words in order, like $p_W(\text{multinomial})p_W(\text{Naive})p_W(\text{Bayes})$. Formally, this is a model for the joint probability $p(\mathbf{w}, y)$, not $p(\mathbf{x}, y)$.

However, as a classifier, it is identical to MNB. The final probabilities are reduced by a factor corresponding to the normalization term in the multinomial, $\frac{(\sum_j x_j)!}{\prod_j x_j!}$. This means that the resulting probabilities for a given \mathbf{x} are different. However, none of this has anything to do with the label y or the parameters ϕ . The ratio of probabilities between any two labels y_1 and y_2 will be identical, as will the maximum likelihood estimates for the parameters μ and ϕ (defined later).

Prediction

The Naive Bayes prediction rule is to choose the label y which maximizes $p(\mathbf{x}, y; \phi, \mu)$:

$$\begin{aligned}
 \hat{y} &= \arg \max_y p(\mathbf{x}, y; \mu, \phi) \\
 &= \arg \max_y p(\mathbf{x} \mid y; \phi) p(y; \mu) \\
 &= \arg \max_y \log p(\mathbf{x} \mid y; \phi) + \log p(y; \mu)
 \end{aligned}$$

Converting to logarithms makes the notation easier. It doesn't change the prediction rule because the log function is monotonically increasing.

Now we can plug in the probability distributions from the generative story.

$$\begin{aligned}
\log p(\mathbf{x}, y; \mu, \phi) &= \arg \max_y \log p(\mathbf{x} \mid y; \phi) + \log p(y; \mu) \\
&= \log \left[\frac{(\sum_j x_j)!}{\prod_j x_j!} \prod_j \phi_{y,j}^{x_j} \right] + \log \mu_y \\
&= \log \frac{(\sum_j x_j)!}{\prod_j x_j!} + \sum_j x_j \log \phi_{y,j} + \log \mu_y \\
&\propto \sum_j x_j \log \phi_{y,j} + \log \mu_y \\
&= \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y),
\end{aligned}$$

where

$$\begin{aligned}
\boldsymbol{\theta} &= [\boldsymbol{\theta}^{(1)\top}, \boldsymbol{\theta}^{(2)\top}, \dots, \boldsymbol{\theta}^{(K)\top}]^\top \\
\boldsymbol{\theta}^{(y)} &= [\log \phi_{y,1} \ \log \phi_{y,2} \ \dots \ \log \phi_{y,M} \ \log \mu_y]^\top
\end{aligned}$$

and $\mathbf{f}(\mathbf{x}, y)$ is a vector of word counts and an offset, padded by zeros for the labels not equal to y (see equations 3.1-3.5). This ensures that the inner product $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$ only activates the features in $\boldsymbol{\theta}^{(y)}$, which are what we need to compute the joint log-probability $\log p(\mathbf{x}, y)$ for each y .

Estimation

The parameters of a multinomial distribution have a simple interpretation: they're the expected frequency for each word. Based on this interpretation, it's tempting to set the parameters empirically, as

$$\phi_{y,j} = \frac{\sum_{i:Y_i=y} x_{i,j}}{\sum_{j'} \sum_{i:Y_i=y} x_{i,j'}} = \frac{\text{count}(y, j)}{\sum_{j'} \text{count}(y, j')} \quad (3.16)$$

In NLP this is called a *relative frequency estimator*. It can be justified more rigorously as a *maximum likelihood estimate*.

As in prediction, we want to maximize the joint likelihood of the data,

$$L = \sum_i \log p_{\text{mult}}(\mathbf{x}_i; \boldsymbol{\phi}_{y_i}) + \log p_{\text{cat}}(y_i; \mu) \quad (3.17)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

Since $p(y)$ is unrelated to ϕ , we can forget about it for now. But before we can just optimize L , we have to deal with a constraint:

$$\sum_j \phi_{y,j} = 1 \quad (3.18)$$

We'll do this by adding a Lagrange multiplier. Here's the resulting Lagrangian:

$$\ell[\phi_y] = \sum_{i:Y_i=y} \sum_j x_{ij} \log \phi_{y,j} + \lambda \left(\sum_j \phi_{y,j} - 1 \right) \quad (3.19)$$

We solve by setting $\frac{\partial \ell}{\partial \phi_j} = 0$.

$$\begin{aligned} 0 &= \sum_{i:Y_i=y} x_{i,j} / \phi_{y,j} - \lambda \\ \lambda \phi_{y,j} &= \sum_{i:Y_i=y} x_{i,j} \\ \phi_{y,j} &\propto \sum_{i:Y_i=y} x_{i,j} = \sum_i \delta(Y_i = y) x_{i,j} \\ &= \frac{\sum_{i:Y_i=y} x_{i,j}}{\sum_{j'} \sum_{i:Y_i=y} x_{i,j'}} \end{aligned}$$

Similarly, $\mu_y \propto \sum_i \delta(Y_i = y)$, where $\delta(Y_i = y) = 1$ if $Y_i = y$ and 0 otherwise.

Smoothing and MAP estimation

If data is sparse, you can end up with values of $\phi = 0$, allowing a single feature to completely veto a label. This is undesirable, because it imposes high **variance**: depending on what data happens to be in the training set, we could get vastly different classification rules.

One solution is Laplace smoothing: adding “pseudo-counts” of α to each estimate, and then normalize.

$$\phi_{y,j} = \frac{\alpha + \sum_{i:Y_i=y} x_{i,j}}{\sum_{j'} \alpha + \sum_{i:Y_i=y} x_{i,j'}} = \frac{\alpha + \text{count}(i, j)}{V\alpha + \sum_{j'} \text{count}(i, j')} \quad (3.20)$$

Laplace smoothing has a nice Bayesian justification, in which we extend the generative story to include ϕ as a random variable (rather than as a parameter). The resulting estimate is called *maximum a posteriori*, or MAP.

Smoothing reduces **variance**, but it takes us away from the maximum-likelihood estimate: it imposes a **bias** (towards uniform probabilities). Machine learning theory shows that errors on held out data result from the sum of bias and variance. Techniques for reducing variance typically increase the bias, so there is a **bias-variance tradeoff**.

- Unbiased classifiers **overfit** the training data, yielding poor performance on unseen data.
- But if we set a very large smoothing value, we can **underfit** instead. In the limit of $\alpha \rightarrow \infty$, we have zero variance: it is the same classifier no matter what data we see! But the bias of such a classifier will be high.
- Navigating this tradeoff is hard. But in general, as you have more data, variance is less of a problem, so you just go for low bias.

Training, testing, and tuning (development) sets

We'll soon talk about more learning algorithms, but whichever one we apply, we will want to report its accuracy. Really, this is an educated guess about how well the algorithm will do on new data in the future.

To do this, we need to hold out a separate “test set” from the data that we use for estimation (i.e., training, learning). Otherwise, if we measure accuracy on the same data that is used for estimation, we will badly overestimate the accuracy we're likely to get on new data. See <http://xkcd.com/1122/> for a cartoon related to this idea.

Many learning algorithms also have “tuning” parameters:

- the smoothing pseudo-counts α in Naive Bayes
- the regularization λ in logistic regression
- the slack weight C in the support-vector machine

All of these tuning parameters really do the same thing: they navigate the bias-variance tradeoff. Where is the best position on this tradeoff curve? It's hard to tell in advance. Sometimes it is tempting to see which tuning parameter gives the best performance on the test set, and then report that performance. Resist this temptation! It will also lead to overestimating accuracy on truly unseen future

data. For that reason, this is a sure way to get your research paper rejected. Instead, you should split off a piece of your training data, called a “development set” (or “tuning set”).

Sometimes, people average across multiple test sets and/or multiple development sets. One way to do this is to divide your data into “folds,” and allow each fold to be the development set one time. This is called **K-fold cross-validation**. In the extreme, each fold is a single data point. This is called **leave-one-out**.

The Naivety of Naive Bayes

Naive Bayes is very simple to work with. Estimation and prediction can be done in closed form, and the nice probabilistic interpretation makes it relatively easy to extend the model in various ways.

But Naive Bayes makes assumptions which seriously limit its accuracy, especially in NLP.

- The multinomial distribution assumes that each word is generated independently of all the others (conditioned on the parameter ϕ_y). Formally, we assume conditional independence:

$$p(\text{naïve}, \text{Bayes}; \phi) = p(\text{naïve}; \phi)p(\text{Bayes}; \phi). \quad (3.21)$$

- But this is clearly wrong, because words “travel together.” Question for you, is it:

$$p(\text{naïve Bayes}) > p(\text{naïve})p(\text{Bayes}) \quad (3.22)$$

or...

$$p(\text{naïve Bayes}) < p(\text{naïve})p(\text{Bayes}) \quad (3.23)$$

Apply the chain rule!

