

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

## Parameter Estimation in Probabilistic Models

### 1. Point Estimate

Point estimate is a computationally cheaper alternative to fully Bayesian Inference. It gives the single best estimate to an optimization problem. There are two types of Point Estimates – Maximum Likelihood Estimation (MLE) and Maximum-a-Posteriori (MAP)

#### 1.1 Point Estimation via MLE

MLE finds the estimate that maximizes the log-Likelihood ( $\log(\mathbb{P}[\mathbf{X} | \theta])$ ). We can write the maximization function as

$$\mathcal{L}[\theta] = \log(\mathbb{P}[\mathbf{X} | \theta])$$

If the data points are independent, we can write this the joint probability as

$$\begin{aligned} \mathbb{P}[\mathbf{X} | \theta] &= \prod_{n=1}^N \mathbb{P}[\mathbf{x}^n | \theta] \\ \Rightarrow \mathcal{L}[\theta] &= \sum_{n=1}^N \log(\mathbb{P}[\mathbf{x}^n | \theta]) \end{aligned}$$

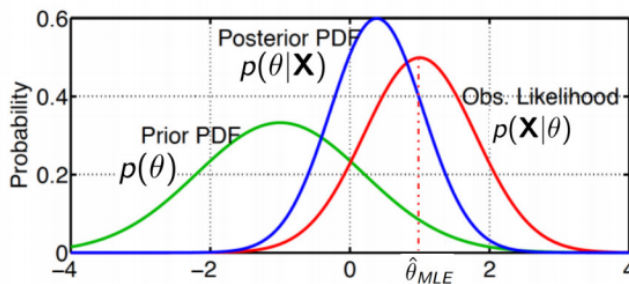


Figure 1: Visual Look at MLE Estimate

**Note.** MLE is consistent, i.e. as  $N \rightarrow \infty$  the value of  $\hat{\theta}_{MLE}$  converges to the real  $\theta$

#### 1.2 Point Estimate via MAP

MAP finds the estimate that maximizes the log-posterior-probability ( $\log(\mathbb{P}[\theta | \mathbf{X}])$ ). We can write the maximization function as

$$\begin{aligned}
\mathcal{L}[\boldsymbol{\theta}] &= \log(\mathbb{P}[\boldsymbol{\theta} \mid \mathbf{X}]) \\
&= \log(\mathbb{P}[\mathbf{X} \mid \boldsymbol{\theta}]) + \log(\mathbb{P}[\boldsymbol{\theta}])
\end{aligned}$$

Again, if the data points are independent, we can write

$$\begin{aligned}
\mathbb{P}[\mathbf{X} \mid \boldsymbol{\theta}] &= \prod_{n=1}^N \mathbb{P}[\mathbf{x}^n \mid \boldsymbol{\theta}] \\
\Rightarrow \mathcal{L}[\boldsymbol{\theta}] &= \sum_{n=1}^N (\log(\mathbb{P}[\mathbf{x}^n \mid \boldsymbol{\theta}])) + \log(\mathbb{P}[\boldsymbol{\theta}])
\end{aligned}$$

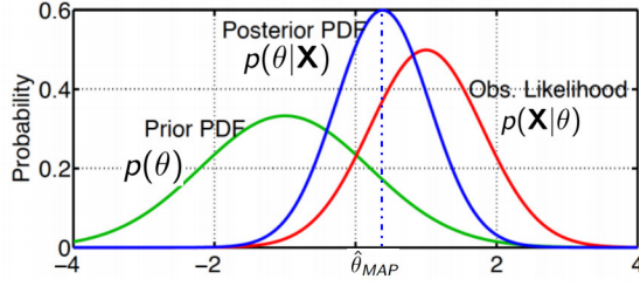


Figure 2: Visual Look at MAP Estimate

**Note.** When the prior is uniform, the MLE and MAP estimates are the same. MAP estimate is not Bayesian, as it still gives a point estimate (can be considered as optimization with regularization). We are still not modelling the uncertainty in the parameter estimation.

### 1.3 Point Estimate via Loss Minimization Function

Instead of maximizing the log-probability, we can also minimize the negative-log-probability. This is the loss function.

