

Machine Learning Lecture Notes

Predrag Radivojac

January 25, 2015

1 Basic Principles of Parameter Estimation

In probabilistic modeling, we are typically presented with a set of observations and the objective is to find a model, or function, \hat{M} that shows good agreement with the data and respects certain additional requirements. We shall roughly categorize these requirements into three groups: (i) the ability to generalize well, (ii) the ability to incorporate prior knowledge and assumptions into modeling, and (iii) scalability. First, the model should be able to stand the test of time; that is, its performance on the previously unseen data should not deteriorate once this new data is presented. Models with such performance are said to generalize well. Second, \hat{M} must be able to incorporate information about the model space \mathcal{M} from which it is selected and the process of selecting a model should be able to accept training “advice” from an analyst. Finally, when large amounts of data are available, learning algorithms must be able to provide solutions in reasonable time given the resources such as memory or CPU power. In summary, the choice of a model ultimately depends on the observations at hand, our experience with modeling real-life phenomena, and the ability of algorithms to find good solutions given limited resources.

An easy way to think about finding the “best” model is through learning parameters of a distribution. Suppose we are given a set of observations $\mathcal{D} = \{x_i\}_{i=1}^n$, where $x_i \in \mathbb{R}$ and have knowledge that \mathcal{M} is a family of all univariate Gaussian distributions, e.g. $\mathcal{M} = \text{Gaussian}(\mu, \sigma^2)$, with $\mu \in \mathbb{R}$ and $\sigma \in \mathbb{R}^+$. In this case, the problem of finding the best model (by which we mean function) can be seen as finding the best parameters μ^* and σ^* , i.e. the problem can be seen as *parameter estimation*. We call this process estimation because the typical assumption is that the data was generated by an unknown model from \mathcal{M} whose parameters we are trying to recover from data.

We will formalize parameter estimation using probabilistic techniques and will subsequently find solutions through optimization, occasionally with constraints in the parameter space. The main assumption throughout this part will be that the set of observations \mathcal{D} was generated (or collected) independently and according to the same distribution $p_X(x)$. The statistical framework for model inference is shown in Figure 1.

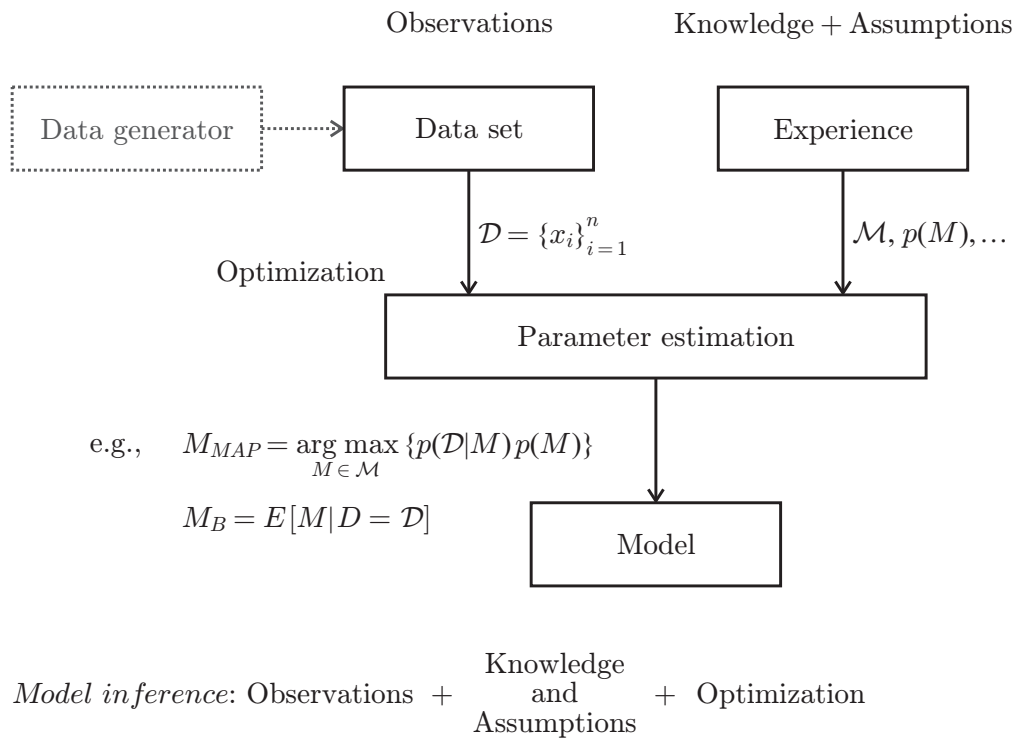


Figure 1: Statistical framework for model inference. The estimates of the parameters are made using a set of observations \mathcal{D} as well as experience in the form of model space \mathcal{M} , prior distribution $p(M)$, or specific starting solutions in the optimization step.

