

Instructors: Piyush Rai  
Authors: Gurpreet Singh  
Date: December 2, 2017

## Hyper-Parameter Estimation using MLE-II

In the previous scribes, we have discussed how to estimate parameters for probabilistic models, and looked at parameter estimation for the coin-toss example as well as the regression problem. In this scribe, we will look at estimating the hyperparameters for Bayesian Regression, using the MLE-II approach.

### 1. Learning Hyperparameters

Let us reconsider the Bayesian Linear Regression model.

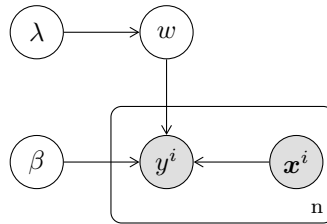


Figure 1: Plate Notation of a Bayesian Linear Regression Model

As discussed earlier, we estimate the hyperparameters using methods like Cross-Validation, however we wish to probabilistically estimate a good value for the hyperparameters. There are various methods to do this, such as

- (i) Doing point estimation of hyperparameters and learning full posterior for the parameters (MLE-II)
- (ii) Using general inference methods such as MCMC and Variational Bayes (Discussed in later scribes)

Before we move on to learning these methods, we must first understand why we need them. Let us try a simple probabilistic approach to estimating hyperparameters (determining the full posterior).

Let us first take a look at the form of the prior for the hyperparameters and the parameters (all can be considered as the parameters of the model)  $\mathbb{P}[\mathbf{w}, \beta, \lambda]$ .

From the plate notation, we can say that  $\beta$  and  $\lambda$  are independent of other parameters, whereas  $\mathbf{w}$  is dependent on  $\lambda$  (which makes intuitive sense as well since  $\lambda$  is the precision value for the assumed gaussian that explains the distribution of  $\mathbf{w}$ ). Therefore,

$$\mathbb{P}[\mathbf{w}, \beta, \lambda] = \mathbb{P}[\mathbf{w} | \lambda] \cdot \mathbb{P}[\beta] \cdot \mathbb{P}[\lambda]$$

Now, let us try to compute the posterior —

$$\mathbb{P}[\mathbf{w}, \beta, \lambda | \mathcal{X}, \mathcal{Y}] = \frac{\mathbb{P}[\mathcal{Y} | \mathcal{X}, \mathbf{w}, \beta, \lambda] \mathbb{P}[\mathbf{w}, \beta, \lambda]}{\mathbb{P}[\mathcal{Y} | \mathcal{X}]}$$

where

$$\mathbb{P}[\mathcal{Y} | \mathcal{X}] = \int_{\mathbf{w}, \beta, \lambda} \mathbb{P}[\mathbb{P}[\mathcal{Y} | \mathcal{X}, \mathbf{w}, \beta]] \mathbb{P}[\mathbf{w}, \beta, \lambda] \quad (1)$$

It can be shown that the above term is intractable. Therefore, we cannot perform Bayesian Inference on this directly. This is the reason we require different approaches to predicting the hyperparameters.

**Exercise 4.1.** Prove that the integral in Equation 1 is indeed intractable.

**Note.** The approaches to estimate hyperparameters are not just limited to hyperparameters, but can also be extended for multi-parameter models

## 2. Learning Hyperparameters

As mentioned in the previous section, we have the prior over  $\mathbf{w}, \beta, \lambda$  with the form

$$\mathbb{P}[\mathbf{w}, \beta, \lambda] = \mathbb{P}[\mathbf{w} \mid \lambda] \cdot \mathbb{P}[\beta] \cdot \mathbb{P}[\lambda] \quad (2)$$

We can also write the posterior (using Bayes Rule) as

$$\mathbb{P}[\mathbf{w}, \beta, \lambda \mid \mathcal{X}, \mathcal{Y}] = \mathbb{P}[\mathbf{w} \mid \mathcal{X}, \mathcal{Y}, \beta, \lambda] \cdot \mathbb{P}[\beta, \lambda] \quad (3)$$

Now, we already know how to compute  $\mathbb{P}[\mathbf{w} \mid X, \mathbf{y}, \beta, \lambda]$  if  $\beta$  and  $\lambda$  are known. However, computing  $\mathbb{P}[\beta, \lambda \mid X, \mathbf{Y}]$  is intractable.

$$\mathbb{P}[\beta, \lambda \mid X, \mathbf{Y}] = \frac{\mathbb{P}[\mathbf{y} \mid X, \beta, \lambda] \cdot \mathbb{P}[\beta] \cdot \mathbb{P}[\lambda]}{\mathbb{P}[\mathbf{y} \mid X]}$$

From the arguments in the discussion in [Section 1](#), we know that this quantity is intractable, therefore, we need a work-around for this.