In case of both MLE and MAP, the loss function is the NLL (Negative Log Likelihood). For MLE, we minimize the following function

$$L(\boldsymbol{\theta}) = \text{NLL}(\boldsymbol{\theta}) = -\log(\mathbb{P}[\mathbf{X} \mid \boldsymbol{\theta}])$$

In case of MAP, we minimize the loss function with a regularizer term

$$L(\boldsymbol{\theta}) = \text{NLL}(\boldsymbol{\theta}) - \log(\mathbb{P}[\boldsymbol{\theta}])$$

**Note.** NLL is the ‘loss function’ and  $-\log(\mathbb{P}[\boldsymbol{\theta}])$  is the regularizer term

Thus, MLE is like empirical risk / loss minimization (ERM) and MAP is like regularized ERM.

### 1.4 Point Estimate Example: Coin-Toss

Consider a simple problem. We toss a coin  $N$  times. The task is to predict the next outcome. Assume consecutive coin tosses to be independent.

Say the outcome of toss is our random variable  $x$ . Since there are two options, the logical choice is a Bernoulli distribution. Hence we can write the probability function of the random variable  $x \in \{0, 1\}$  (where 1 represents ‘Head’ and 0 represents ‘Tail’) conditioned on the bias  $\theta$  as follows

$$\mathbb{P}[x] = \theta^x (1 - \theta)^{1-x}$$

Since this is the probability of our data points, we can write the log-likelihood as follows

$$\log(\mathbf{X} \mid \theta) = \sum_{n=1}^N x^n \log(\theta) + (1 - x^n) \log(1 - \theta)$$

Maximizing the log-likelihood, we get the MLE estimate  $(\hat{\theta}_{MLE})$

$$\hat{\theta}_{MLE} = \frac{\sum_{n=1}^N x^n}{N}$$

In this case, MLE is simply the fraction of heads. The problem with MLE is that it does not express our prior beliefs about  $\theta$ . This can be problematic if the number of samples is very less (the number of heads might be zero)

To avoid this, we use the MAP estimate. Here, we are faced with another problem of inductive bias. We need to choose the prior distribution. In our case, we need to find a probability function for the random variable  $\theta$ .

Since  $\theta \in (0, 1)$ , we can assume a Beta prior

$$\mathbb{P}[\theta \mid \alpha, \beta] = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

Here,  $\alpha$  and  $\beta$  are hyperparameters and are generally assumed to be known, however can also be estimated either using MLE-II or EM Method (Discussed Later). Now, for the MAP estimation, we have the following minimizing term

$$L = \sum_{n=1}^N x^n (\log(\theta) + (1 - x^n) \log(1 - \theta)) + (\alpha - 1) \log(\theta) + (\beta - 1) \log(1 - \theta)$$

Therefore, we have the MAP estimate as follows

$$\hat{\theta}_{MAP} = \frac{\sum_{n=1}^N x^n + \alpha - 1}{N + \alpha + \beta - 2}$$

Intuitively, the hyperparameters represent the pseudo-observations *i.e.*  $\alpha - 1$  and  $\beta - 1$  are the number of heads and tails, respectively, before actually seeing any data.

## 2. Bayesian Inference

Although harder, for certain problems, we can compute the posterior probability of the parameters. This gives us a more complete picture than point estimates, as it allows us to model the uncertainty in the predicted value of the parameters.

The cases in which it is possible to do Bayesian Inference are when the prior and the likelihood probabilities are *conjugate*. In other cases, we might be able to approximate the posterior using approximate Bayesian Inference methods such as MCMC and Variational Bayes.

## 2.1 Parameter Estimation via Bayesian Inference

We can use the Bayes Rule to compute the posterior distribution over the parameters

$$\mathbb{P}[\boldsymbol{\theta} \mid \mathbf{X}] = \frac{\mathbb{P}[\mathbf{X} \mid \boldsymbol{\theta}] \mathbb{P}[\boldsymbol{\theta}]}{\mathbb{P}[\mathbf{X}]}$$

The denominator *i.e.*  $\mathbb{P}[\mathbf{X}]$  is the marginal likelihood. Since the marginal likelihood is independent of  $\boldsymbol{\theta}$ , we generally only compute the numerator and then integrate (either exact or approximate) it to compute the marginal likelihood.

