

# Notes on some basic probability stuff

David Meyer

dmm@{1-4-5.net,uoregon.edu,brocade.com,...}

September 30, 2014

## 1 Introduction

Note well: There are likely to be many mistakes in this document. That said...

Much of what is described here follows from two simple rules:

$$\text{Sum Rule: } P(\mathcal{X}) = \sum_y P(\mathcal{X}, \mathcal{Y}) \quad (1)$$

$$\text{Product Rule: } P(\mathcal{X}, \mathcal{Y}) = P(\mathcal{X}|\mathcal{Y})P(\mathcal{Y}) \quad (2)$$

- The Sum Rule is sometimes called marginalization
- The Product Rule is part of the proof of the Hammersley-Clifford Theorem

Remember also that  $\mathcal{X} = \{x_i\}_{i=1}^{|\mathcal{X}|}$ , where each  $x_i$  is a realization of the random variable  $x^1$ . Lets also say that the set  $\Theta$  of probability distribution parameters can be used to explain the *evidence*  $\mathcal{X}$ . Then we say that the "manner in which the evidence  $\mathcal{X}$  depends on the parameters  $\Theta$ " is the *observation model*. The analytic form of the observation model is the likelihood  $P(\mathcal{X}|\Theta)$ .

---

<sup>1</sup>Each observation  $x_i$  is, in general, a data point in a multidimensional space.

## 2 Estimating the parameters $\Theta$ with Bayes' Theorem

Note that

$$P(\Theta, \mathcal{X}) = P(\mathcal{X}, \Theta) \quad (3)$$

$$P(\Theta, \mathcal{X}) = P(\Theta|\mathcal{X})P(\mathcal{X}) \quad (4)$$

$$P(\mathcal{X}, \Theta) = P(\mathcal{X}|\Theta)P(\Theta) \quad (5)$$

$$P(\Theta|\mathcal{X})P(\mathcal{X}) = P(\mathcal{X}|\Theta)P(\Theta) \quad (6)$$

Solving for  $P(\Theta|\mathcal{X})$  we get Bayes' Theorem

$$P(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta)P(\Theta)}{P(\mathcal{X})} \quad (7)$$

$$(8)$$

Said another way

$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}} \quad (9)$$

You might also see Bayes' Theorem written using the *Law of Total Probability*<sup>2</sup> which is sometimes written as follows:

$$P(A) = \sum_n P(A \cap B_n) \quad \# \text{ by the } \textit{Sum Rule} \text{ (Equation 1)} \quad (10)$$

$$= \sum_n P(A, B_n) \quad \# \text{ in the notation used in Equation 1} \quad (11)$$

$$= \sum_n P(A|B_n)P(B_n) \quad \# \text{ by the } \textit{Product Rule} \text{ (Equation 2)} \quad (12)$$

so that the posterior distribution  $P(\mathcal{C}_1|\mathbf{x})$  for two classes  $\mathcal{C}_1$  and  $\mathcal{C}_2$  given input vector  $\mathbf{x}$  would look like

$$P(\mathcal{C}_1|\mathbf{x}) = \frac{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1)}{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1) + P(\mathbf{x}|\mathcal{C}_2)P(\mathcal{C}_2)} \quad (13)$$

---

<sup>2</sup>The Law of Total Probability is a combination of the Sum and Product Rules

Interestingly, the posterior distribution is related to logistic regression as follows: First recall that the posterior  $P(\mathcal{C}_1|\mathbf{x})$  is

$$P(\mathcal{C}_1|\mathbf{x}) = \frac{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1)}{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1) + P(\mathbf{x}|\mathcal{C}_2)P(\mathcal{C}_2)} \quad (14)$$

Now, if we set

$$a = \ln \frac{P(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1)}{P(\mathbf{x}|\mathcal{C}_2)P(\mathcal{C}_2)} \quad (15)$$

we can see that

$$P(\mathcal{C}_1|\mathbf{x}) = \frac{1}{1 + e^{-a}} = \sigma(a) \quad (16)$$

that is, the sigmoid function.

## 2.1 Maximum Likelihood Estimation (MLE)

Given all of that, for the MLE we seek the value of  $\Theta$  that maximizes the likelihood  $P(\mathcal{X}|\Theta)$  for our observations  $\mathcal{X}$ . Remembering that  $\mathcal{X} = \{x_1, x_2, \dots\}$  and that the  $x_i$  are iid, the value of  $\Theta$  we seek maximizes

$$\prod_{x_i \in \mathcal{X}} P(x_i|\Theta) \quad (17)$$

Because of the product it is easier to use the  $\log^3$ , we use the log likelihood  $\mathcal{L}$ :

$$\mathcal{L} = \sum_{x_i \in \mathcal{X}} \log P(x_i|\Theta) \quad (18)$$

and define  $\hat{\Theta}_{ML}$  as follows:

$$\hat{\Theta}_{ML} = \operatorname{argmax}_{\Theta} \mathcal{L} \quad (19)$$

---

<sup>3</sup>and since  $\log(x)$  is monotonically increasing it doesn't effect the argmax

The maximization is obtained by (calculus tricks):

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = 0 \quad \forall \theta_i \in \Theta \quad (20)$$

