

STATISTICAL MODELLING AND INFERENCE

Sergio-Yersi Villegas Pelegrín

SEMINAR 4: Problems 4 & 5

Problem 4

We want to derive the Bayesian EM algorithm for maximizing the posterior probability $P(\theta, \pi | x)$ of the studied Bernoulli mixture model, together its corresponding priors for θ and π . To begin with, we summarize the Bayesian Bernoulli mixture model we have: for a random sample $x = \{x_1, \dots, x_n\}$, for each x_i we have:

$$x_i | \theta_k, z_i = k \sim \text{Bern}(\theta_k)$$

$$z_i | \pi \sim \text{Multinomial}(\pi)$$

$$\pi | \alpha \sim \text{Dir}(\alpha)$$

$$\theta_k | a_k, b_k \sim \text{Beta}(a_k, b_k)$$

Now, in order to be able to analyze our posterior distribution of interest, we must know the conditional independence relationships between all parameters:

- π just depends on α as a Dirichlet distribution (one of the added priors in this problem).
- z depends on π as a Multinomial distribution and then, it also depends on α , since π does.
- θ depends on the $\{a_k, b_k\}_{k=1}^K$ parameters of the Beta distribution (the other prior added in this problem).
- Finally, x depends on all the aforementioned parameters

Clearly, directly trying to compute $P(\theta, \pi | x)$ is infeasible. However, we can use Gibbs sampling and construct a Markov chain that allows us to approximate this distribution arbitrarily accurate. Specifically, we will iterate between all the parameters and sample each one of them given the rest. In the following steps, we will keep using both Bayes rule (to alternate between conditional probabilities and the probability of the intersection) and the conditional independence property (simplifying conditional probabilities by excluding the independent parameters in each case), in order to get to the solution.

Step 1: sample π given $x, z, \alpha, \theta_k, \{a_k, b_k\}_{k=1}^K$:

$$\begin{aligned} P(\pi | z, \{a_k, b_k\}_{k=1}^K, x, \alpha, \theta_k) &= P(\pi | z, \alpha) \sim P(\pi, z, \alpha) = P(z | \pi, \alpha) \cdot P(\pi, \alpha) \sim P(z | \pi, \alpha) \cdot P(\pi | \alpha) \sim \\ &\sim \text{Multinomial}(\pi) \cdot \text{Dir}(\alpha) \sim \text{Dir}(\alpha + n) \end{aligned}$$

where $n_j = \sum_{i=1}^n \mathbf{1}(z_i = j)$, meaning that it sums 1 whenever $z_i = j$.

Step 2: sample θ_k given $x, z, \pi, \alpha, \{a_k, b_k\}_{k=1}^K$:

$$\begin{aligned} P(\theta_k | z, \{a_k, b_k\}_{k=1}^K, x, \alpha, \pi) &= P(\theta_k | z_i, a_k, b_k, x_i) \sim P(\theta_k, z_i, a_k, b_k, x_i) = P(x_i | \theta_k, z_i, a_k, b_k) \cdot P(\theta_k, z_i, a_k, b_k) = \\ &= P(x_i | \theta_k, z_i) \cdot P(\theta_k, z_i, a_k, b_k) \sim P(x_i | \theta_k, z_i) \cdot P(\theta_k | z_i, a_k, b_k) = P(x_i | \theta_k, z_i) \cdot P(\theta_k | a_k, b_k) \sim \text{Bern}(\theta_k) \cdot \text{Beta}(a_k, b_k) \sim \\ &\sim \prod_{j=1}^d \theta_{kj}^{x_{ij}} \cdot (1 - \theta_{kj})^{1 - x_{ij}} \cdot \theta_{kj}^{a_{kj} - 1} \cdot (1 - \theta_{kj})^{b_{kj} - 1} \sim \prod_{j=1}^d \theta_{kj}^{a_{kj} + x_{ij} - 1} \cdot (1 - \theta_{kj})^{b_{kj} - x_{ij}} \sim \text{Beta}(a_k + x_i, 1 + b_k - x_i) \end{aligned}$$

Step 3: sample z given $x, \pi, \alpha, \theta_k, \{a_k, b_k\}_{k=1}^K$:

For each $i, z_i \in \{1, \dots, K\}$ given $x, \pi, \alpha, \theta_k, \{a_k, b_k\}_{k=1}^K$ follows a multinomial distribution such that:

$$P(z_i | \theta_k, \{a_k, b_k\}_{k=1}^K, x_i, \alpha, \pi) = \frac{\pi_k \cdot \text{Bern}(x_i | \theta_k)}{\sum_{j=1}^K \pi_k \cdot \text{Bern}(x_j | \theta_{kj})}$$

Now that we know the posterior distributions in steps 1 – 3, we can run the Gibbs sampling algorithm, where the approximate posterior probabilities can be obtained by Monte Carlo approximation.

