

Assignment 1

CSE 517: Natural Language Processing – University of Washington

Winter 2020 – Due: January 20, 2020, 11:59 pm

From the Textbook

Complete these exercises from the Eisenstein book. They are ordered roughly the same way the content is ordered in lectures.

- Language modeling: 6.4, 6.5, and 6.6 (pp. 135–6)
- Linear classification: 2.5 and 2.6 (pp. 43–44)
- Implementation: 4.6 (p. 89); the last part, performing statistical significance tests, is extra credit
- Nonlinear classification: 3.4, 3.5, and 3.6 (p. 65)

Maximum *a Posteriori* Estimation

In the noisy channel paradigm, we consider two random variables, X and Y , and imagine that the observable X is a noisy or corrupted version of Y .

$$\boxed{\text{source}} \longrightarrow Y \longrightarrow \boxed{\text{channel}} \longrightarrow X$$

For example, if we observe a Chinese sentence (a value for X) but believe it was “originally” an English sentence (a value for Y), we can build source and channel models to capture, respectively, the distribution over English sentences and the translation of English into Chinese. Then, we use Bayes’ rule to “decode”:

$$\hat{y} = \operatorname{argmax}_y p(y | x) = \operatorname{argmax}_y \frac{p(y) \cdot p(x | y)}{p(x)} = \operatorname{argmax}_y p(y) \cdot p(x | y)$$

Question 1: Why can we drop the $p(x)$ factor in the denominator?

Interestingly, this same “noisy channel” idea can also be applied when we are estimating the parameters of a probabilistic model, though it goes by a different name: maximum *a posteriori* (MAP) estimation. The story looks like this:

$$\boxed{\text{source}} \longrightarrow \Theta \longrightarrow \boxed{\text{channel}} \longrightarrow X$$

Instead of Y , we use the random variable Θ for the unknown parameters of our model. X corresponds to the training data.

Applying Bayes' rule gives us:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(\theta) \cdot p(x | \theta)$$

Here, we're using the notation $p(x | \theta)$ instead of $p_{\theta}(x)$, since we're thinking of θ as a random variable (its name would be Θ) with its own probability distribution, but both are usually fine and they mean the same thing.

It's helpful to think about this for a specific model, so let's take the simplest one we know: the unigram model. The parameters for the unigram model are a vector $\theta \in \Delta^V$.¹ The evidence is a sequence of words in a text, denoted x , though—as discussed in class—we really only need to know the count of each word, not the order in which they occur.

Let's consider the “channel” model first; it defines the probability of x given the parameters. Recall this model from class:

$$p(x | \theta) = \prod_{i=1}^n \theta_{x_i} = \prod_{v \in \mathcal{V}} \theta_v^{\text{count}_x(v)} \quad (1)$$

where $\text{count}_x(v)$ is the count of word v in the data x .

We discussed maximum likelihood estimation (MLE), in which we select $\hat{\theta}$ to maximize the above quantity, and proved that relative frequency estimation accomplishes MLE. You should think of MAP as an *alternative* to MLE.

The “source” model is a bit harder to think about at first. We need a probability distribution over Δ^V . An attractive candidate for $p(\Theta)$ is the **Dirichlet** distribution.²

Here is the form of the Dirichlet distribution:

$$p_{\alpha}(\theta) = \frac{\Gamma(\sum_{v \in \mathcal{V}} \alpha_v)}{\prod_{v \in \mathcal{V}} \Gamma(\alpha_v)} \prod_{v \in \mathcal{V}} \theta_v^{(\alpha_v - 1)} \quad (2)$$

This is somewhat daunting, so we'll break it down:

- The parameters $\alpha \in \mathbb{R}^V$ are a vector of (strictly) positive values. If we normalize them (divide each by their sum), we get the mean value of the Dirichlet; if we sampled many different values of Θ , their average would be close to this normalized version of α . As these values get larger, the tendency of draws from the Dirichlet to be close to the mean increases. When they are smaller, draws from the Dirichlet will be more diffuse.
- The Γ function is an extension of the factorial function to the nonnegative reals. For natural numbers x , $\Gamma(x) = (x - 1)!$. The more general form is not important for this assignment; if you are curious, read more at <http://mathworld.wolfram.com/GammaFunction.html>.
- Fortunately for us, the Γ factors are constant with respect to θ , so they don't affect our estimation procedure.