Traffic lights Dan Klein makes this point with an example about traffic lights. In his hometown of Pittsburgh, there is a 1/7 chance that the lights will be broken, and both lights will be red. There is a 3/7 chance that the lights will work, and the north-south lights will be green; there is a 3/7 chance that the lights work and the east-west lights are green.

The *prior* probability that the lights are broken is 1/7. If they are broken, the conditional likelihood of each light being red is 1. The prior for them not being broken is 6/7. If they are not broken, the conditional likelihood of each being light being red is 1/2.

Now, suppose you see that both lights are red. According to Naive Bayes, the probability that the lights are broken is $1/7 \times 1 \times 1 = 1/7 = 4/28$. The probability that the lights are not broken is $6/7 \times 1/2 \times 1/2 = 6/28$. So according to naive Bayes, there is a 60% chance that the lights are not broken!

What went wrong? We have made an independence assumption to factor the probability $P(R, R \mid \text{not-broken}) = P_{\text{north-south}}(R \mid \text{not-broken})P_{\text{east-west}}(R \mid \text{not-broken})$. But this independence assumption is clearly incorrect, because $P(R, R \mid \text{not-broken}) = 0$.

Less Naive Bayes? Of course we could decide not to make the naive Bayes assumption, and model $P(R, R)$ explicitly. But this idea does not scale when the feature space is large (as it often is in NLP). The number of possible feature configurations grows exponentially, so our ability to estimate accurate parameters will suffer from high variance. With an infinite amount of data, we'd be fine (in theory, maybe not in practice); but we never have that. Naive Bayes accepts some bias (because of the incorrect modeling assumption) in exchange for lower variance.

3.3 Recap

- Bag-of-words representation $\mathbf{f}(\mathbf{x}, y)$
- Classification as a dot-product $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$
- Naive Bayes
 - Define $p(\mathbf{x}, y)$ via a *generative model*
 - Prediction: $\hat{y} = \arg \max_y p(\mathbf{x}_i, y)$
 - Learning:

$$\begin{aligned}\boldsymbol{\theta} &= \arg \max_{\boldsymbol{\theta}} p(\mathbf{x}, y; \boldsymbol{\theta}) \\ p(\mathbf{x}, y; \boldsymbol{\theta}) &= \prod_i p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) = \prod_i p(\mathbf{x}_i | y_i) p(y_i) \\ \phi_{y,j} &= \frac{\sum_{i: Y_i=y} x_{ij}}{\sum_{i: Y_i=y} \sum_j x_{ij}} \\ \mu_y &= \frac{\text{count}(Y = y)}{N}\end{aligned}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

This gives the maximum-likelihood estimator (MLE; same as relative frequency estimator)

- Bias-variance tradeoff: MLE is high-variance, so add smoothing pseudo counts α . This reduces variance but adds bias.

Chapter 4

Sentiment analysis

Todo: add notes about sentiment analysis here

Chapter 5

Discriminative learning

5.1 Features

Naive Bayes is a simple classifier, where the weights are learned based on the joint probability of labels and words. It includes an independence assumption: all features are mutually independent, conditioned on the label.

- We have defined a **feature function** $f(x, y)$, which corresponds to “bag-of-words” features. While these features do violate the independence assumption, the violation is relatively mild.
- We may be interested in other features, which violate independence more severely. Can you think of any?
 - Prefixes, e.g. *anti-*, *im-*, *un-*
 - Punctuation and capitalization
 - Bigrams, e.g. *not good*, *not bad*, *least terrible*, ...

Rich feature sets generally cannot be combined with Naive Bayes because the distortions resulting from violations of the independence assumption overwhelm the additional power of better features.

$$p(\text{not bad food}|y) \approx p(\text{not}|y)p(\text{bad}|y)p(\text{food}|y) \quad (5.1)$$

$$p(\text{not bad food}|y) \not\approx p(\text{not}|y)p(\text{bad}|y)p(\text{not bad}|y)p(\text{food}|y) \quad (5.2)$$

To use these features, we will need learning algorithms that do not rely on an independence assumption.

5.2 Perceptron

In NB, the weights can be interpreted as parameters of a probabilistic model. But this model requires an independence assumption that usually does not hold, and limits our choice of features.

Why not forget about probability and learn the weights in an error-driven way?

- Until converged, at each iteration t
 - Select an instance i
 - Let $\hat{y} = \arg \max_y \boldsymbol{\theta}_t^\top \mathbf{f}(\mathbf{x}_i, y)$
 - If $\hat{y} = y_i$, do nothing
 - If $\hat{y} \neq y_i$, set $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$

Basically we are saying: if you make a mistake, increase the weights for features which are active with the correct label y_i , and decrease the weights for features which are active with the guessed label \hat{y} .

This seems like a cheap heuristic, right? Will it really work? In fact, there is some nice theory for the perceptron.

- If there is a set of weights that correctly separates your data, then your data is **separable**.
- Formally, your data is (linearly) separable if there exists a set of weights $\boldsymbol{\theta}$ such that

$$\forall \mathbf{x}_i, y_i, \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) > \max_{y' \neq y_i} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y') \quad (5.3)$$

- If your data is linearly separable, it can be proven that the perceptron algorithm will eventually find a separator.
- What if your data is not separable?
 - the number of errors is bounded...
 - but the algorithm will thrash. That is, the weights will cycle between different values, and will never converge.

The perceptron is an **online** learning algorithm.

- This means that it adjusts the weights after every example.

- This is different from Naïve Bayes, which computes corpus statistics and then sets the weights in a single operation. This is a **batch learning** algorithm.
- Other algorithms are **iterative**, in that they perform multiple updates to the weights, but are also **batch**, in that they have to use all the training data to compute the update. We'll mention two of those algorithms later.

Voted (averaged) perceptron

One solution to thrashing is to average the weights across all iterations:

$$\bar{\theta} = \frac{1}{T} \sum_t \theta_t$$

$$y = \arg \max_y \bar{\theta}^\top f(\mathbf{x}, y)$$

There is some analysis showing that voting can improve generalization (Freund and Schapire, 1999; Collins, 2002). However, this rule as described here is not practical. Can you see why not, and how to fix it?

5.3 Loss functions and large-margin classification

Naive Bayes chooses the weights θ by maximizing the likelihood $p(\mathbf{x}, \mathbf{y})$. This can be seen, equivalently, as maximizing the log-likelihood (due to the monotonicity of the log function), and as **minimizing** the negative log-likelihood. This negative log-likelihood can therefore be viewed as a **loss function**, which is minimized:

$$\log p(\mathbf{x}, \mathbf{y}; \theta) = \sum_i \log p(\mathbf{x}_i, y_i; \theta) \quad (5.4)$$

$$\ell_{\text{NB}}(\theta; \mathbf{x}_i, y_i) = -\log p(\mathbf{x}_i, y_i; \theta) \quad (5.5)$$

$$\hat{\theta} = \arg \min_{\theta} \sum_i \ell_{\text{NB}}(\theta, \mathbf{x}_i, y_i) \quad (5.6)$$

This may seem confusing and backwards, but loss functions provide a very general framework in which to compare many approaches to machine learning. For example, even though the perceptron is not a probabilistic model, it is also trying to minimize a **loss function**:

$$\ell_{\text{perceptron}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = \begin{cases} 0, & y_i = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y) \\ 1, & \text{otherwise} \end{cases} \quad (5.7)$$

This loss function has some pros and cons in comparison with Naive Bayes.

- ℓ_{NB} can suffer **infinite** loss on a single example, which suggests it will overemphasize some examples, and underemphasize others.
- $\ell_{\text{perceptron}}$ treats all errors equally. It only cares if the example is correct, and not about how confident the classifier was. Since we usually evaluate on accuracy, this is a better match.
- $\ell_{\text{perceptron}}$ is non-convex¹ and discontinuous. Finding the global optimum is intractable when the data is not separable.

We can fix this last problem by defining a loss function that behaves more nicely. To do this, let's define the *margin* as

$$\gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \max_{y \neq y_i} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y) \quad (5.8)$$

Then we can write a convex and continuous “hinge loss” as

$$\ell_{\text{hinge}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = \begin{cases} 0, & \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i) \geq 1, \\ 1 - \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i), & \text{otherwise} \end{cases} \quad (5.9)$$

Equivalently, we can write $\ell_{\text{hinge}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = (1 - \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i))_+$, where $(x)_+$ indicates the positive part of x .

Essentially, we want a *margin* of at least 1 between the score for the true label and the best-scoring alternative, which we have written \hat{y} .

The hinge and perceptron loss functions are shown in Figure 5.1.