1.1 Maximum a posteriori and maximum likelihood estimation

The idea behind the *maximum a posteriori* (MAP) estimation is to find the most probable model for the observed data. Given the data set \mathcal{D} , we formalize the MAP solution as

$$M_{MAP} = \arg \max_{M \in \mathcal{M}} \{p(M|\mathcal{D})\},$$

where $p(M|\mathcal{D})$ is the *posterior distribution* of the model given the data. In discrete model spaces, $p(M|\mathcal{D})$ is the probability mass function and the MAP estimate is exactly the most probable model. Its counterpart in continuous spaces is the model with the largest value of the posterior density function. Note that we use words *model*, which is a function, and its *parameters*, which are the coefficients of that function, somewhat interchangeably. However, we should keep in mind the difference, even if only for pedantic reasons.

To calculate the posterior distribution we start by applying the Bayes rule as

$$p(M|\mathcal{D}) = \frac{p(\mathcal{D}|M) \cdot p(M)}{p(\mathcal{D})}, \quad (1)$$

where $p(\mathcal{D}|M)$ is called the *likelihood* function, $p(M)$ is the *prior* distribution of the model, and $p(\mathcal{D})$ is the *marginal* distribution of the data. Notice that we use \mathcal{D} for the observed data set, but that we usually think of it as a realization of a multidimensional random variable D drawn according to some distribution $p(\mathcal{D})$. Using the formula of total probability, we can express $p(\mathcal{D})$ as

$$p(\mathcal{D}) = \begin{cases} \sum_{M \in \mathcal{M}} p(\mathcal{D}|M)p(M) & M : \text{discrete} \\ \int_{\mathcal{M}} p(\mathcal{D}|M)p(M)dM & M : \text{continuous} \end{cases}$$

Therefore, the posterior distribution can be fully described using the likelihood and the prior. The field of research and practice involving ways to determine this distribution and optimal models is referred to as *inferential statistics*. The posterior distribution is sometimes referred to as inverse probability.

Finding M_{MAP} can be greatly simplified because $p(\mathcal{D})$ in the denominator does not affect the solution. We shall re-write Eq. (1) as

$$p(M|\mathcal{D}) = \frac{p(\mathcal{D}|M) \cdot p(M)}{p(\mathcal{D})} \\ \propto p(\mathcal{D}|M) \cdot p(M),$$

where \propto is the proportionality symbol. Thus, we can find the MAP solution by solving the following optimization problem

$$M_{MAP} = \arg \max_{M \in \mathcal{M}} \{p(\mathcal{D}|M)p(M)\}.$$

In some situations we may not have a reason to prefer one model over another and can think of $p(M)$ as a constant over the model space \mathcal{M} . Then, the maximum a posteriori estimation reduces to the maximization of the likelihood function; i.e.

$$M_{ML} = \arg \max_{M \in \mathcal{M}} \{p(\mathcal{D}|M)\}.$$

We will refer to this solution as the *maximum likelihood* solution. Formally speaking, the assumption that $p(M)$ is constant is problematic because a uniform distribution cannot be always defined (say, over \mathbb{R}). Thus, it may be useful to think of the maximum likelihood approach as a separate technique, rather than a special case of MAP estimation, with a conceptual caveat.

Observe that MAP and ML approaches report solutions corresponding to the mode of the posterior distribution and the likelihood function, respectively. We shall later contrast this estimation technique with the view of the Bayesian statistics in which the goal is to minimize the posterior risk. Such estimation typically results in calculating conditional expectations, which can be complex integration problems. From a different point of view, MAP and ML estimates are called *point estimates*, as opposed to estimated that report confidence intervals for particular group of parameters.

Example 8. Suppose data set $\mathcal{D} = \{2, 5, 9, 5, 4, 8\}$ is an i.i.d. sample from a Poisson distribution with an unknown parameter λ_t . Find the maximum likelihood estimate of λ_t .

