

Expectation Maximization*

Notes for the bAG seminar

MIGUEL DE BENITO

Universität Augsburg

Oct. 14th, 2014

ABSTRACT. We revisit some of the ideas from Philipp's talk from the point of view of latent variables, then explain how an iterative algorithm, Expectation-Maximization, appears naturally for the estimation of the parameters. We apply it to mixtures of Gaussian and Bernoulli variables. We also say a few words about the Kullback-Leibler divergence to be able to show why EM works.

WARNING! These notes are sloppy, incomplete, inconsistent and sloppy. Also, they contain errors a.s. Did I mention they are sloppy? They are also a rip-off of [Bis06, Chapters 9 and 1.6].

TABLE OF CONTENTS

1. Why we do this	2
2. Notation and conventions	2
3. Gaussian Mixtures with latent variables	3
4. The EM algorithm	5
5. Back to Gaussian Mixtures	7
6. An example with mixtures of Bernoullis	9
7. Where to go from here	10
8. Appendix	10
8.1. A word about conditional densities	10
8.2. Problems related to the log likelihood approach	10
8.3. A few sketchy ideas from information theory	12
8.4. Maximization of the log likelihood for the complete dataset	13
9. References	14

*. This document has been written using the GNU $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$ text editor (see www.texmacs.org).

1. WHY WE DO THIS

As before, when we studied K -means, we will model a random variable representing each data point as a mixture of known densities, though not necessarily Gaussians.

1. We would very much like to generalize/improve the algorithm we saw last week, as well as understand and prove its convergence. Furthermore, not using gradient descent frees us from the issue with the parameter for step length.
2. The technique we will develop can be used for *maximum a posteriori* estimation in order to fix the problem of overfitting without recourse to ad-hoc fixes or heuristics.
3. In a fully Bayesian framework we can use EM for automatic *model selection* (e.g. in our examples with mixtures, determining the right number K of different distributions).
4. An online version of EM runs in $O(1)$ time. See [Bis06, §9.4].
5. Understanding EM sets the grounds for a larger class of techniques for deterministic parameter estimation, known as variational inference.
6. Because we can.

2. NOTATION AND CONVENTIONS

Random variables take values in \mathbb{R} or \mathbb{R}^d and are denoted with capital letters X, Z . Their **realizations** are denoted with x, z . We will always use X for the RV originating the data. Z will be reserved for the **latent variables**, e.g. those “artificially” added to model our assumptions about the internals (e.g. hidden state) of the system studied. In our particular examples Z will take values in $\{e_1, \dots, e_k\}$ (see §3). The **probability measure** is denoted by a capital P , **densities of random variables** by a small p . If X is a discrete RV, $P(X = x) = p_X(x)$. If no confusion can arise we will use the arguments of densities to distinguish among them: $p(x) = p_X(x)$ will be the density of X evaluated at x , $p(z|x) = p_{Z|X}(z|x) = p_{Z|X=x}(z)$ the density of Z “given $X = x$ ” at z , etc. (see §8.1 for a few words on conditioning on an event with probability zero). $\mathcal{N}(\mu, \Sigma)$ denotes a **Gaussian distribution** of mean μ and covariance matrix Σ , but if we add an argument x , then $\mathcal{N}(x|\mu, \Sigma)$ denotes its density, i.e.

$$\mathcal{N}(x|\mu, \Sigma) := \frac{1}{(2\pi)^{D/2}} \frac{1}{\sqrt{\det \Sigma}} \exp \left\{ -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right\}.$$

We subsume all **distribution parameters** into the letter θ , which may refer to different parameters even in the same line, but is always understood as “the parameters for the distribution where θ is used”. Conditioning wrt. θ is to be understood (for now) as a notational reminder that there are some parameters. We observe N **data points** in \mathbb{R}^d . They are i.i.d. observations of X , i.e. we have the **joint random variable** $\mathbf{X} = (X_1, \dots, X_N)$ with **realization** $\mathbf{x} = (x_1, \dots, x_N)$. \mathbf{x} is the actual data we observe. Finally, C is a constant (hopefully) independent of the relevant quantities. It may change from line to line.

3. GAUSSIAN MIXTURES WITH LATENT VARIABLES

In a first example to connect with the previous talk about K -means, we postulate a convex combination of normal distributions (a **Gaussian mixture**) as the model for each of our data points:

$$\sum_{k=1}^K \pi_k \mathcal{N}(\mu_k, \Sigma_k), \quad (1)$$

for some given $K \in \mathbb{N}$ and $\sum_{k=1}^K \pi_k = 1$, $\pi_k \in [0, 1]$. Our final objective is of course to fit the parameters $\theta = (\pi_1, \dots, \pi_k, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k)$ to the data $\mathbf{x} = (x_1, \dots, x_N)$, which is a realization of the joint $\mathbf{X} = (X_1, \dots, X_N)$, but we will now arrive at this model from a new point of view.

We introduce the **latent variable** Z with values in $\{e_1, \dots, e_K\}$ with $e_k \in \mathbb{R}^K$ basis vectors. For each data point, i.e. for each realization x_n of X , we will have an unobserved z_n . The event $Z = e_k$ should be interpreted as: the corresponding realization of X “comes from” the distribution $\mathcal{N}(\mu_k, \Sigma_k)$, i.e. we want to set

$$p(x|Z = e_k) = \mathcal{N}(x|\mu_k, \Sigma_k).$$

and we choose $p(Z = e_k) = \pi_k$, for some $\pi_k \in [0, 1]$ such that $\sum_{k=1}^K \pi_k = 1$, to be determined.