Large-margin online classification

Note that we can write $\boldsymbol{\theta} = s\mathbf{u}$, where $\|\mathbf{u}\|_2 = 1$. Think of s as the magnitude and \mathbf{u} as the direction of the vector $\boldsymbol{\theta}$. If the data is separable, there are many values

¹As a reminder, a function f is convex iff $\alpha f(x_i) + (1 - \alpha)f(x_j) \geq f(\alpha x_i + (1 - \alpha)x_j)$, for all $\alpha \in [0, 1]$ and for all x_i and x_j on the domain of the function. Convexity implies that any local minimum is also a global minimum, and there are a wide array of techniques for optimizing convex functions (Boyd and Vandenberghe, 2004)



Figure 5.1: Hinge and perceptron loss functions

of s which attain zero hinge loss. For generality, we will try to make the smallest magnitude change to θ possible.²

At step t , we optimize:

$$\theta_{t+1} = \arg \min_{\theta} \frac{1}{2} \|\theta - \theta_t\|^2 \text{ s.t. } \ell_{\text{hinge}}(\theta; \mathbf{x}_i, y_i) = 0 \quad (5.10)$$

Assuming that the constraint can be satisfied (i.e., the problem is linearly separable), the optimal solution is found at,

$$\theta_{t+1} = \theta_t + \tau_t (\mathbf{f}(y_i, \mathbf{x}_i) - \mathbf{f}(\hat{y}, \mathbf{x}_i)) \quad (5.11)$$

$$\tau_t = \frac{\ell(\theta; \mathbf{x}_i, y_i)}{\|\mathbf{f}(x_i, y_i) - \mathbf{f}(x_i, \hat{y})\|^2}, \quad (5.12)$$

where again \hat{y} is the best scoring y according to θ_t . This solution can be obtained by introducing τ_t as a Lagrange multiplier for the constraint in (5.10).

²In the support vector machine (without slack variables), we choose the smallest magnitude weights that satisfy the constraint of zero hinge loss. Pegasos is an online algorithm for training SVMs (Shwartz et al., 2007); it is similar to Passive-Aggressive.

If the data is not linearly separable, there will be instances for which we can't meet this constraint. To deal with this, we introduce a “slack” variable ξ_i . We use the slack variable to trade off between the constraint (having a large margin) and the objective (having a small change in θ). The tradeoff is controlled by a parameter C .

$$\begin{aligned} \min w \frac{1}{2} \|\theta - \theta_t\|^2 + C\xi_t \\ \text{s.t. } \ell_{\text{hinge}}(\theta; \mathbf{x}_i, y_i) \leq \xi_t, \xi_t \geq 0 \end{aligned} \quad (5.13)$$

The solution to 5.13 is,

$$\theta_{t+1} = \theta_t + \tau_t (\mathbf{f}(y_i, \mathbf{x}_i) - \mathbf{f}(\hat{y}, \mathbf{x}_i)) \quad (5.14)$$

$$\tau_t = \min \left(C, \frac{\ell(\theta; \mathbf{x}_i, y_i)}{\|\mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})\|^2} \right), \quad (5.15)$$

- If C is 0, then infinite slack is permitted, and the weights will never change.
- As $C \rightarrow \infty$, no slack is permitted, and the optimization is identical to equation 5.10 and 5.12.

This algorithm is called “Passive-Aggressive” (PA; Crammer et al., 2006), because it is passive when the margin constraint is satisfied, but it aggressively changes the weights to satisfy the constraints if necessary.³

- PA is error-driven like the perceptron, but is more stable to violations of separability, like the averaged perceptron.
- PA allows more explicit control than the Averaged Perceptron, due to the C parameter. When C is small, we make very conservative adjustments to θ from each instance, because the slack variables aren't very expensive. When C is large, we make large adjustments to avoid using the slack variables.
- You can also apply weight averaging to PA.
- **Support vector machines** (SVMs) are another learning algorithm based on the hinge loss (Burges, 1998), but they try to minimize the norm of the weights, rather than the norm of the change in the weights. They are typically trained

³A related algorithm without slack variables is called MIRA, for Margin-Infused Relaxed Algorithm (Crammer and Singer, 2003).

in **batch** style, meaning that they have to read all the training instances in to compute each update. However, SVMs can also be trained in an online fashion (Shwartz et al., 2007). The LXMLS lab guide provides a simpler online learning algorithm, based on stochastic subgradient descent (Figueiredo et al., 2013).

Pros and cons of Perceptron and PA

- Perceptron and PA are error-driven, which means they usually do better in practice than naive Bayes.
- They are also online, which means we can learn without having our whole dataset in memory at once. NB can also be estimated online, in the sense that you can stream the data and store the counts.
- The original perceptron doesn't behave well if the data is not separable, and doesn't make it easy to control model complexity.
- All these models lack a probabilistic interpretation. Probabilities are useful because they quantify the classification certainty, allowing us to compute expected utility, and to incorporate the classifier in more complex probabilistic models.

5.4 Logistic regression

Logistic regression is error-driven like the perceptron, but probabilistic like Naive Bayes. This is useful in case we want to quantify the uncertainty about a classification decision.

Recall that NB selects weights to optimize the joint probability $p(y, \mathbf{x})$.

- In NB, we factor this as $p(y, \mathbf{x}) = p(\mathbf{x}|y)p(y)$.
- But we could equivalently write $p(y, \mathbf{x}) = p(y|\mathbf{x})p(\mathbf{x})$.

Since we always know \mathbf{x} , we really care only about $p(y|\mathbf{x})$. Logistic regression optimizes this directly. To do this, we have to define the probability function

differently. We define the conditional probability directly, as,

$$p(y|\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y))}{\sum_{y'} \exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y'))} \quad (5.16)$$

$$\log p(y|\mathbf{x}) = \sum_i \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y'} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y') \quad (5.17)$$

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_i \log p(y_i|\mathbf{x}_i; \boldsymbol{\theta}) \quad (5.18)$$

Inside the sum, we have the (additive inverse of) the **logistic loss**.

- In binary classification, we can write this as

$$\ell_{\text{logistic}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = -(y_i \boldsymbol{\theta}^\top \mathbf{x}_i - \log(1 + \exp \boldsymbol{\theta}^\top \mathbf{x}_i)) \quad (5.19)$$

- In multi-class classification, we have,⁴

$$\ell_{\text{logistic}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = -(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y'} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y')) \quad (5.20)$$

The logistic loss is shown in Figure 5.2. Because it is smooth and convex, we can optimize it through gradient steps:

$$\ell = \sum_i \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y'} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y') \quad (5.21)$$

$$\frac{\partial \ell}{\partial \boldsymbol{\theta}} = \sum_i \mathbf{f}(\mathbf{x}_i, y_i) - \frac{\sum_{y'} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y') \mathbf{f}(\mathbf{x}_i, y')}{\sum_{y''} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y'')} \quad (5.22)$$

$$= \sum_i \mathbf{f}(\mathbf{x}_i, y_i) - \sum_{y'} \frac{\exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y')}{\sum_{y''} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y'')} \mathbf{f}(\mathbf{x}_i, y') \quad (5.23)$$

$$= \sum_i \mathbf{f}(\mathbf{x}_i, y_i) - \sum_y' p(y'|\mathbf{x}_i; \boldsymbol{\theta}) \mathbf{f}(\mathbf{x}_i, y') \quad (5.24)$$

$$= \sum_i \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y')] \quad (5.25)$$

⁴The log-sum-exp term is very common in machine learning. It is numerically instable because you can underflow if the inner product is small, and overflow if the inner product is large. Libraries like `scipy` contain special functions for computing `logsumexp`, but with some thought, you should be able to see how to create an implementation that is numerically stable.

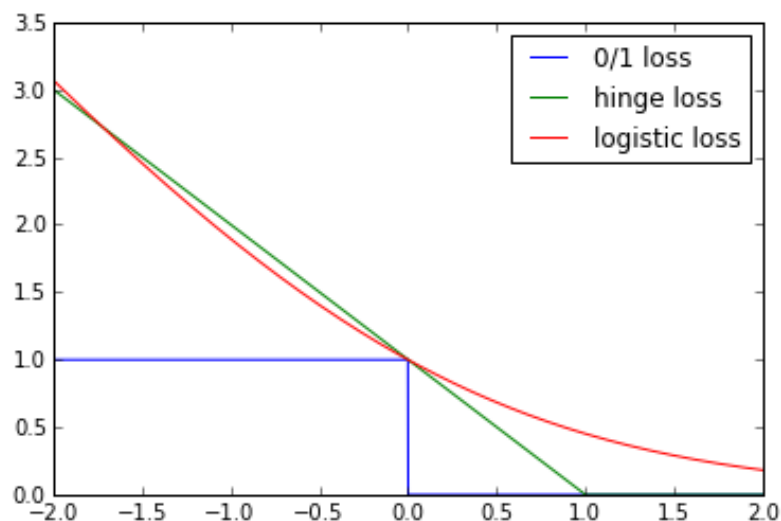


Figure 5.2: Hinge, perceptron, and logistic loss functions

This gradient has a very pleasing interpretation as the difference between the observed counts and the expected counts.⁵ Compare this gradient with the perceptron and PA update rules.

