

# Finding Informative Features

36-350: Data Mining

10 September 2008

READINGS: David P. Feldman, “Introduction to Information Theory”, chapter 1 (on Blackboard or at <http://hornacek.coa.edu/dave/Tutorial/index.html>) *Principles of Data Mining*, sections 10.1, 10.2, 10.6 and 10.8

As I mentioned last time, everything we have learned how to do so far — similarity searching, nearest-neighbor and prototype classification, multidimensional scaling — relies on our having a vector of **features** or **attributes** for each object in data set. (These are also called **dimensions**, because the number of dimensions in the vector space is the number of features.) The success of our procedures depends on our choosing good features, but I’ve said very little about how to do this, aside from qualitative considerations about invariance. This week we’ll look at one way of picking out presumably-useful features, using information theory.

The basic idea, remember, is that the features are the aspects of the data which show up in our representation. However, they’re not what we *really* care about, which is rather something we don’t, or can’t, directly represent, like the **class** of the object (is it a post about cars or motorcycles? a picture of a flower or a tiger?). We use the observable features to make a guess (formally, an **inference**) about the unobservable thing, like the class. Good features are ones which let us make better guesses — ones which reduce our **uncertainty** about the unobserved class.

Good features are therefore **informative**, **discriminative** or **uncertainty-reducing**. This means that they need to *differ* across the different classes, at least statistically. I said before that the number of occurrences of the word “the” in an English document isn’t a useful feature, because it occurs about as often in all kinds of text. This means that looking at that count leaves us exactly as uncertain about which class of document we’ve seen as we were before. Similarly, the word “narthex” is going to be equally *rare* whether the topic is cars or motorcycles, so it’s also uninformative. On the other hand, the word “seatbelt” is going to be more common in posts about cars than in ones about motorcycles, so counting its occurrences *is* going to reduce our uncertainty. The important thing is that the distribution of the feature *differ* across the classes.

Let’s try and make these ideas about uncertainty, discrimination, and reduction in uncertainty precise. (This is where the information theory comes in.)

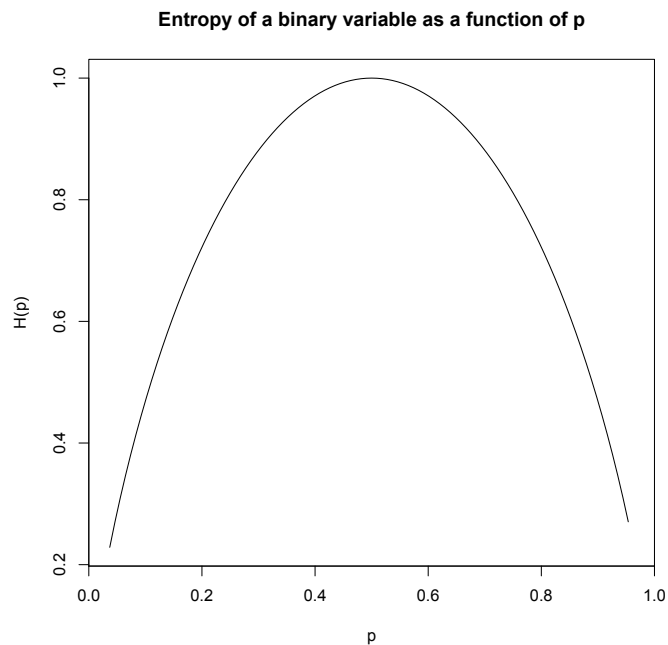


Figure 1: Entropy of a binary variable as a function of the probability of (either) class value. Note that it is symmetric around  $p = 1/2$ , where it is maximal.

Let  $X$  be some feature of the data in our representation, and  $x$  a particular value. How uncertain are we about  $X$ ? Well, one way to measure this is the **entropy** of  $X$ :

$$H[X] = - \sum_x \Pr(X = x) \log_2 \Pr(X = x)$$

The entropy, in bits, equals the average number of yes-or-no questions we'd have to ask to figure out the value of  $X$ . (This is also the number of bits of computer memory needed to store the value of  $X$ .) If there are  $n$  possible values for  $X$ , and they are all equally likely, then our uncertainty is maximal, and  $H[X] = \log_2 n$ , the maximum possible value. If  $X$  can take only one value, we have no uncertainty, and  $H[X] = 0$ .

Similarly, our uncertainty about the class  $C$ , in the absence of any other information, is just the entropy of  $C$ :

$$H[C] = - \sum_c \Pr(C = c) \log_2 \Pr(C = c)$$

Now suppose we observe the value of the feature  $X$ . This will, in general, change

our distribution for  $C$ , since we can use Bayes's Rule:

$$\Pr(C = c|X = x) = \frac{\Pr(C = c, X = x)}{\Pr(X = x)} = \frac{\Pr(X = x|C = c)\Pr(C = c)}{\Pr(X = x)}$$

$\Pr(X = x)$  tells us the frequency of the value  $x$  is over the whole population.  $\Pr(X = x|C = c)$  tells us the frequency of that value is when the class is  $c$ . If the two frequencies are not equal, we should change our estimate of the class, making it larger if that feature is more common in  $c$ , and making it smaller if that feature is rarer. Generally, our uncertainty about  $C$  is going to change, and be given by the **conditional entropy**:

$$H[C|X = x] = - \sum_c \Pr(C = c|X = x) \log_2 \Pr(C = c|X = x)$$

The difference in entropies,  $H[C] - H[C|X = x]$ , is how much our uncertainty about  $C$  has changed, conditional on seeing  $X = x$ . This change in uncertainty is **information** or **conditional information**:

$$I[C; X = x] = H[C] - H[C|X = x]$$

Notice that the conditional information can be negative. For a simple example, suppose that  $C$  is "it will rain today", and that it normally rains only one day out of seven. Then  $H[C] = 0.59$  bits. If however we look at the weather forecast, and it tells us that it will rain with 50% probability,  $H[C|X = x] = 1$  bit, so our uncertainty has increased by 0.41 bits.