Note finally that in a Generalized Linear Regression setting, we have

$$\eta = \mathbf{w}^T \mathbf{x} + b \quad (21)$$

$$p(y|\mathbf{x}) = p(y|g(\eta); \theta) \quad (22)$$

where  $g(\cdot)$  is an *inverse link function*, also referred to as an activation function. For example, if the link function is the logistic function, then the inverse link function  $g(\eta) = \frac{1}{1+e^{-\eta}}$  and the negative log-likelihood  $\mathcal{L}$  is

$$\mathcal{L} = -\log p(y|g(\eta); \theta) \quad (23)$$

## 2.2 Maximum a Posteriori (MAP) Estimation of $\Theta$

Recall that

$$P(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta)P(\Theta)}{P(\mathcal{X})} \quad (24)$$

We are seeking the value of  $\Theta$  that maximizes  $P(\Theta|\mathcal{X})$ , so the solution can be stated as

$$\hat{\Theta}_{MAP} = \operatorname{argmax}_{\Theta} P(\Theta|\mathcal{X}) \quad (25)$$

$$= \operatorname{argmax}_{\Theta} \frac{P(\mathcal{X}|\Theta) \cdot P(\Theta)}{P(\mathcal{X})} \quad (26)$$

However, since  $P(\mathcal{X})$  does not depend on  $\Theta$ , we can write

$$\hat{\Theta}_{MAP} = \operatorname{argmax}_{\Theta} P(\mathcal{X}|\Theta) \cdot P(\Theta) \quad (27)$$

$$= \prod_{x_i \in \mathcal{X}} P(x_i|\Theta) \cdot P(\Theta) \quad (28)$$

If we again take the log, we get

$$\hat{\Theta}_{MAP} = \operatorname{argmax}_{\Theta} \left( \sum_{x_i \in \mathcal{X}} \log P(x_i|\Theta) + \log P(\Theta) \right) \quad (29)$$

## 2.3 Notes

- Both MLE and MAP are point estimates for  $\Theta$  (contrast probability distributions)
- MLE notoriously overfits
- MAP allows us to take into account knowledge about the prior (which is a sort of a regularizer)
- Bayesian estimation, by contrast, calculates the full posterior distribution  $P(\Theta|\mathcal{X})$

## 2.4 Bayesian Estimation

Recall that Bayesian estimation calculates the full posterior distribution  $P(\Theta|\mathcal{X})$ , where

$$P(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{P(\mathcal{X})} \quad (30)$$

In this case, however, the denominator  $P(\mathcal{X})$  cannot be ignored, and we know from the *sum* and *product* rules that

$$P(\mathcal{X}) = \int_{\Theta} P(\mathcal{X}, \Theta) d\Theta \quad (31)$$

$$= \int_{\Theta} P(\mathcal{X}|\Theta) P(\Theta) d\Theta \quad (32)$$

putting it all together we get

$$P(\Theta|\mathcal{X}) = \frac{P(\mathcal{X}|\Theta) P(\Theta)}{\int_{\Theta} P(\mathcal{X}|\Theta) P(\Theta) d\Theta} \quad (33)$$

If we want to be able to derive an algebraic form for the posterior  $P(\Theta|\mathcal{X})$ , the most challenging part will be finding the integral in the denominator. This is where the idea of *conjugate priors* and approximate inference approaches (*Monte Carlo Integration* and *Variational Bayesian methods*<sup>4</sup>) are useful. Need to further expand this...

---

<sup>4</sup>Variational Bayesian methods are a family of techniques for approximating intractable integrals arising in Bayesian inference and machine learning. They are typically used in complex statistical models consisting of observed variables (usually termed "data") as well as unknown parameters and latent variables, with various sorts of relationships among the three types of random variables, as might be described by a graphical model.

### 3 Monte Carlo Integration

Suppose we have a distribution  $p(\theta)$  (perhaps a posterior) the we want to sample quantities of interest from. To do this analytically, we need to take an integral of the form

$$I = \int_{\Theta} g(\theta) p(\theta) d\theta \tag{34}$$

where  $g(\theta)$  is some function of  $\theta$  (typically  $g(\theta) = \theta$  (the mean), etc). Need a deeper analysis here (note to self), but the punchline is that you can estimate  $I$  using *Monte Carlo Integration* as follows: Sample  $M$  values  $(\theta^i)$  from  $p(\theta)$  and calculate

$$\hat{I}_M = \frac{1}{M} \sum_{i=1}^M g(\theta^i) \tag{35}$$

Note that this works fine if the samples from  $p(\theta)$  are iid<sup>5</sup> but if not, we can use a Markov Chain to draw "slightly dependent" samples and depend on the *Ergodic Theorem* (see Section 5.2).

### 4 Acknowledgements

---

<sup>5</sup>We know this by the Strong Law of Large Numbers, see Section 5.1.

## References

### 5 Appendix

#### 5.1 Strong Law of Large Numbers

Let  $X_1, X_2, \dots, X_M$  be a sequence of **independent** and **identically distributed** random variables, each having a finite mean  $\mu_i = E[X_i]$ .

Then with probability 1

$$\frac{1}{M} \sum_{i=1}^M X_i \rightarrow E[X] \quad (36)$$

as  $M \rightarrow \infty$ .

#### 5.2 Ergodic Theorem

Let  $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(M)}$  be  $M$  samples from a Markov chain that is *aperiodic*, *irreducible*, and *positive recurrent*<sup>6</sup>, and  $E[g(\theta)] < \infty$ .

Then with probability 1

$$\frac{1}{M} \sum_{i=1}^M g(\theta_i) \rightarrow E[g(\theta)] = \int_{\Theta} g(\theta) \pi(\theta) d\theta \quad (37)$$

as  $M \rightarrow \infty$  and where  $\pi$  is the stationary distribution of the Markov chain.

---

<sup>6</sup>In this case, the chain is said to be *ergodic*.