The bias-variance tradeoff is handled by penalizing large θ in the objective, adding a term of $\frac{\lambda}{2} \|\theta\|_2^2$. This is called L2 regularization, because of the L2 norm. It can be viewed as placing a 0-mean Gaussian prior on θ .

This penalty contributes a term of $\lambda\theta$ to the gradient, so we have,

$$\ell = \sum_i \theta^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y'} \exp \theta^\top \mathbf{f}(\mathbf{x}_i, y') + \frac{\lambda}{2} \|\theta\|_2^2$$

$$\frac{\partial \ell}{\partial \theta} = \sum_i \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y')] - \lambda \theta.$$

Optimization

Batch optimization In batch optimization, you keep all the data in memory and iterate over it many times.

⁵Recall that the definition of an expected value $E[f(x)] = \sum_x f(x)p(x)$

- The logistic loss is smooth and convex, so we can find the global optimum using gradient descent. But in practice, this can be very slow.
- Second-order (Newton) optimization would incorporate the inverse Hessian. The Hessian is

$$H_{i,j} = \frac{\partial^2}{\partial w_i \partial w_j} \ell, \quad (5.26)$$

but this matrix is usually too big to deal with.

- In practice, people usually apply **quasi-Newton optimization**, which approximates the Hessian matrix. The specific method that is particularly popular is L-BFGS⁶ NLP people usually treat L-BFGS as a black box; you will typically pass it a pointer to a function that computes the likelihood and gradient. L-BFGS is provided in `scipy.optimize`.

Online optimization In online optimization, you consider one example (or a “mini-batch” of a few examples) at a time. *Stochastic gradient descent* makes a stochastic online approximation to the overall gradient:

$$\begin{aligned} \boldsymbol{\theta}^{(t+1)} &\leftarrow \boldsymbol{\theta}^{(t)} - \eta_t \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}^{(t)}, \mathbf{x}, \mathbf{y}) \\ &= \boldsymbol{\theta}^{(t)} - \eta_t (\lambda \boldsymbol{\theta}^{(t)} - \sum_i^N \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y')]) \\ &= (1 - \lambda \eta_t) \boldsymbol{\theta}^{(t)} + \eta_t \sum_i^N \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y')] \\ &\approx (1 - \lambda \eta_t) \boldsymbol{\theta}^{(t)} + N \eta_t (\mathbf{f}(\mathbf{x}_{i(t)}, y_{i(t)}) - E[\mathbf{f}(\mathbf{x}_{i(t)}, y')]) \end{aligned}$$

where η_t is the **stepsize** at iteration t , and $\langle \mathbf{x}_{i(t)}, y_{i(t)} \rangle$ is the instance selected at iteration t .

- Note how similar this update is to the perceptron!
- If we set $\eta_t = \eta_0 t^{-\alpha}$ for $\alpha \in [1, 2]$, we have guaranteed convergence.
- We can also just fix η_t to a small value, like 10^{-3} . (This is what we will do in the problem set.)

⁶A friend of mine told me you can remember the order of the letters as “Large Big Friendly Giants.” Does this help you?

- In either case, we could tune this parameter on a development set. However, it would be acceptable to just find a value that gives a good regularized log-likelihood on the training set, since this parameter relates to the quality of the optimization, and not the generalization capability of the classifier.
- In theory, we select $\langle \mathbf{x}_{i(t)}, y_{i(t)} \rangle$ at random, but in practice we usually just iterate through the dataset.
- We can fold N into η and λ , so that $\eta^* = N\eta$ and $\lambda^* = \lambda \frac{\eta^*}{N}$. This gives the more compact form,

$$(1 - \lambda^* \eta_t^*) \boldsymbol{\theta}^{(t)} + \eta_t^* (\mathbf{f}(\mathbf{x}_{i(t)}, y_{i(t)}) - E[\mathbf{f}(\mathbf{x}_{i(t)}, y')]) \quad (5.27)$$

For more on stochastic gradient descent, as applied to a number of different learning algorithms, see (Zhang, 2004) and (Bottou, 1998). Murphy (2012) traces SGD to a 1978 paper by GT's own Arkadi Nemirovski (Nemirovski and Yudin, 1978). You can find several recent chapters about online optimization in the edited volume by Sra et al. (2012).

Adagrad Recent work has shown that you can often learn more quickly by using an **adaptive** step-size, which is different for every feature (Duchi et al., 2011). Specifically, in the **Adagrad** algorithm (adaptive gradient), you keep track of the sum of the squares of the gradients for each feature, and rescale the learning rate by its inverse:

$$\mathbf{g}_t = -\mathbf{f}(\mathbf{x}_i, y_i) + \sum_{y'} p(y' | \mathbf{x}_i) \mathbf{f}(\mathbf{x}_i, y_i) + \lambda \boldsymbol{\theta} \quad (5.28)$$

$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \frac{\eta}{\sqrt{\sum_{t'} g_{t,j}^2}} g_{t,j}, \quad (5.29)$$

where j iterates over features in $\mathbf{f}(\mathbf{x}, y)$. The effect of this is that features with consistently large gradients are updated more slowly. Another way to view this update is that rare features are taken more seriously, since their sum of squared gradients will be smaller. Adagrad seems to require less careful tuning of η , and Dyer (2014) reports that $\eta = 1$ works for a wide range of problems.

Note that the Adagrad update can apply to any smooth loss function, including the hinge loss defined in Equation 5.9.

Names

Logistic regression is so named because in the binary case where $y \in \{0, 1\}$, we are performing a regression of x against y , after passing the inner product $\theta^\top x$ through a logistic transformation. You could always do a linear regression, but this would ignore the fact that the y is limited to a few values.

- Logistic regression is also called **maximum conditional likelihood** (MCL), because it maximizes... the conditional likelihood $p(y | x)$.
- Logistic regression can be viewed as part of a larger family, called **generalized linear models**. If you use R, you are probably familiar with `glmnet`.
- Logistic regression is also called **maximum entropy**, especially in the earlier NLP literature (Berger et al., 1996). This is due to an alternative formulation, which tries to find the maximum entropy probability function that satisfies moment-matching constraints.

The moment matching constraints specify that the empirical counts of each label-feature pair should match the expected counts:

$$\forall j, \sum_i f_j(x_i, y_i) = \sum_i \sum_y p(y | x_i; \theta) f_j(x_i, y) \quad (5.30)$$

Note that this constraint will be met exactly when the derivative of the likelihood function (equation 5.25) is equal to zero. However, this will be true for many values of θ . Which should we choose?

The entropy of a conditional likelihood function $P(Y|X)$ is

$$H(P) = - \sum_x \tilde{p}(x) \sum_y p(y|x) \log p(y|x), \quad (5.31)$$

where $\tilde{p}(x)$ is the *empirical probability* of x . We compute an empirical probability by summing over all the instances in training set.

If the entropy is large, this function is smooth across possible values of y ; if it is small, the function is sharp. The entropy is zero if $p(y|x) = 1$ for some particular $Y = y$ and zero for everything else. By saying we want maximum-entropy classifier, we are saying we want to make the least commitments possible, while satisfying the moment-matching constraints:

(c) Jacob Eisenstein 2014-2015. Work in progress.

$$\begin{aligned} \max_{\boldsymbol{\theta}} \quad & - \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}) \sum_y p(y|\mathbf{x}; \boldsymbol{\theta}) \log p(y|\mathbf{x}; \boldsymbol{\theta}) \\ \text{s.t.} \quad & \forall j, \sum_i f_j(\mathbf{x}_i, y_i) = \sum_i \sum_y p(y|\mathbf{x}_i; \boldsymbol{\theta}) f_j(\mathbf{x}_i, y) \end{aligned}$$

Now, the solution to this constrained optimization problem is identical to the maximum conditional likelihood (logistic-loss) formulation we've considered in the previous section.

This view of logistic regression is arguably a little dated, but it's useful to understand what's going on. The information-theoretic concept of entropy will pop up again a few times in the course. For a tutorial on maximum entropy, see <http://www.cs.cmu.edu/afs/cs/user/abberger/www/html/tutorial/tutorial.html>.