$$\mathbb{P}[\mathbf{X}] = \int_{\boldsymbol{\theta}} \mathbb{P}[\mathbf{X} \mid \boldsymbol{\theta}] \mathbb{P}[\boldsymbol{\theta}] d\boldsymbol{\theta}$$

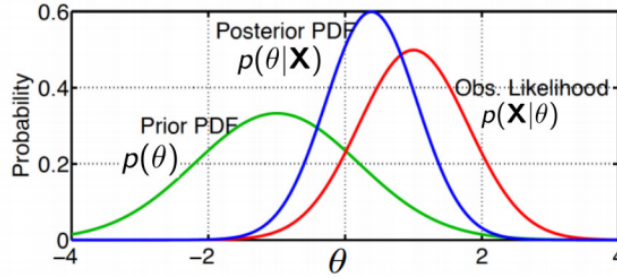


Figure 3: Visual Look at Bayesian Inference

In general, the posterior is hard to compute, as  $\mathbb{P}[\mathbf{X}]$  can be intractable. It is, however, easier if the prior and the likelihood terms and conjugate.

## 2.2 Inferring the Full Posterior: Coin-Toss Example

The problem is the same as described earlier, and we use the same terms for the likelihood (Bernoulli  $(x^n \mid \theta)$ ) and the prior (Beta  $(\theta \mid \alpha, \beta)$ ).

The posterior can then be computed as follows

$$\begin{aligned} \mathbb{P}[\theta \mid \mathbf{X}] &= \frac{\mathbb{P}[\mathbf{X} \mid \theta] \mathbb{P}[\theta]}{\mathbb{P}[\mathbf{X}]} \\ &\propto \mathbb{P}[\mathbf{X} \mid \theta] \mathbb{P}[\theta] \\ &\propto \theta^{\sum_{n=1}^N x^n + \alpha - 1} (1 - \theta)^{\beta + N - 1 - \sum_{n=1}^N x^n} \end{aligned}$$

After normalizing (or through normal inspection), we can say that the posterior is also Beta

$$\mathbb{P}[\theta \mid \mathbf{X}] = \text{Beta}\left(\theta \mid \sum_{n=1}^N x^n + \alpha - 1, \beta + N - 1 - \sum_{n=1}^N x^n\right)$$

Here, the posterior has the same form as the prior (both Beta): *property of conjugate priors*

## 2.3 Using the Posterior for Prediction

We need to predict the probability of getting a head in the next try. Essentially we need to predict the value of  $\mathbb{P}[x = 1 \mid \mathbf{X}]$

In order to make predictions, we can use posterior averaging, *i.e.* compute the predicted probabilities (of getting a head) over all values of  $\theta$  weighted on their posterior probability (Expected Value)

$$\begin{aligned}
\mathbb{P}[x = 1 \mid \mathbf{X}] &= \int_{\theta} \mathbb{P}[x = 1 \mid \theta, \mathbf{X}] \mathbb{P}[\theta \mid \mathbf{X}] d\theta \\
&= \int_{\theta} \theta \cdot \mathbb{P}[\theta \mid \mathbf{X}] d\theta \\
&= \mathbb{E}_{\mathbb{P}[\theta \mid \mathbf{X}]}[\theta] \\
&= \frac{\sum_{n=1}^N x^n + \alpha}{\alpha + \beta + N}
\end{aligned}$$

Hence,  $x \sim \text{Bernoulli}\left(x \mid \mathbb{E}_{\mathbb{P}[\theta \mid \mathbf{X}]}[\theta]\right)$

**Note.** The predicted distribution doesn't depend on only one point estimate, and is an expectation over the complete distribution

<++>

## 2.4 Examples of Conjugate Priors

Many pairs of distributions are conjugate to each other such as

- Bernoulli (likelihood) + Beta (prior)  $\longrightarrow$  Beta posterior
- Binomial (likelihood) + Beta (prior)  $\longrightarrow$  Beta posterior
- Multinomial (likelihood) + Dirichlet (prior)  $\longrightarrow$  Dirichlet posterior
- Poisson (likelihood) + Gamma (prior)  $\longrightarrow$  Gamma posterior
- Gaussian (likelihood) + Gaussian (prior)  $\longrightarrow$  Gaussian posterior