A possible solution is (as already stated) take only the point estimate of  $\beta$  and  $\lambda$ . Therefore, we only need to find an estimate of  $\beta$  and  $\lambda$  and plug that into the posterior for simple Bayesian Linear Regression (discussed in the previous scribe) *i.e.* assume  $\alpha$  and  $\beta$  to be known.

$$\mathbb{P}[\mathbf{w} \mid \mathbf{y}, X, \hat{\beta}, \hat{\lambda}] = \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}, \Sigma) \quad (4)$$

where

$$\begin{aligned} \Sigma &= (\hat{\lambda} \mathbf{I} + \hat{\beta} X X^T)^{-1} \\ \boldsymbol{\mu} &= \Sigma \left( \hat{\beta} \sum_{i \in [n]} y^i \mathbf{x}^i \right) \end{aligned}$$

Therefore, we now find the point estimate of  $\beta$  and  $\lambda$

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} \mathbb{P}[\beta, \lambda \mid X, \mathbf{y}]$$

Note that this is MAP estimation. Therefore, we need to assume some prior on  $\beta$  and  $\lambda$ . However, we simply assume that the priors on  $\beta$  and  $\lambda$  are uniform and therefore uninformative. Thus, the MAP estimation boils down to MLE estimation.

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} \mathbb{P}[\mathbf{y} \mid X, \beta, \lambda] \quad (5)$$

**Exercise 4.2.** Prove that if the prior is a uniform distribution, then MAP estimate is equal to the MLE estimate.

**Remark.** We can also do MAP estimation on  $\beta$  and  $\lambda$ , if we assume some informative prior for them. The method could be called MAP-II

At this point, it may seem like we are done with the estimation of the hyperparameters. However, we are missing a point. The RHS in Equation 5 does not have a closed form. This is evident from the following equations.

$$\mathbb{P}[\mathbf{y} \mid X, \beta, \lambda] = \int_{\mathbf{w}} \mathbb{P}[\mathbf{y} \mid X, \mathbf{w}, \beta] \cdot \mathbb{P}[\mathbf{w} \mid \lambda] d\mathbf{w}$$

We can integrate this, using properties of Gaussian Distributions, to get the following term, the intermediate steps for which are left as an exercise.

$$\mathbb{P}[\mathbf{y} \mid X, \beta, \lambda] = \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \beta^{-1} \mathbf{I} + \lambda^{-1} X X^T) \quad (6)$$

Now MLE solution of this is generally not in closed form, hence, we require an alternate technique, known as Alternating Minimization, discussed in the next section.

### 3. Alternating Optimization

Suppose we already have an estimate of  $\beta$  and  $\lambda$ . Then we know we can easily compute  $\mathbb{P}[\mathbf{w} \mid X, \mathbf{y}, \hat{\beta}, \hat{\lambda}]$ . Also, if we have this probability, then we can easily perform MLE estimation on the term shown in Equation 6.

The above suggests that we can have an alternating approach to solving for the optimal  $\beta$  and  $\lambda$  and at the same time computing the posterior  $\mathbf{w} \mid \mathbf{y}, X$  i.e. assume some  $\alpha^0$  and  $\beta^0$ , compute posterior in Equation 4 and then solve for Equation 6 iteratively. We give the algorithm for the alternating optimization below

#### Algorithm 1: Alternating Optimization for MLE-II

1. Initialize  $\{\hat{\beta}\}^0$  and  $\{\hat{\lambda}\}^0$

2. Repeat until convergence for  $t = 1, 2 \dots$

(a) Estimate the posterior over  $\{\mathbf{w}\}^t$  as

$$\mathbb{P}[\{\mathbf{w}\}^t \mid \mathbf{y}, X, \{\hat{\beta}\}^{t-1}, \{\hat{\lambda}\}^{t-1}] = \mathcal{N}(\{\mathbf{w}\}^t \mid \boldsymbol{\mu}, \Sigma)$$

where  $\boldsymbol{\mu}$  and  $\Sigma$  are as defined in Equation 4

(b) Estimate  $\hat{\lambda}^t$  and  $\hat{\beta}^t$  as

$$\begin{aligned} \{\hat{\lambda}\}^t &= \frac{\gamma}{\boldsymbol{\mu}^T \boldsymbol{\mu}} \\ \{\hat{\beta}\}^t &= \frac{N - \gamma}{\sum_{i \in [n]} (y^i - \boldsymbol{\mu}^T \mathbf{x}^i)^2} \end{aligned}$$

where  $\gamma$  is a function of the eigenvalues of  $\Sigma$

Note that we do not discuss the intricacies of the algorithm or even the exact form of the algorithm, however the general idea of Alternating Optimization should be clear from this example.

**Remark.** The optimization function in Equation 4 is a convex function, and hence, we can expect that function to always converge. Therefore, we say that Alternating Optimization for this example will always converge to a minima (local or global)

<+>