5.5 Summary of learning algorithms

- **Naive Bayes.** pros: easy and probabilistic. cons: arguably optimizes wrong objective; usually has poor accuracy, especially with overlapping features.
- **Perceptron and PA.** pros: easy, online, and error-driven. cons: not probabilistic. this can be bad in pipeline architectures, where the output of one system becomes the input for another.
- **Logistic regression.** pros: error-driven and probabilistic. cons: batch learning requires black-box software; hinge loss sometimes yields better accuracy than logistic loss.

What about non-linear classification?

The feature spaces that we consider in NLP are usually huge, so non-linear classification can be quite difficult. When the feature dimension V is larger than the number of instances N — often the case in NLP — you can always learn a linear classifier that will perfectly classify your training instances.⁷ This makes selecting an appropriate **non-linear** classifier especially difficult. Nonetheless, there are some approaches to non-linear learning in NLP:

⁷Assuming your feature matrix is full-rank.

- You can add **features**, such as bigrams, which are non-linear combinations of other features. For example, the base feature $\langle \text{coffee house} \rangle$ will not fire unless both features $\langle \text{coffee} \rangle$ and $\langle \text{house} \rangle$ also fire.
- Another option is to apply non-linear transformations to the feature vector. Recall that the feature function $f(x, y)$ may be composed of a vector of word counts, padded by zeros. We can think of these word counts as basic features, and apply non-linear transformations, such as $x \circ x$ or $|x|$.
- There is some work in NLP on using kernels for strings, bags-of-words, sequences, trees, etc. Kernelized learning algorithms are outside the scope of this class (Collins and Duffy, 2001; Zelenko et al., 2003). Kernel-based learning can be seen as a generalization of algorithms such k -nearest-neighbors, which classifies instances by considering the labels of the k most similar instances in the training set (Hastie et al., 2009).
- Boosting (Freund et al., 1999) and decision tree algorithms (Schmid, 1994) sometimes do well on NLP tasks, but they are used less frequently these days, especially as the field increasingly emphasizes big data and simple classifiers.
- More recent work has shown how **deep learning** can perform non-linear classification. One way to use deep learning in NLP is by learning word representations while jointly learning how these representations combine to classify instances (Collobert and Weston, 2008). This approach is very hot at the moment, so I will discuss it towards the end of the semester.

5.6 Summary of classifiers

So now we've talked about four different classifiers. That's it! No more classifiers in this class. Yay? Anyway, let's review.

	Naive Bayes	Logistic Regression	Perceptron	PA
Objective	Joint likelihood	Conditional likelihood	0-1 loss	Hinge loss
estimation	$\max \sum_i \log \mathbf{p}(\mathbf{x}_i, y_i)$	$\max \sum_i \log \mathbf{p}(y_i \mathbf{x}_i)$	$\min \sum_i \delta(y_i, \hat{y})$	$\sum_i [1 - \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i)] +$
tuning	$\theta_{ij} = \frac{c(\mathbf{x}_i, y=j) + \alpha}{c(y=j) + V\alpha}$	$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \sum_i \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y)]$	$\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$	$\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \tau_k(\mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y}))$
complexity	smoothing α	regularizer $\lambda \ \boldsymbol{\theta}\ _2^2$	weight averaging	slack penalty C
easy?	$\mathcal{O}(NV)$	$\mathcal{O}(NVT)$	$\mathcal{O}(NVT)$	$\mathcal{O}(NVT)$
probabilities?	very	not really	yes	yes
features?	yes	yes	no	no
	no	yes	yes	yes

Table 5.1: Comparison of classifiers. N = number of examples, V = number of features, T = number of instances.

Chapter 6

Word-sense disambiguation

Todo: add notes about WSD here

Chapter 7

Learning without supervision

So far we've assumed the following setup:

- A **training set** where you get observations x_i and labels y_i
- A **test set** where you only get observations x_i

What if you never get labels y_i ?

For example, you get a bunch of text, and you suspect that there are at least two different meanings for the word *concern*.¹

The immediate context includes two groups of words:

- services, produces, banking, pharmaceutical, energy, electronics
- about, said, that, over, in, with, had

Suppose we plot each instance of *concern* on a graph

- x-axis is the density of words in group 1
- y-axis is the density of words in group 2

Two blobs might emerge. These blobs would correspond to two different sense of *concern*.

- But in reality, we don't know the word groupings in advance.

¹example from Pedersen and Bruce (1997)

- We have to try to apply the same idea in a very high dimensional space, where every word gets its own dimension (and most dimensions are irrelevant!)
- Or we have to automatically find a low-dimensional projection. More on that much later in the course.

Here's a related scenario:

- You look at thousands of news articles from today
- Plot them on a graph of *Miley* vs *Syria*
- Three clumps emerge (Miley, Syria, others)
- Those clumps correspond to natural document classes
- Again, in reality this is a hugely high-dimensional graph

So these examples show that we can find structure in data, even without labels.

7.1 K-means clustering

You might know about classic clustering algorithms like K-means. These algorithms are iterative:

1. Guess the location of cluster centers.
2. Assign each point to the nearest center.
3. Re-estimate the centers as the mean of the assigned points.
4. Goto 2.

This is an algorithm for finding coherent “blobs” of documents. There is a variant called “soft k-means.”

- Instead of assigning each point x_i to a specific cluster z_i
- You assign it a distribution over clusters $q_i(z_i)$

We're now going to explore a more principled, statistical version of soft K-means, called EM clustering.

By understanding the statistical principles underlying the algorithm, we can extend it in a number of cool ways.

7.2 The Expectation-Maximization Algorithm

Let's go back to the Naive Bayes model:

$$\log p(\mathbf{x}, \mathbf{y}; \phi, \mu) = \sum_i \log p(\mathbf{x}_i | y_i; \phi) P(y_i; \mu)$$

For example, \mathbf{x} can describe the documents that we see today, and \mathbf{y} can correspond to their labels. But suppose we never observe y_i ? Can we still do something?

Since we don't know \mathbf{y} , let's marginalize it:

$$\log p(\mathbf{x}) = \log \sum_{\mathbf{y}} p(\mathbf{x} | \mathbf{y}; \phi) p(\mathbf{y}; \mu) \quad (7.1)$$

$$= \log \sum_{\mathbf{y}} \prod_i p(\mathbf{x}_i | y_i; \phi) p(y_i; \mu) \quad (7.2)$$

$$= \sum_i \log \sum_{y_i} p(\mathbf{x}_i | y_i; \phi) p(y_i; \mu) \quad (7.3)$$

Now we introduce an auxiliary variable q_i , for each y_i . We have the usual constraints: $\sum_y q_i(y) = 1$ and $\forall y, q_i(y) \geq 0$. In other words, q_i defines a probability distribution over Y , for each instance i .

Now since $\frac{q_i(y)}{q_i(y)} = 1$,

$$\begin{aligned} \log p(\mathbf{x}) &= \sum_i \log \sum_{y_i} p(\mathbf{x}_i | y_i; \phi) p(y_i; \mu) \frac{q_i(y)}{q_i(y)} \\ &= \sum_i \log E_q \left[\frac{p(\mathbf{x}_i | y; \phi) p(y; \mu)}{q_i(y)} \right], \end{aligned}$$

by the definition of expectation. (Note that E_q just means the expectation under the distribution q .)

Now we apply *Jensen's inequality*. Jensen's equality says that because \log is concave, we can push it inside the expectation, and obtain a lower bound.

$$\begin{aligned} \log p(\mathbf{x}) &\geq \sum_i E_q \left[\log \frac{p(\mathbf{x}_i | y; \phi) p(y; \mu)}{q_i(y)} \right] \\ \mathcal{J} &= \sum_i E_q [\log p(\mathbf{x}_i | y; \phi)] + E_q [\log p(y; \mu)] - E_q [q_i(y)] \end{aligned}$$

By maximizing \mathcal{J} , we are maximizing a lower bound on the joint log-likelihood $\log p(\mathbf{x})$.

Now, \mathcal{J} is a function of two arguments:

- the distributions $q_i(\mathbf{y})$ for each i
- the parameters μ and ϕ

We'll optimize with respect to each of these in turn, holding the other one fixed.