The probability density function of a Poisson distribution is expressed as $p(x|\lambda) = \lambda^x e^{-\lambda}/x!$, with some parameter $\lambda \in \mathbb{R}^+$. We will estimate this parameter as

$$\lambda_{ML} = \arg \max_{\lambda \in (0, \infty)} \{p(\mathcal{D}|\lambda)\}.$$

We can write the likelihood function as

$$\begin{aligned} p(\mathcal{D}|\lambda) &= p(\{x_i\}_{i=1}^n | \lambda) \\ &= \prod_{i=1}^n p(x_i|\lambda) \\ &= \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!}. \end{aligned}$$

To find λ that maximizes the likelihood, we will first take a logarithm (a monotonic function) to simplify the calculation, then find its first derivative with respect to λ , and finally equate it with zero to find the maximum. Specifically, we express the log-likelihood $l(\mathcal{D}, \lambda) = \ln p(\mathcal{D}|\lambda)$ as

$$l(\mathcal{D}, \lambda) = \ln \lambda \sum_{i=1}^n x_i - n\lambda - \sum_{i=1}^n \ln(x_i!)$$

and proceed with the first derivative as

$$\begin{aligned} \frac{\partial l(\mathcal{D}, \lambda)}{\partial \lambda} &= \frac{1}{\lambda} \sum_{i=1}^n x_i - n \\ &= 0. \end{aligned}$$

By substituting $n = 6$ and values from \mathcal{D} , we can compute the solution as

$$\begin{aligned} \lambda_{ML} &= \frac{1}{n} \sum_{i=1}^n x_i \\ &= 5.5, \end{aligned}$$

which is simply a sample mean. The second derivative of the likelihood function is always negative because λ must be positive; thus, the previous expression indeed maximizes the likelihood. We have ignored the anomalous cases when \mathcal{D} contains all zeros.

□

Example 9. Let $\mathcal{D} = \{2, 5, 9, 5, 4, 8\}$ again be an i.i.d. sample from $\text{Poisson}(\lambda_t)$, but now we are also given additional information. Suppose the prior knowledge about λ_t can be expressed using a gamma distribution $\Gamma(x|k, \theta)$ with parameters $k = 3$ and $\theta = 1$. Find the maximum a posteriori estimate of λ_t .

First, we write the probability density function of the gamma family as

$$\Gamma(x|k, \theta) = \frac{x^{k-1} e^{-\frac{x}{\theta}}}{\theta^k \Gamma(k)},$$

where $x > 0$, $k > 0$, and $\theta > 0$. $\Gamma(k)$ is the gamma function that generalizes the factorial function; when k is an integer, we have $\Gamma(k) = (k-1)!$. The MAP estimate of the parameters can be found as

$$\lambda_{MAP} = \arg \max_{\lambda \in (0, \infty)} \{p(\mathcal{D}|\lambda)p(\lambda)\}.$$

As before, we can write the likelihood function as

$$p(\mathcal{D}|\lambda) = \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!}$$

and the prior distribution as

$$p(\lambda) = \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)}.$$

Now, we can maximize the logarithm of the posterior distribution $p(\lambda|\mathcal{D})$ using

$$\begin{aligned} \ln p(\lambda|\mathcal{D}) &\propto \ln p(\mathcal{D}|\lambda) + \ln p(\lambda) \\ &= \ln \lambda \left(k - 1 + \sum_{i=1}^n x_i\right) - \lambda \left(n + \frac{1}{\theta}\right) - \sum_{i=1}^n \ln x_i! - k \ln \theta - \ln \Gamma(k) \end{aligned}$$

to obtain

$$\begin{aligned} \lambda_{MAP} &= \frac{k - 1 + \sum_{i=1}^n x_i}{n + \frac{1}{\theta}} \\ &= 5 \end{aligned}$$

after incorporating all data.

A quick look at λ_{MAP} and λ_{ML} suggests that as n grows, both numerators and denominators in the expressions above become increasingly more similar. Unfortunately, equality between λ_{MAP} and λ_{ML} when $n \rightarrow \infty$ cannot be guaranteed because it could occur by (a distant) chance that $s_n = \sum_{i=1}^n x_i$ does not grow with n . One way to proceed is to calculate the expectation of the difference between λ_{MAP} and λ_{ML} and then investigate what happens when $n \rightarrow \infty$. Let us do this formally.

We shall first note that both estimates can be considered to be random variables. This follows from the fact that \mathcal{D} is assumed to be an i.i.d. sample from a Poisson distribution with some true parameter λ_t . With this, we can write $\lambda_{MAP} = (k - 1 + S_n)/(n + \frac{1}{\theta})$ and $\lambda_{ML} = S_n/n$, where $S_n = \sum_{i=1}^n X_i$ and $X_i \sim \text{Poisson}(\lambda_t)$. We shall now prove the following convergence (in mean) result

$$\lim_{n \rightarrow \infty} E_D [|\lambda_{MAP} - \lambda_{ML}|] = 0,$$

where the expectation is performed with respect to the random variable D . Let us first express the absolute difference between the two estimates as