¹The probability simplex in V dimensions, denoted Δ^V , is defined as the set of all V -length vectors that are (i) nonnegative and (ii) sum to one.

²We mentioned this family briefly when we talked about latent Dirichlet allocation, but did not go into details.

The fascinating thing about Dirichlet distributions and categorical distributions like the unigram model's θ is that the form of the posterior distribution $p(\Theta \mid \mathbf{x})$ is also a Dirichlet distribution. When this happens, Bayesian statisticians get very excited and use the term “conjugate prior” to denote the special relationship between a family of priors (here, Dirichlet distributions) and likelihood functions (here, categorical distributions). In fact, this property is why the Dirichlet was chosen as the prior in latent Dirichlet allocation in the first place.

Consider the logarithm of this quantity as a function of θ :

$$\log p(\theta \mid \mathbf{x}) = \log p_{\alpha}(\theta) + \log p(\mathbf{x} \mid \theta) - \log p(\mathbf{x}) \quad (3)$$

$$= \log p_{\alpha}(\theta) + \log p(\mathbf{x} \mid \theta) + \text{constant} \quad (4)$$

$$= \log \frac{\Gamma(\sum_{v \in \mathcal{V}} \alpha_v)}{\prod_{v \in \mathcal{V}} \Gamma(\alpha_v)} \prod_{v \in \mathcal{V}} \theta_v^{(\alpha_v - 1)} + \log \prod_{v \in \mathcal{V}} \theta_v^{\text{count}_{\mathbf{x}}(v)} + \text{constant} \quad (5)$$

$$= \log \prod_{v \in \mathcal{V}} \theta_v^{(\alpha_v - 1)} + \log \prod_{v \in \mathcal{V}} \theta_v^{\text{count}_{\mathbf{x}}(v)} + \text{constant} \quad (6)$$

where we've just plugged in the expressions from (1) and (2). Starting on line (4), we fold all terms that do not depend on θ into “constant.”³

Question 2: Conjugacy means that the above posterior can be expressed as a Dirichlet distribution with some parameters; call them α'_v . Write α'_v as a function of α and the counts $\text{count}_{\mathbf{x}}(*)$.

MAP estimation requires that we find the most probable value of θ under this posterior—in other words, its mode. The mode has a closed form whenever all $\alpha_v > 1$ for all $v \in \mathcal{V}$. The closed form is:

$$\left[\underset{\theta}{\operatorname{argmax}} p_{\alpha}(\theta) \cdot p(\mathbf{x} \mid \theta) \right]_v = \frac{\alpha'_v - 1}{\sum_{v' \in \mathcal{V}} \alpha'_{v'} - V} \quad (7)$$

for each $v \in \mathcal{V}$.

Question 3: For an appropriate choice of α , the Dirichlet can be “uninformative,” meaning that it assigns equal probability to all possible θ . In this case, MAP is the same as MLE! Derive the values of α that will make this happen.

Question 4: For a (different) appropriate choice of α , we make MAP estimation equivalent to Laplace smoothing.⁴ Derive the values of α that will make this happen.

Question 5: A *symmetric* Dirichlet distribution is one where all α_v take the same value (say, a). Show that, for any a that is sufficiently large, we can write the MAP estimate as an interpolation between the MLE (weighted by some λ) and a uniform distribution over \mathcal{V} (weighted by $1 - \lambda$). Derive λ as a function of a .

³What we mean here is only that they are constant *with respect to* θ and therefore do not affect our objective of finding θ to maximize the (log-)likelihood.

⁴In lecture, we may have mentioned Laplace smoothing casually but probably not in detail or with the name “Laplace smoothing.” It's what the instructor would have described as “the simplest smoothing approach you can think of.” Before normalizing counts into probabilities (as you do in MLE), take the count of every $v \in \mathcal{V}$ (including the v that have a count of zero) and add one to it. So now, everything you didn't see, you're pretending you saw once. Everything you *did* see, you're pretending you saw it one more time than you did. If, instead of adding one, you add a (positive) value λ , then this method is called “add- λ ” smoothing.