The E-step

First, we expand the expectation in the lower bound as:

$$\begin{aligned}\mathcal{J} &= \sum_i E_q[\log p(\mathbf{x}_i|y; \phi)] + E_q[\log p(y; \mu)] - E_q[q_i(y)] \\ &= \sum_i \sum_y q_i(y) (\log p(\mathbf{x}_i|Y_i = y; \phi) + \log p(y; \mu) - \log q_i(y))\end{aligned}$$

As in relative frequency estimation of Naive Bayes, we need to add a Lagrange multiplier to ensure $\sum_y q_i(y) = 1$, so

$$\begin{aligned}\mathcal{J} &= \sum_i \sum_y q_i(y) (\log p(\mathbf{x}_i|Y_i = y; \phi) + \log p(y; \mu) - \log q_i(y)) + \lambda_i(1 - \sum_y q_i(y)) \\ \frac{\partial \mathcal{J}}{\partial q_i(y)} &= \log p(\mathbf{x}_i|Y_i = y; \phi) + \log p(y; \mu) - \log q_i(y) - 1 - \lambda_i \\ \log q_i(y) &= \log p(\mathbf{x}_i|Y_i = y; \phi) + \log p(y; \mu) - 1 - \lambda_i \\ q_i(y) &\propto p(\mathbf{x}_i|Y_i = y; \phi)p(y; \mu) \\ &\propto p(\mathbf{x}_i, y; \phi, \mu) \\ q_i(y) &= \frac{p(\mathbf{x}_i, y; \phi, \mu)}{\sum_{y'} p(\mathbf{x}_i, y'; \phi, \mu)} \\ &= P(Y_i = y|\mathbf{x}_i; \theta, \phi)\end{aligned}$$

After normalizing, each $q_i(y)$ – which is the soft distribution over clusters for data \mathbf{x}_i – is set to the conditional probability $P(y_i|\mathbf{x}_i)$ under the current parameters μ, ϕ .

This is called the E-step, or “expectation step,” because it is derived from updating the expected likelihood under $q(\mathbf{y})$.

The M-step

Next, we hold $q(\mathbf{y})$ fixed and update the parameters. Let's do ϕ , which parametrizes $p(\mathbf{x}|\mathbf{y})$. Again, we start by adding Lagrange multipliers to the lower bound,

$$\begin{aligned}\mathcal{J} &= \sum_i \sum_y q_i(y) (\log p(\mathbf{x}_i | Y_i = y; \phi) + \log p(y; \mu) - \log q_i(y)) + \sum_y \lambda_y (1 - \sum_j \phi_{y,j}) \\ \frac{\partial \mathcal{J}}{\partial \phi_{y,j}} &= \sum_i q_i(y) \frac{x_{i,j}}{\phi_{y,j}} - \lambda_y \\ \lambda_y \phi_{y,j} &= \sum_i q_i(y) x_{i,j} \\ \phi_{y,j} &= \frac{\sum_i q_i(y) x_{i,j}}{\sum_{j'} \sum_i q_i(y) x_{i,j'}} = \frac{E_q[\text{count}(y, j)]}{E_q[\text{count}(y)]}\end{aligned}$$

So ϕ_y is now equal to the relative frequency estimate of the **expected counts** under the distribution $q(y)$.

- As in supervised Naïve Bayes, we can apply smoothing to add α to all these counts
- The update for μ is identical: $\mu_y \propto \sum_i q_i(y)$, the expected proportion of cluster $Y = y$. If needed, we can add smoothing here too.
- So, everything in the M-step is just like Naive Bayes, except we used expected counts rather than observed counts.

Coordinate ascent

Algorithms that alternate between updating various subsets of the parameters are called “coordinate-ascent” algorithms.

The objective function \mathcal{J} is **biconvex**, meaning that it is separately convex in $q(\mathbf{y})$ and $\langle \mu, \phi \rangle$, but it is not jointly convex.

- Each step is guaranteed not to decrease \mathcal{J}
- This is called hill-climbing: you never go down.
- Specifically, EM is guaranteed to converge to a **local optima** – a point which is as good or better than any of its immediate neighbors. But there may be many such points.

- But the overall procedure is **not** guaranteed to find a global maximum.
- This means that initialization is important: where you start can determine where you finish.
- This is not true in most of the supervised learning algorithms that we have considered, such as logistic regression; in that case, we are optimizing $\log p(\mathbf{y}|\mathbf{x}; \boldsymbol{\theta})$, which is defined so as to be convex with respect to the parameter $\boldsymbol{\theta}$. This means that for logistic regression (and many other supervised learning algorithms), we don't need to worry about initialization, because it won't affect our ultimate solution: we are guaranteed to reach the global minimum.

7.3 Applications of EM

EM is not really an “algorithm” like, say, quicksort. Rather, it's a framework for learning with missing data. The recipe for using EM on a problem of interest to you is something like this:

- Introduce latent variables \mathbf{z} , such that it's easy to write the probability $P(\mathcal{D}, \mathbf{z})$, where \mathcal{D} is your observed data, and easy to estimate the associated parameters.
- Derive the E-step updates for $q(\mathbf{z})$, which is typically factored as $q(\mathbf{z}) = \prod_i q_{z_i}(z_i)$.

Some applications of this basic setup are presented here.

Word sense clustering

In the “demos” folder, you can find a demonstration of expectation-maximization for word sense clustering. I assume we know that there are two senses, and that the senses can be distinguished by the contextual information in the document. The basic framework is identical to the clustering model of EM as presented above.

Semi-supervised learning

Nigam et al. (2000) offer another application of EM: **semi-supervised learning**. They apply this idea to document classification in the classic “20 Newsgroup” dataset.

- In this setting, we have labels for some of the instances, $\langle \mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)} \rangle$, but not for others, $\langle \mathbf{x}^{(u)} \rangle$.
- Can unlabeled data improve learning?

We will choose parameters to maximize the joint likelihood,

$$\log p(\mathbf{x}^{(\ell)}, \mathbf{x}^{(u)}, \mathbf{y}^{(\ell)}) = \log p(\mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)}) + \log p(\mathbf{x}^{(u)}) \quad (7.4)$$

- We treat the labels of the unlabeled documents as missing data. In the E-step we impute $q(y)$ for the unlabeled documents only.
- The M-step computes estimates of μ and ϕ from the sum of the observed counts from $\langle \mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)} \rangle$ and the expected counts from $\langle \mathbf{x}^{(u)} \rangle$ and $q(\mathbf{y})$.
- We can further parametrize this approach by weighting the unlabeled documents by a scalar λ , which is a tuning parameter.

Multi-component modeling

- One of the classes in 20 newsgroups is `comp.sys.mac.hardware`.
- Suppose that there are two kinds of posts: reviews of new hardware, and question-answer posts about hardware problems.
- The language in these **components** of the `mac.hardware` class might have little in common.
- So we might do better if we model these components separately.

We can envision a new generative process here:

- For each document i ,
 - draw the label $y_i \sim \text{Categorical}(\theta)$
 - draw the component $z_i | y_i \sim \text{Categorical}(\psi_{y_i})$
 - draw the vector of counts $\mathbf{x}_i | z_i \sim \text{Multinomial}(\phi_{z_i})$

Our labeled data includes $\langle \mathbf{x}_i, y_i \rangle$, but not z_i , so this is another case of missing data.

$$\begin{aligned} p(\mathbf{x}_i, y_i) &= \sum_z p(\mathbf{x}_i, y_i, z) \\ &= p(\mathbf{x}_i | z; \phi) p(z | y_i; \psi) p(y_i; \mu) \end{aligned}$$

Again, we can apply EM

- We need a distribution over the missing data, $q_i(z)$. This is updated during the E-step.
- During the m-step, we compute:

$$\begin{aligned} \psi_{y,z} &= \frac{E_q[\text{count}(y, z)]}{\sum_{z'} E_q[\text{count}(y, z')]} \\ \phi_{j,y,z} &= \frac{E_q[\text{count}(z, j)]}{\sum_{j'} E_q[\text{count}(z, j')]} \end{aligned}$$

- Suppose we assume each class y is associated with K components, \mathcal{Z}_y . We can add a constraint to the E-step so that $q_i(z) = 0$ if $z \notin \mathcal{Z}_y \wedge Y_i = y$.

Chapter 8

Language models

A **language model** is used to compute the probability of a sequence of text. Why would we want to do this? Thus far, we have considered problems where text is the **input**, and we want to select an output, such as a document class or a word sense. But in many of the most prominent problems in language technology, text itself is the output:

- machine translation
- speech recognition
- summarization

As we will soon see, we can produce more **fluent** text output by computing the probability of the text.

Specifically, suppose we have a vocabulary of word types

$$\mathcal{V} = \{aardvark, abacus, \dots, zither\} \quad (8.1)$$

Given a sequence of word tokens w_1, w_2, \dots, w_M , with $w_i \in \mathcal{V}$, we would like to compute the probability $p(w_1, w_2, \dots, w_M)$. We will do this in a data-driven way, assuming we have a **corpus** of text.

- For now, we'll assume that the vocabulary \mathcal{V} covers all the word tokens that we will ever see. Of course, we can enforce this by allocating a special token ♠ for unknown words. However, this might not be a great solution, as we will see later.
- Language models typically make an independence assumption across sentences, $p(s_1, s_2, \dots) = \prod_j p(s_j)$, where each sentence $s_j = [w_1, w_2, \dots, w_{N_j}]$.