$$\begin{aligned} |\lambda_{MAP} - \lambda_{ML}| &= \left| \frac{k - 1 + S_n}{n + 1/\theta} - \frac{S_n}{n} \right| \\ &= \left| \frac{k - 1}{n + 1/\theta} - \frac{S_n}{n(n + 1/\theta)} \right| \\ &\leq \frac{|k - 1|}{n + 1/\theta} + \frac{S_n}{n(n + 1/\theta)} \\ &= \epsilon \end{aligned}$$

where ϵ is a random variable that bounds the absolute difference between λ_{MAP} and λ_{ML} . We shall now express the expectation of ϵ as

$$\begin{aligned}
E[\epsilon] &= \frac{1}{n(n+1/\theta)} \cdot E\left[\sum_{i=1}^n X_i\right] + \frac{|k-1|}{n+1/\theta} \\
&= \frac{1}{n+1/\theta} \cdot \lambda_t + \frac{|k-1|}{n+1/\theta}
\end{aligned}$$

and therefore calculate that

$$\lim_{n \rightarrow \infty} E[\epsilon] = 0.$$

This result shows that the MAP estimate approaches the ML solution for large data sets. In other words, large data diminishes the importance of prior knowledge. This is an important conclusion because it simplifies mathematical apparatus necessary for practical inference. \square

1.1.1 The relationship with Kullback-Leibler divergence

We now investigate the relationship between maximum likelihood estimation and Kullback-Leibler divergence. Kullback-Leibler divergence between two probability distributions $p(x)$ and $q(x)$ is defined as

$$D_{KL}(p||q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx.$$

In information theory, Kullback-Leibler divergence has a natural interpretation of the inefficiency of signal compression when the code is constructed using a suboptimal distribution $q(x)$ instead of the correct (but unknown) distribution $p(x)$ according to which the data has been generated. However, more often than not, Kullback-Leibler divergence is simply considered to be a measure of divergence between two probability distributions. Although this divergence is not a metric (it is not symmetric and does not satisfy the triangle inequality) it has important theoretical properties in that (i) it is always non-negative and (ii) it is equal to zero if and only if $p(x) = q(x)$.

Consider now a divergence between an estimated probability distribution $p(x|\theta)$ and an underlying (true) distribution $p(x|\theta_t)$ according to which the data set $\mathcal{D} = \{x_i\}_{i=1}^n$ was generated. The Kullback-Leibler divergence between $p(x|\theta)$ and $p(x|\theta_t)$ is

$$\begin{aligned}
D_{KL}(p(x|\theta_t)||p(x|\theta)) &= \int_{-\infty}^{\infty} p(x|\theta_t) \log \frac{p(x|\theta_t)}{p(x|\theta)} dx \\
&= \int_{-\infty}^{\infty} p(x|\theta_t) \log \frac{1}{p(x|\theta)} dx - \int_{-\infty}^{\infty} p(x|\theta_t) \log \frac{1}{p(x|\theta_t)} dx.
\end{aligned}$$

The second term in the above equation is simply the entropy of the true distribution and is not influenced by our choice of the model θ . The first term, on the other hand, can be expressed as

$$\int_{-\infty}^{\infty} p(x|\theta_t) \log \frac{1}{p(x|\theta)} dx = -E[\log p(x|\theta)]$$

Therefore, maximizing $E[\log p(x|\theta)]$ minimizes the Kullback-Leibler divergence between $p(x|\theta)$ and $p(x|\theta_t)$. Using the strong law of large numbers, we know that

$$\frac{1}{n} \sum_{i=1}^n \log p(x_i|\theta) \xrightarrow{a.s.} E[\log p(x|\theta)]$$

when $n \rightarrow \infty$. Thus, when the data set is sufficiently large, maximizing the likelihood function minimizes the Kullback-Leibler divergence and leads to the conclusion that $p(x|\theta_{\text{ML}}) = p(x|\theta_t)$, if the underlying assumptions are satisfied. Under reasonable conditions, we can infer from it that $\theta_{\text{ML}} = \theta_t$. This will hold for families of distributions for which a set of parameters uniquely determines the probability distribution; e.g. it will not generally hold for mixtures of distributions but we will discuss this situation later. This result is only one of the many connections between statistics and information theory.

1.2 Parameter estimation for mixtures of distributions

We now investigate parameter estimation for mixture models, which is most commonly carried out using the *expectation-maximization (EM) algorithm*. As before, we are given a set of i.i.d. observations $\mathcal{D} = \{x_i\}_{i=1}^n$, with the goal of estimating the parameters of the mixture distribution

$$p(x|\theta) = \sum_{j=1}^m w_j p(x|\theta_j).$$

In the equation above, we used $\theta = (w_1, w_2, \dots, w_m, \theta_1, \theta_2, \dots, \theta_m)$ to combine all parameters. Just to be more concrete, we shall assume that each $p(x_i|\theta_j)$ is an exponential distribution with parameter λ_j , i.e. $p(x|\theta_j) = \lambda_j e^{-\lambda_j x}$, where $\lambda_j > 0$. Finally, we shall assume that m is given and will address simultaneous estimation of θ and m later.