We can also look at the *expected* information a feature gives us about the class:

$$I[C; X] = H[C] - H[C|X] = H[C] - \sum_x \Pr(X = x) H[C|X = x]$$

The expected information is never negative. In fact, it's not hard to show that the only way it can even be zero is if  $X$  and  $C$  are **statistically independent** — if the distribution of  $X$  is the same for all classes  $c$ ,

$$\Pr(X|C = c) = \Pr(X)$$

It's also called the **mutual information**, because it turns out that  $H[C] - H[C|X] = H[X] - H[X|C]$ . (You might want to try to prove this to yourself, using Bayes's rule and the definitions.)

For example, suppose we pick a random position in a random document. Let  $X$  be 1 if the word is "car", and 0 otherwise. The frequencies for the 10 auto/moto documents are

$c$	$X$	
	"car"	not "car"
auto	13	598
not auto	0	696

For this table,

$$\begin{aligned}H[C] &= 0.997 \\H[C|X = \text{“car”}] &= 0 \\H[C|X = \text{not “car”}] &= 0.996 \\Pr(X = \text{“car”}) &= 0.01 \\I[C; X] &= 0.997 - (0.01 * 0) - (0.99 * 0.996) = 0.01\end{aligned}$$

In words, when we see the word “car”, we can be certain that the post is about automobiles ( $H[C|X = \text{“car”}] = 0$  bits). On the other hand, when it is absent we are almost as uncertain as if we flipped a fair coin ( $H[C|X = \text{not “car”}] = 0.996$  bits), which is almost how much uncertainty as we’d have if we didn’t look at the words at all ( $H[C] = 0.997$  bits). Since most documents do not contain the word “car” ( $Pr(X = \text{“car”}) = 0.01$ ), the *expected* reduction in uncertainty is small but non-zero ( $I[C; X] = 0.01$ ).

## Finding the Important Words

Here’s one information-theoretic procedure for finding the important words.

1. Collect counts for each class,  $1 \dots K$
2. For each word, make the  $2 \times K$  table of word occurrences by classes
3. Compute the mutual information in each table. Alternatively, compute the information for the word having occurred.

This just looks at whether the presence or absence on the word is informative; we could also look at whether how often it appears is informative, but we’d need bigger tables. (And nothing hinges on this being words; we could do the same thing for colors in a bag-of-colors representation of pictures, etc.)

Calculating the expected information is actually very similar to performing a  $\chi^2$  test for independence. (Remember that mutual information is 0 if and only if the two variables are statistically independent.) In fact, if the sample size is large enough and the variables really are independent, the sample mutual information *has* a  $\chi^2$  distribution.

The figure shows the results of doing these calculations for the automobile/motorcycle posts.

- Conditional and expected information tend to rise together. (Why might this be so?)
- Expected information values are much smaller than conditional information values. (Why?)
- Some of the high-information words make sense (“bike”, “car”, “cars”), but others are weird (“are”, “your”).

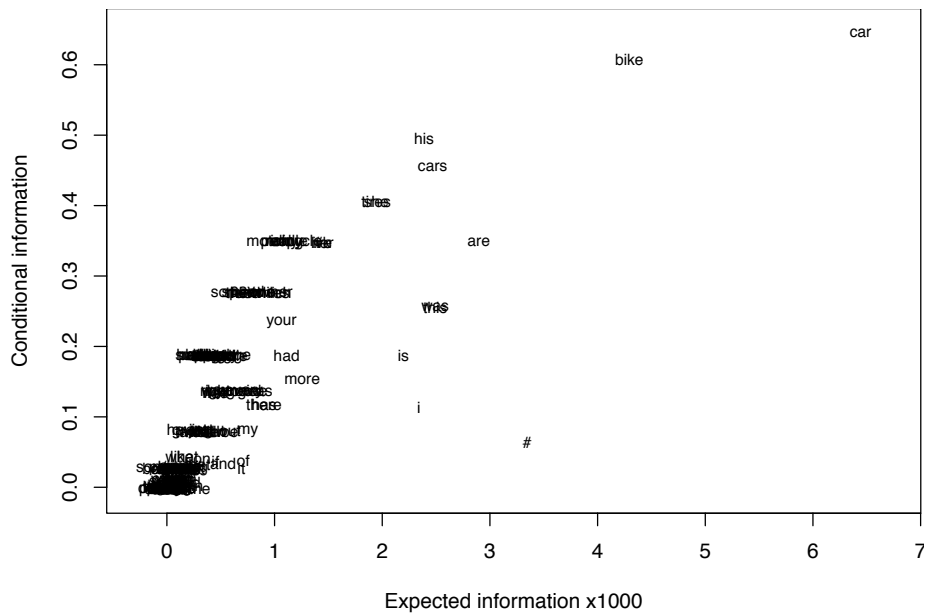


Figure 2: Reduction in uncertainty about the class (automobile/motorcycle) for the newsgroups example when a word is present (“conditional information”, vertical axis) versus the average reduction in uncertainty from checking for that word (“mutual information”, horizontal axis). Note the difference in scales.

All of this is just looking at one feature at a time, so it ignores the possibility that certain *combinations* of features are useful, or that some features are **redundant** given others. We will look at this sort of **interaction** in the next lecture.