So for our purposes, it is sufficient to compute the probability of sentences. The justification for this assumption is that the probability of words that are not in the same sentence don't depend on each other too much. Clearly this isn't true: once I mention *Manuel Noriega* once in a document, I'm far more likely to mention him again (Church, 2000). But the dependencies between words within a sentence are usually even stronger, and are more relevant to the fluency considerations inherent in applications such as translation and speech recognition (which are typically evaluated at the sentence level anyway).

So how can we compute the probability of a sentence? The simplest idea would be to apply a **relative frequency estimator**:

$$p(\text{Computers are useless, they can only give you answers}) \quad (8.2)$$

$$= \frac{\text{count}(\text{Computers are useless, they can only give you answers})}{\text{count}(\text{all sentences ever spoken})} \quad (8.3)$$

It's useful to think about this estimator in terms of bias and variance.

- In the theoretical limit of infinite data, it might work. But in practice, we are asking for accurate counts over an infinite number of events, since sentences can be arbitrarily long.
- Even if we set an aggressive upper bound of, say, $n = 20$, the number of possible sentences is $\#|\mathcal{V}|^{20}$. A small vocabulary for English would have $\#|\mathcal{V}| = 10^4$, so we would have 10^{80} possible sentences.
- Clearly, this estimator is extremely data-hungry. We need to introduce bias to have a chance of making reliable estimates.

Are language models meaningful? What are the probabilities of the following two sentences?

- *Colorless green ideas sleep furiously*
- *Furiously sleep ideas green colorless*

Noam Chomsky used this pair of examples to argue that the probability of a sentence is a meaningless concept:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Any English speaker can tell that the first sentence is grammatical but the second sentence is not.
- Yet neither sentence, nor their substrings, had ever appeared at the time that Chomsky wrote this article (they have appeared lots since then).
- Thus, he argued, empirical probabilities can't distinguish grammatical from ungrammatical sentences.

Pereira (2000) showed that by identifying *classes* of words (e.g., noun, verb, adjective, adverb — but not necessarily these grammatical categories), it is easy to show that the first sentence is more probable than the second. We will talk about class-based language models later.

Are language models useful? Suppose we want to translate a sentence from Spanish:

- *El cafe negro me gusta mucho.*
- Word-for-word: *The coffee black me pleases much.*
- But a good language model of English will tell us:

$$P(\text{The coffee black me pleases much}) < P(\text{I like black coffee a lot}) \quad (8.4)$$

- How can we use this fact?

Warren Weaver on translation as decoding:

When I look at an article in Russian, I say: 'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.'

This motivates a generative model (like Naive Bayes!):

- English sentence $\mathbf{w}^{(e)}$ generated from language model $p_e(\mathbf{w}^{(e)})$
- Spanish sentence $\mathbf{w}^{(s)}$ generated from noisy channel $p_{s|e}(\mathbf{w}^{(s)}|\mathbf{w}^{(e)})$

(picture)

Then the **decoding** problem is: $\max_{\mathbf{w}^{(e)}} p(\mathbf{w}^{(e)}|\mathbf{w}^{(s)}) \propto p(\mathbf{w}^{(s)}, \mathbf{w}^{(e)}) = p(\mathbf{w}^{(e)})p(\mathbf{w}^{(s)}|\mathbf{w}^{(e)})$

- The **translation model** is $p(w^{(s)}|w^{(e)})$. This ensures the **adequacy** of the translation.
- The **language model** is $p(w^{(e)})$. This ensures the **fluency** of the translation.

What else can we model with a noisy channel?

- Speech recognition (original = words; encoded = sound)
- Spelling correction (original = well-spelled text; encoded = text with spelling mistakes)
- Part of speech tagging (original = tags; encoded = words)
- Parsing (original = parse tree; encoded = words)
- ...

The noisy channel model allows us to decompose NLP systems into two parts:

- The translation model, which we need labeled data to estimate.
- The language model, which we need only *unlabeled* data to estimate.

Since there is always more unlabeled data, this means we can improve NLP systems just by improving $p_e(w)$.

8.1 N-gram language models

We began with the relative frequency estimator,

$$p(\text{Computers are useless, they can only give you answers}) \quad (8.5)$$

$$= \frac{\text{count}(\text{Computers are useless, they can only give you answers})}{\text{count}(\text{all sentences ever spoken})} \quad (8.6)$$

We'll define the probability of a sentence as the probability of the words (in order): $p(w) = p(w_1, w_2, \dots, w_M)$. We can apply the chain rule:

$$\begin{aligned} p(w) &= p(w_1, w_2, \dots, w_M) \\ &= p(w_1)p(w_2 | w_1)p(w_3 | w_2, w_1) \dots p(w_M | w_{M-1}, \dots, w_1) \end{aligned}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

Each element in the product is the probability of a word given all its predecessors. We can think of this as a *word prediction* task: *Computers are [BLANK]*. The relative frequency estimate:

$$p(\text{useless} | \text{computers are}) = \frac{\text{count}(\text{computers are useless})}{\sum_x \text{count}(\text{computers are } x)} = \frac{\text{count}(\text{computers are useless})}{\text{count}(\text{computers are})}$$

Note that we haven't made any approximations yet, and we could have applied the chain rule in reverse order, $p(\mathbf{w}) = p(w_M)p(w_{M-1}|w_M)\dots$, or in any other order. But this means that we also haven't really improved anything either: to compute the conditional probability $P(W_M | W_{M-1}, W_{M-2}, \dots)$, we need to model $\#|\mathcal{V}|^{N-1}$, with $\#|\mathcal{V}|$ events. We can't even **store** this probability distribution, let alone reliably estimate it.

N-gram models

N-gram models make a simple approximation: condition on only the past $n - 1$ words.

$$p(w_m | w_{m-1} \dots w_1) \approx P(w_m | w_{m-1}, \dots, w_{m-n+1})$$

This means that the probability of a sentence \mathbf{w} can be computed as

$$p(w_1, \dots, w_M) \approx \prod_m p(w_m | w_{m-1}, \dots, w_{m-n+1})$$

- To compute the probability of a whole sentence, it's convenient to pad the beginning and end with special symbols \diamond and \square . Then the bigram ($n = 2$) approximation to the probability of *I like black coffee* is:

$$p(I | \diamond)p(\text{like} | I)p(\text{black} | \text{like})p(\text{coffee} | \text{black})p(\square | \text{coffee}) \quad (8.7)$$

- In this model, we have to estimate and store the probability of only $\#|\mathcal{V}|^n$ events. A very common choice is a trigram model, in which $n = 3$.
- The n-gram probabilities can be determined by relative frequency estimation,

$$p(a|b, c) = \frac{\text{count}(a, b, c)}{\text{count}(b, c)} = \frac{\text{count}(a, b, c)}{\sum_{a'} \text{count}(a', b, c)} \quad (8.8)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

There could be too problems with an n -gram language model:

- **n is too small.** In this case, we are missing important linguistic context. Consider the following sentences:
 - Gorillas *always like to groom* **THEIR** friends.
 - The computer *that's on the 3rd floor of our office building* **CRASHED**.

The bolded words depend crucially on their predecessors in italics: *their* depends on knowing that *gorillas* is plural, and *crashed* depends on knowing that the subject is a *computer*. The resulting model would offer probabilities that are too low for these sentences, and too high for sentences that fail basic linguistic tests like number agreement.

- **n is too big.** In this case, we can't make good estimates of the n -gram parameters from our dataset. See the slides for some examples of this.
- These two problems point to another **bias/variance** tradeoff. Can you see how it works?
- In reality, we often have **both** problems! Language is full of long-range dependencies, and datasets are small.

We will seek approaches to keep n large, while still making low-variance estimates of the underlying parameters. To do this, we will introduce a different sort of bias: **smoothing**. But before we talk about that, let's consider how we can evaluate language models.

Evaluating language models

- Because language models are typically components of larger systems (language modeling is not really an application itself), we would prefer **extrinsic evaluation**: does the LM help the task (translation or whatever). But this is often hard to do, and depends on details of the overall system which may be irrelevant to language modeling.
- **Intrinsic evaluation** is task-neutral. Better performance on intrinsic metrics may be expected to improve extrinsic metrics across a variety of tasks (unless we are over-optimizing the intrinsic metric).

Held-out likelihood

A popular intrinsic metric is the **held-out likelihood**.

- We obtain a test corpus, and compute the (log) probability according to our model. It is crucial that the words in this corpus were not used in estimating the model itself.
- A good model should assign high probability to this held-out data.
- Specifically, we compute

$$\ell(\mathbf{w}) = \sum_i \sum_m \log p(w_m^{(i)} | w_{m-1}^{(i)}, \dots, w_{m-n+1}^{(i)}), \quad (8.9)$$

for all sentences $\mathbf{w}^{(i)}$ in the held-out corpus.