Let us attempt to find the maximum likelihood solution first. By plugging the formula for $p(x|\theta)$ into the likelihood function we obtain

$$\begin{aligned} p(\mathcal{D}|\theta) &= \prod_{i=1}^n p(x_i|\theta) \\ &= \prod_{i=1}^n \left(\sum_{j=1}^m w_j p(x_i|\theta_j) \right), \end{aligned} \tag{2}$$

which, unfortunately, is difficult to maximize using differential calculus (why?). Note that although $p(\mathcal{D}|\theta)$ has $O(m^n)$ terms, it can be calculated in $O(mn)$ time as a log-likelihood.

Before introducing the EM algorithm, let us for a moment present two hypothetical scenarios that will help us to understand the algorithm. First, suppose that information is available as to which mixing component generated which data point. That is, suppose that $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ is an i.i.d. sample from some distribution $p_{XY}(x, y)$, where $y \in \mathcal{Y} = \{1, 2, \dots, m\}$ specifies the mixing component. How would the maximization be performed then? Let us write the likelihood function as

$$\begin{aligned} p(\mathcal{D}|\theta) &= \prod_{i=1}^n p(x_i, y_i|\theta) \\ &= \prod_{i=1}^n p(x_i|y_i, \theta)p(y_i|\theta) \\ &= \prod_{i=1}^n w_{y_i}p(x_i|\theta_{y_i}), \end{aligned} \tag{3}$$

where $w_j = p_Y(j) = P(Y = j)$. The log-likelihood is

$$\begin{aligned} \log p(\mathcal{D}|\theta) &= \sum_{i=1}^n (\log w_{y_i} + \log p(x_i|\theta_{y_i})) \\ &= \sum_{j=1}^m n_j \log w_j + \sum_{i=1}^n \log p(x_i|\theta_{y_i}), \end{aligned}$$

where n_j is the number of data points in \mathcal{D} generated by the j -th mixing component.

It is useful to observe here that when $\mathbf{y} = (y_1, y_2, \dots, y_n)$ is known, the internal summation operator in Eq. (2) disappears. More importantly, it follows that Eq. (3) can be maximized in a relatively straightforward manner. Let us show how. To find $\mathbf{w} = (w_1, w_2, \dots, w_m)$ we need to solve a constrained optimization problem, which we will do by using the method of Lagrange multipliers. We shall first form the Lagrangian function as

$$L(\mathbf{w}, \alpha) = \sum_{j=1}^m n_j \log w_j + \alpha \left(\sum_{j=1}^m w_j - 1 \right)$$

where α is the Lagrange multiplier. Then, by setting $\frac{\partial}{\partial w_k} L(\mathbf{w}, \alpha) = 0$ for every $k \in \mathcal{Y}$ and $\frac{\partial}{\partial \alpha} L(\mathbf{w}, \alpha) = 0$, we derive that $w_k = -\frac{n_k}{\alpha}$ and $\alpha = -n$. Thus,

$$w_k = \frac{1}{n} \sum_{i=1}^n I(y_i = k),$$

where $I(\cdot)$ is the indicator function. To find all θ_j , we recall that we assumed a

mixture of exponential distributions. Thus, we proceed by setting

$$\frac{\partial}{\partial \lambda_k} \sum_{i=1}^n \log p(x_i | \lambda_{y_i}) = 0,$$

for each $k \in \mathcal{Y}$. We obtain that

$$\lambda_k = \frac{n_k}{\sum_{i=1}^n I(y_i = k) \cdot x_i},$$

which is simply the inverse mean over those data points generated by the k -th mixture component. In summary, we observe that if the mixing component designations \mathbf{y} are known, the parameter estimation is greatly simplified. This was achieved by decoupling the estimation of mixing proportions and all parameters of the mixing distributions.

In the second hypothetical scenario, suppose that parameters θ are known, and that we would like to estimate the best configuration of the mixture designations \mathbf{y} (one may be tempted to call them “class labels”). This task looks like clustering in which cluster memberships need to be determined based on the known set of mixing distributions and mixing probabilities. To do this we can calculate the posterior distribution of \mathbf{y} as

$$\begin{aligned} p(\mathbf{y} | \mathcal{D}, \theta) &= \prod_{i=1}^n p(y_i | x_i, \theta) \\ &= \prod_{i=1}^n \frac{w_{y_i} p(x_i | \theta_{y_i})}{\sum_{j=1}^m w_j p(x_i | \theta_j)} \end{aligned} \quad (4)$$

and subsequently find the best configuration out of m^n possibilities. Obviously, because of the i.i.d. assumption each element y_i can be estimated separately and, thus, this estimation can be completed in $O(mn)$ time. The MAP estimate for y_i can be found as

$$\hat{y}_i = \arg \max_{y_i \in \mathcal{Y}} \left\{ \frac{w_{y_i} p(x_i | \theta_{y_i})}{\sum_{j=1}^m w_j p(x_i | \theta_j)} \right\}$$

for each $i \in \{1, 2, \dots, n\}$.