Problem 5

As seen in *Problem 3*, the K-means algorithm arises as a particular limit of the Gaussian mixture EM-algorithm. We want to show that, for $\epsilon \rightarrow 0$, maximizing the expected complete-data log likelihood (ECL) for this model is equivalent to minimizing the distortion measure (J) for the K-means algorithm. Therefore, we will compute both derivatives and see that, when $\epsilon \rightarrow 0$, we get the same result:

Distortion measure (J)

J can be written as the following:

$$J = \sum_{i=1}^n \sum_{k=1}^K r_{ik} \cdot \|x_i - \mu_k\|^2$$

Therefore, we can minimize it in respect of μ_k , by fixing $k = k'$ (therefore not computing the sum over all k) and obtain the following result:

$$\left. \frac{\partial J}{\partial \mu_k} \right|_{\mu_{k'}} = \sum_{i=1}^n r_{ik'} \cdot 2 \cdot \|x_i - \mu_{k'}^J\| \cdot (-1) = 0 \Rightarrow \sum_{i=1}^n r_{ik'} \cdot x_i = \mu_{k'}^J \cdot \sum_{i=1}^n r_{ik'} \Rightarrow \mu_{k'}^J = \frac{\sum_{i=1}^n r_{ik'} \cdot x_i}{\sum_{i=1}^n r_{ik'}}$$

Expected complete-data log likelihood (ECL)

We can write our distribution function, since we are working with a Gaussian Mixture Model, as the following:

$$P(x_i, z_i | \theta) = \sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k)$$

where μ_k and Σ_k are the parameters of the normal distribution. Then, maximizing the expected complete-data log likelihood can be expressed as:

$$\theta = \arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{z_i | x_i} \cdot \log(P(x_i, z_i | \theta)) = \arg \max_{\mu_k, \Sigma_k} \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik} \cdot \log\left(\sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k)\right)$$

However, as seen in *Problem 3*, the Σ_k parameter can be written as $\Sigma_k = \epsilon \mathbb{I}_d$, where \mathbb{I}_d is the identity matrix for d dimensions. Therefore, Σ_k is no longer a parameter, but a constant $\Sigma_k = \epsilon^d$. Together with the fact that, again, we can compute the corresponding derivative by fixing $k = k'$ (therefore not computing the sum over all k), we could now write the maximization expression, just in respect of μ_k , as:

$$\mu_{k'} = \arg \max_{\mu_{k'}} \sum_{i=1}^n \gamma_{ik'} \cdot \left(\log\left(\frac{\pi_{k'}}{\sqrt{2\pi\epsilon^d}}\right) - \frac{\|\mathbf{x}_i - \mu_{k'}\|^2}{2\epsilon^d} \right) \equiv \arg \max_{\mu_{k'}} ECL$$

Then, we obtain the following result:

$$\begin{aligned} \frac{\partial ECL}{\partial \mu_k} \Big|_{\mu_{k'}^{ECL}} &= \sum_{i=1}^n -\frac{\gamma_{ik'} \cdot 2 \cdot \|\mathbf{x}_i - \mu_{k'}^{ECL}\|}{2\epsilon^d} = \sum_{i=1}^n -\frac{\gamma_{ik'} \cdot \|\mathbf{x}_i - \mu_{k'}^{ECL}\|}{\epsilon^d} = 0 \Rightarrow \sum_{i=1}^n \gamma_{ik'} \cdot \mathbf{x}_i = \mu_{k'}^{ECL} \cdot \sum_{i=1}^n \gamma_{ik'} \\ \mu_{k'}^{ECL} &= \frac{\sum_{i=1}^n \gamma_{ik'} \cdot \mathbf{x}_i}{\sum_{i=1}^n \gamma_{ik'}} \end{aligned}$$

As we saw in *Problem 3*: $\epsilon \rightarrow 0 \Rightarrow \gamma_{ik} \rightarrow r_{ik}$. Therefore, $\mu_{k'}^{ECL} = \mu_{k'}^J$, proving that maximizing the expected complete-data log likelihood for this model is equivalent to minimizing the distortion measure for the K-means algorithm.