Perplexity

Perplexity is a transformation of the held-out likelihood, into an information-theoretic quantity. Specifically, we compute

$$PP(\mathbf{w}) = 2^{-\frac{\ell(\mathbf{w})}{M}}, \quad (8.10)$$

where M is the total number of tokens in the held-out corpus.

- After this transformation, we now prefer lower values. In the limit, we obtain probability 1 for our held-out corpus, with $PP = 2^{-\log 1} = 1$.
- Assume a uniform, unigram model in which $P(s_i) = \frac{1}{V}$ for all V words in the vocabulary. Then,

$$\begin{aligned} PP(\mathbf{w}) &= \left[\left(\frac{1}{V} \right)^M \right]^{-\frac{1}{M}} \\ &= \left(\frac{1}{V} \right)^{-1} = V \end{aligned}$$

- We can think of perplexity as the *weighted branching factor* at each word in the sentence.
 - If we have solved the word prediction problem perfectly, $PP(\mathbf{w}) = 1$, because there is only one possible choice.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- If we have only a uniform model that assigns equal probability to every word, $PP(\mathbf{w}) = V$.
- Most models fall somewhere in between.
- Here's how you remember: lower perplexity is better, because you are less perplexed.

Example On 38M tokens of WSJ, $V \approx 20K$, (Jurafsky and Martin, 2009, page 97) obtain these perplexities on a 1.5M token test set.

- Unigram: 962
- Bigram: 170
- Trigram: 109

Will it keep going down? See slides from (Manning and Schütze, 1999).

Information theory*

Perplexity is very closely related to the concept of entropy, the expected value of the information contained in each word.

$$H(P) = - \sum_w p(\mathbf{w}) \log p(\mathbf{w}) \quad (8.11)$$

The true entropy of English (or any real language) is unknown. Claude Shannon, one of the founders of information theory, wanted to compute upper and lower bounds. He would read passages of 15 characters to his wife, and ask her to guess the next character, recording the number of guesses it took for her to get the correct answer. As a fluent speaker of English, his wife could provide a reasonably tight bound on the number of guesses needed per character. **Question: is this an upper bound or a lower bound?**

Cross-entropy is a relationship between two probability distributions, the true one $P(W)$ and an estimate $Q(W)$.

(c) Jacob Eisenstein 2014-2015. Work in progress.

$$\begin{aligned}
H(P, Q) &= E_P[\log Q] \\
&= - \sum_{\mathbf{w}} p(\mathbf{w}) \log q(\mathbf{w}) \\
&= \sum_{\mathbf{w}} p(\mathbf{w}) \log \frac{p(\mathbf{w})}{q(\mathbf{w})} - p(\mathbf{w}) \log p(\mathbf{w}) \\
&= D_{KL}(P||Q) + H(Q)
\end{aligned}$$

So the cross-entropy is the KL-divergence between P and Q – a non-symmetric distance measure between distributions, which we will see again later in the course – plus the entropy of P . Since P is the language itself, we can only control Q , and minimizing the cross-entropy is equivalent to minimizing the KL-divergence.

We do not have access to the true $P(W)$, just a sequence $\mathbf{w} = \{w_1, w_2, \dots\}$, which is sampled from $P(W)$. In the limit, the length of \mathbf{w} is infinite, so we have,

$$\begin{aligned}
H(P, Q) &= - \sum_{\mathbf{w}} p(\mathbf{w}) \log q(vw) \\
&= - \lim_{M \rightarrow \infty} \frac{1}{M} \log q(\mathbf{w}) \\
&\approx - \frac{1}{M} \log q(\mathbf{w}) \\
PP(S) &= 2^{-\frac{1}{M} \log q(vw)}
\end{aligned}$$

A good language model has low cross-entropy with $P(W)$, and thus low perplexity.

Further aside : A related topic in psycholinguistics is the “constant entropy rate hypothesis,” also called the “uniform information density hypothesis.” The hypothesis is that speakers should prefer linguistic choices that convey a uniform amount of information over time (Jaeger, 2010). Some evidence:

- Speakers shorten predictable words, lengthen unpredictable ones
- High-entropy sentences take longer to read
- Syntactic reductions (e.g., *I’m* versus *I am*) are more likely when the reducible word contains less information.

(c) Jacob Eisenstein 2014-2015. Work in progress.

8.2 Smoothing and discounting

We want to estimate $P(W)$ from sparse statistics, avoiding $P(W) = 0...$

Bibliography

- Berger, A. L., Pietra, V. J. D., and Pietra, S. A. D. (1996). A maximum entropy approach to natural language processing. *Computational linguistics*, 22(1):39–71.
- Bottou, L. (1998). Online learning and stochastic approximations. *On-line learning in neural networks*, 17:9.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, New York, NY, USA.
- Burges, C. J. (1998). A tutorial on support vector machines for pattern recognition. *Data mining and knowledge discovery*, 2(2):121–167.
- Church, K. W. (2000). Empirical estimates of adaptation: the chance of two Norie-gas is closer to $p/2$ than p^2 . In *Proceedings of the 18th conference on Computational linguistics-Volume 1*, pages 180–186.
- Collins, M. (2002). Discriminative training methods for hidden markov models: theory and experiments with perceptron algorithms. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1–8.
- Collins, M. (2013). Notes on natural language processing. <http://www.cs.columbia.edu/~mcollins/notes-spring2013.html>.
- Collins, M. and Duffy, N. (2001). Convolution kernels for natural language. In *Advances in neural information processing systems*, pages 625–632.
- Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167. ACM.
- Crammer, K., Dekel, O., Keshet, J., Shalev-Shwartz, S., and Singer, Y. (2006). On-line passive-aggressive algorithms. *The Journal of Machine Learning Research*, 7:551–585.

- Crammer, K. and Singer, Y. (2003). Ultraconservative online algorithms for multi-class problems. *The Journal of Machine Learning Research*, 3:951–991.
- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159.
- Dyer, C. (2014). Notes on adagrad. www.ark.cs.cmu.edu/cdyer/adagrad.pdf.
- Figueiredo, M., Graça, J., Martins, A., Almeida, M., and Coelho, L. P. (2013). LXMLS lab guide. <http://lxmls.it.pt/2013/guide.pdf>.
- Freund, Y., Schapire, R., and Abe, N. (1999). A short introduction to boosting. *Journal-Japanese Society For Artificial Intelligence*, 14(771-780):1612.
- Freund, Y. and Schapire, R. E. (1999). Large margin classification using the perceptron algorithm. *Machine learning*, 37(3):277–296.
- Hastie, T., Tibshirani, R., Friedman, J., Hastie, T., Friedman, J., and Tibshirani, R. (2009). *The elements of statistical learning*, volume 2. Springer.
- Jaeger, T. F. (2010). Redundancy and reduction: Speakers manage syntactic information density. *Cognitive psychology*, 61(1):23–62.
- Jurafsky, D. and Martin, J. H. (2009). *Speech and Language Processing (2nd Edition) (Prentice Hall Series in Artificial Intelligence)*. Prentice Hall, 2 edition.
- Manning, C. D. and Schütze, H. (1999). *Foundations of statistical natural language processing*. MIT press.
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. The MIT Press.
- Nemirovski, A. and Yudin, D. (1978). On Cezari’s convergence of the steepest descent method for approximating saddle points of convex-concave functions. *Soviet Math. Dokl.*, 19.
- Nigam, K., McCallum, A. K., Thrun, S., and Mitchell, T. (2000). Text classification from labeled and unlabeled documents using em. *Machine learning*, 39(2-3):103–134.
- Pereira, F. (2000). Formal grammar and information theory: together again? *Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 358(1769):1239–1253.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Schmid, H. (1994). Probabilistic part-of-speech tagging using decision trees. In *Proceedings of the international conference on new methods in language processing*, volume 12, pages 44–49. Manchester, UK.
- Shwartz, S. S., Singer, Y., and Srebro, N. (2007). Pegasos: Primal estimated sub-GrAdient SOLver for SVM. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 807–814.
- Sra, S., Nowozin, S., and Wright, S. J. (2012). *Optimization for machine learning*. MIT Press.
- Tausczik, Y. R. and Pennebaker, J. W. (2010). The psychological meaning of words: Liwc and computerized text analysis methods. *Journal of Language and Social Psychology*, 29(1):24–54.
- Zelenko, D., Aone, C., and Richardella, A. (2003). Kernel methods for relation extraction. *The Journal of Machine Learning Research*, 3:1083–1106.
- Zhang, T. (2004). Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proceedings of the twenty-first international conference on Machine learning*, page 116. ACM.