In reality, neither “class labels” \mathbf{y} nor the parameters θ are known. Fortunately, we have just seen that the optimization step is relatively straightforward if one of them is known. Therefore, the intuition behind the EM algorithm is to form an iterative procedure by *assuming* that either \mathbf{y} or θ is known and calculate the other. For example, we can initially pick some value for θ , say $\theta^{(0)}$, and then estimate \mathbf{y} by computing $p(\mathbf{y} | \mathcal{D}, \theta^{(0)})$ as in Eq. (4). We can refer to this estimate as $\mathbf{y}^{(0)}$. Using $\mathbf{y}^{(0)}$ we can now refine the estimate of θ to $\theta^{(1)}$ using Eq. (3). We can then iterate these two steps until convergence. In the case of mixture of exponential distributions, we arrive at the following algorithm:

1. Initialize $\lambda_k^{(0)}$ and $w_k^{(0)}$ for $\forall k \in \mathcal{Y}$
2. Calculate $y_i^{(0)} = \arg \max_{k \in \mathcal{Y}} \left\{ \frac{w_k^{(0)} p(x_i | \lambda_k^{(0)})}{\sum_{j=1}^m w_j^{(0)} p(x_i | \lambda_j^{(0)})} \right\}$ for $\forall i \in \{1, 2, \dots, n\}$
3. Set $t = 0$
4. Repeat until convergence
 - (a) $w_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n I(y_i^{(t)} = k)$
 - (b) $\lambda_k^{(t+1)} = \frac{\sum_{i=1}^n I(y_i^{(t)} = k)}{\sum_{i=1}^n I(y_i^{(t)} = k) \cdot x_i}$
 - (c) $t = t + 1$
 - (d) $y_i^{(t+1)} = \arg \max_{k \in \mathcal{Y}} \left\{ \frac{w_k^{(t)} p(x_i | \lambda_k^{(t)})}{\sum_{j=1}^m w_j^{(t)} p(x_i | \lambda_j^{(t)})} \right\}$
5. Report $\lambda_k^{(t)}$ and $w_k^{(t)}$ for $\forall k \in \mathcal{Y}$

This procedure is not quite yet the EM algorithm; rather, it is a version of it referred to as *classification EM algorithm*. In the next section we will introduce the EM algorithm.

1.2.1 The expectation-maximization algorithm

The previous procedure was designed to iteratively estimate class memberships and parameters of the distribution. In reality, it is not necessary to compute \mathbf{y} ; after all, we only need to estimate θ . To accomplish this, at each step t , we can use $p(\mathbf{y} | \mathcal{D}, \theta^{(t)})$ to maximize the *expected log-likelihood* of both \mathcal{D} and \mathbf{y}

$$E_{\mathbf{Y}}[\log p(\mathcal{D}, \mathbf{y} | \theta) | \theta^{(t)}] = \sum_{\mathbf{y}} \log p(\mathcal{D}, \mathbf{y} | \theta) p(\mathbf{y} | \mathcal{D}, \theta^{(t)}), \quad (5)$$

which can be carried out by integrating the log-likelihood function of \mathcal{D} and \mathbf{y} over the posterior distribution for \mathbf{y} in which the current values of the parameters $\theta^{(t)}$ are assumed to be known. We can now formulate the expression for the parameters in step $t + 1$ as

$$\theta^{(t+1)} = \arg \max_{\theta} \left\{ E[\log p(\mathcal{D}, \mathbf{y} | \theta) | \theta^{(t)}] \right\}. \quad (6)$$

The formula above is all that is necessary to create the update rule for the EM algorithm. Note, however, that inside of it we always have to re-compute $E_{\mathbf{Y}}[\log p(\mathcal{D}, \mathbf{y} | \theta) | \theta^{(t)}]$ function because the parameters $\theta^{(t)}$ have been updated from the previous step. We then can perform maximization. Hence the name “expectation-maximization”, although it is perfectly valid to think of the EM algorithm as an iterative maximization of expectation from Eq. (5), i.e. “expectation maximization”.

We now proceed as follows

$$\begin{aligned}
E[\log p(\mathcal{D}, \mathbf{y}|\theta)|\theta^{(t)}] &= \sum_{y_1=1}^m \cdots \sum_{y_n=1}^m \log p(\mathcal{D}, \mathbf{y}|\theta) p(\mathbf{y}|\mathcal{D}, \theta^{(t)}) \\
&= \sum_{y_1=1}^m \cdots \sum_{y_n=1}^m \sum_{i=1}^n \log p(x_i, y_i|\theta) \prod_{l=1}^n p(y_l|x_l, \theta^{(t)}) \\
&= \sum_{y_1=1}^m \cdots \sum_{y_n=1}^m \sum_{i=1}^n \log (w_{y_i} p(x_i|\theta_{y_i})) \prod_{l=1}^n p(y_l|x_l, \theta^{(t)}).
\end{aligned}$$

After several simplification steps, that we omit for space reasons, the expectation of the likelihood can be written as

$$E[\log p(\mathcal{D}, \mathbf{y}|\theta)|\theta^{(t)}] = \sum_{i=1}^n \sum_{j=1}^m \log (w_j p(x_i|\theta_j)) p_{Y_i}(j|x_i, \theta^{(t)}),$$

from which we can see that \mathbf{w} and $\{\theta_j\}_{j=1}^m$ can be separately found. In the final two steps, we will first derive the update rule for the mixing probabilities and then by assuming the mixing distributions are exponential, derive the update rules for their parameters.

To maximize $E[\log p(\mathcal{D}, \mathbf{y}|\theta)|\theta^{(t)}]$ with respect to \mathbf{w} , we observe that this is an instance of constrained optimization because it must hold that $\sum_{i=1}^m w_i = 1$. We will use the method of Lagrange multipliers; thus, for each $k \in \mathcal{Y}$ we need to solve

$$\frac{\partial}{\partial w_k} \left(\sum_{j=1}^m \log w_j \sum_{i=1}^n p_{Y_i}(j|x_i, \theta^{(t)}) + \alpha \left(\sum_{j=1}^m w_j - 1 \right) \right) = 0,$$

where α is the Lagrange multiplier. It is relatively straightforward to show that

$$w_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n p_{Y_i}(k|x_i, \theta^{(t)}). \quad (7)$$

Similarly, to find the solution for the parameters of the mixture distributions, we obtain that

$$\lambda_k^{(t+1)} = \frac{\sum_{i=1}^n p_{Y_i}(k|x_i, \theta^{(t)})}{\sum_{i=1}^n x_i p_{Y_i}(k|x_i, \theta^{(t)})} \quad (8)$$

for $k \in \mathcal{Y}$. As previously shown, we have

$$p_{Y_i}(k|x_i, \theta^{(t)}) = \frac{w_k^{(t)} p(x_i|\lambda_k^{(t)})}{\sum_{j=1}^m w_j^{(t)} p(x_i|\lambda_j^{(t)})}, \quad (9)$$

which can be computed and stored as an $n \times m$ matrix. In summary, for the mixture of m exponential distributions, we summarize the EM algorithm by combining Eqs. (7-9) as follows:

1. Initialize $\lambda_k^{(0)}$ and $w_k^{(0)}$ for $\forall k \in \mathcal{Y}$
2. Set $t = 0$
3. Repeat until convergence
 - (a) $p_{Y_i}(k|x_i, \theta^{(t)}) = \frac{w_k^{(t)} p(x_i|\lambda_k^{(t)})}{\sum_{j=1}^n w_j^{(t)} p(x_i|\lambda_j^{(t)})}$ for $\forall(i, k)$
 - (b) $w_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n p_{Y_i}(k|x_i, \theta^{(t)})$
 - (c) $\lambda_k^{(t+1)} = \frac{\sum_{i=1}^n x_i p_{Y_i}(k|x_i, \theta^{(t)})}{\sum_{i=1}^n x_i p_{Y_i}(k|x_i, \theta^{(t)})}$
 - (d) $t = t + 1$
4. Report $\lambda_k^{(t)}$ and $w_k^{(t)}$ for $\forall k \in \mathcal{Y}$

Similar update rules can be obtained for different probability distributions; however, a separate derivatives have to be found.

It is important to observe and understand the difference between the CEM and the EM algorithms.

1.3 Bayesian estimation

Maximum a posteriori and maximum likelihood approaches report the solution that corresponds to the mode of the posterior distribution and the likelihood function, respectively. This approach, however, does not consider the possibility of skewed distributions, multimodal distributions or simply large regions with similar values of $p(M|\mathcal{D})$. Bayesian estimation addresses those concerns.

The main idea in Bayesian statistics is minimization of the *posterior risk*

$$R = \int_{\mathcal{M}} \ell(M, \hat{M}) \cdot p(M|\mathcal{D}) dM$$

where \hat{M} is our estimate and $\ell(M, \hat{M})$ is some loss function between two models. When $\ell(M, \hat{M}) = (M - \hat{M})^2$ (ignore the abuse of notation), we can minimize the posterior risk as follows

$$\begin{aligned} \frac{\partial}{\partial \hat{M}} R &= 2\hat{M} - 2 \int_{\mathcal{M}} M \cdot p(M|\mathcal{D}) dM \\ &= 0 \end{aligned}$$

from which it can be derived that the minimizer of the posterior risk is the posterior mean function, i.e.

$$\begin{aligned} M_B &= \int_{\mathcal{M}} M \cdot p(M|\mathcal{D}) dM \\ &= E_M[M|\mathcal{D}]. \end{aligned}$$

We shall refer to M_B as the Bayes estimator. It is important to mention that computing the posterior mean usually involves solving complex integrals. In some situations, these integrals can be solved analytically; in others, numerical integration is necessary.

Example 11. Let $\mathcal{D} = \{2, 5, 9, 5, 4, 8\}$ yet again be an i.i.d. sample from $\text{Poisson}(\lambda_t)$. Suppose the prior knowledge about the parameter of the exponential distribution can be expressed using a gamma distribution with parameters $k = 3$ and $\theta = 1$. Find the Bayesian estimate of λ_t .

We want to find $E[\lambda|\mathcal{D}]$. Let us first write the posterior distribution as

$$\begin{aligned} p(\lambda|\mathcal{D}) &= \frac{p(\mathcal{D}|\lambda)p(\lambda)}{p(\mathcal{D})} \\ &= \frac{p(\mathcal{D}|\lambda)p(\lambda)}{\int_0^\infty p(\mathcal{D}|\lambda)p(\lambda)d\lambda}, \end{aligned}$$

where, as shown in previous examples, we have that

$$p(\mathcal{D}|\lambda) = \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!}$$

and

$$p(\lambda) = \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)}.$$

Before calculating $p(\mathcal{D})$, let us first note that

$$\int_0^\infty x^{\alpha-1} e^{-\beta x} dx = \frac{\Gamma(\alpha)}{\beta^\alpha}.$$

Now, we can derive that

$$\begin{aligned} p(\mathcal{D}) &= \int_0^\infty p(\mathcal{D}|\lambda)p(\lambda)d\lambda \\ &= \int_0^\infty \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!} \cdot \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)} d\lambda \\ &= \frac{\Gamma(k + \sum_{i=1}^n x_i)}{\theta^k \Gamma(k) \prod_{i=1}^n x_i! (n + \frac{1}{\theta})^{\sum_{i=1}^n x_i + k}} \end{aligned}$$

and subsequently that

$$\begin{aligned}
p(\lambda|\mathcal{D}) &= \frac{p(\mathcal{D}|\lambda)p(\lambda)}{p(\mathcal{D})} \\
&= \frac{\lambda^{\sum_{i=1}^n x_i} \cdot e^{-n\lambda}}{\prod_{i=1}^n x_i!} \cdot \frac{\lambda^{k-1} e^{-\frac{\lambda}{\theta}}}{\theta^k \Gamma(k)} \cdot \frac{\theta^k \Gamma(k) \prod_{i=1}^n x_i! (n + \frac{1}{\theta})^{\sum_{i=1}^n x_i + k}}{\Gamma(k + \sum_{i=1}^n x_i)} \\
&= \frac{\lambda^{k-1 + \sum_{i=1}^n x_i} \cdot e^{-\lambda(n + \frac{1}{\theta})} \cdot (n + \frac{1}{\theta})^{\sum_{i=1}^n x_i + k}}{\Gamma(k + \sum_{i=1}^n x_i)}.
\end{aligned}$$

Finally,

$$\begin{aligned}
E[\lambda|\mathcal{D}] &= \int_0^\infty \lambda p(\lambda|\mathcal{D}) d\lambda \\
&= \frac{k + \sum_{i=1}^n x_i}{n + \frac{1}{\theta}} \\
&= 5.14
\end{aligned}$$

which is nearly the same solution as the MAP estimate found in Example 9. \square

It is evident from the previous example that selection of the prior distribution has important implications on calculation of the posterior mean. We have not picked the gamma distribution by chance; that is, when the likelihood was multiplied by the prior, the resulting distribution remained in the same class of functions as the likelihood. We shall refer to such prior distributions as *conjugate priors*. Conjugate priors are also simplifying the mathematics; in fact, this is a major reason for their consideration. Interestingly, in addition to the Poisson distribution, the gamma distribution is a conjugate prior to the exponential distribution as well as the gamma distribution itself.