Remark: What we are doing here is determining the joint distribution $p_{X,Z}$ from the conditional $p_{X|Z}$ and the marginal p_Z . Although in this case this is no more than the product rule, in more complicated situations it is essential to study how joint densities factorize as products of conditional probabilities to reduce the complexity of the problem. We can depict the relationship between X, Z as in Figure 1.

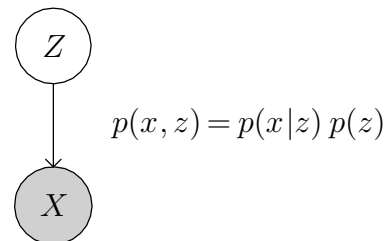


Figure 1. Graphical representation for our latent variable model.

Notice that using the fact that each realization $z = (z_1, \dots, z_K)$ of Z is one of the e_k , we may write the distribution of Z as

$$p(z) = \prod_{k=1}^K \pi_k^{z_k},$$

where of course the π_k are unknown, and the conditional $X|Z = z$ as

$$p(x|z) = \prod_{k=1}^K \mathcal{N}(x|\mu_k, \Sigma_k)^{z_k}.$$

The density of X may now be written as the marginalization of the joint $p_{X,Z}$:

$$\begin{aligned} p(x) &= \sum_{z \in \{e^1, \dots, e^K\}} p(z) p(x|z) \\ &= \sum_{k=1}^K \prod_{j=1}^K \pi_j^{\delta_{jk}} \prod_{l=1}^K \mathcal{N}(x|\mu_l, \Sigma_l)^{\delta_{lk}} \\ &= \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k). \end{aligned}$$

So we see that we recover the density (1) as we wished. From here we could proceed by attempting to maximize the likelihood of the data given the parameters, that is using the marginal:

$$\log p(\mathbf{x}|\theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k), \quad (2)$$

(recall that the X_n are i.i.d.), differentiating wrt. the parameters π_k, μ_k, Σ_k , equating to zero and solving for the parameters, but we will meet with the following two problems, already discussed in Philipp's talk, and explained in more detail below.

1. We don't obtain closed-form estimators when we differentiate (2), so we need to devise and justify an iterative scheme (e.g. Newton).
2. We suffer overfitting with collapsing variances (actually "ultrafitting"!).

So, why the did we introduce Z ? We now have a joint probability distribution $p_{X,Z}$ to work with which may (as is the case with Gaussian distributions) be more easily tractable and lead to closed form estimators when maximizing its log likelihood (see §5). However, the values z are part of the problem, so we will end up computing an expectation wrt. the posterior $p_{Z|X}$. We will show that this provides an iterative scheme (which coincides with the Newton scheme for K -means) guaranteed to increase the log likelihood.

For comparison purposes here is how Expectation-Maximization looks for a mixture of Gaussians. You'll notice the exact analogy with K -means, where the *assignment step* is now the *E-step* and the *update step* is now the *M-step*. We will derive this algorithm in §5.

Algorithm No-questions-asked-EM for Gaussian mixtures

1. Initialize the parameters $\theta = (\mu_k, \Sigma_k, \pi_k)_{k=1}^K$.
2. **E-Step:** compute the posterior $p_{\mathbf{Z}|\mathbf{X}=\mathbf{x},\theta^{\text{old}}}(\mathbf{z})$ or *responsibilities*.
3. **M-Step:** maximize $E_{\mathbf{Z}|\mathbf{x},\theta}[\log p(\mathbf{x}, \mathbf{Z}|\theta)]$ wrt. the parameters μ_k, Σ_k, π_k to obtain new values $\mu_k^{\text{new}}, \Sigma_k^{\text{new}}, \pi_k^{\text{new}}$.
4. Evaluate log likelihood, check for convergence, go to 2 if necessary.

4. THE EM ALGORITHM

The **key idea** is to use the following decomposition, true for any strictly positive probability density q :

$$\log p(\mathbf{x}|\theta) = \mathcal{L}(q, \theta) + \text{KL}(q\|p_{\mathbf{Z}|\mathbf{x},\theta}), \quad (3)$$

where KL is the Kullback-Leibler divergence (see §8.3), $p_{\mathbf{Z}|\mathbf{x},\theta} = p_{\mathbf{Z}|\mathbf{X}=\mathbf{x},\theta}$ is the distribution of the posterior of $\mathbf{Z} = (Z_1, \dots, Z_N)$ given the dataset $\mathbf{x} \in \mathbb{R}^{N \times d}$ and the choice of parameters θ and

$$\begin{aligned} \mathcal{L}(q, \theta) &:= \sum_{\mathbf{z}' \in \{e_1, \dots, e_K\}^N} q(\mathbf{z}') \log \frac{p(\mathbf{x}, \mathbf{z}'|\theta)}{q(\mathbf{z}')} \\ \text{KL}(q\|p_{\mathbf{Z}|\mathbf{x}}) &:= - \sum_{\mathbf{z}' \in \{e_1, \dots, e_K\}^N} q(\mathbf{z}') \log \frac{p(\mathbf{z}'|\mathbf{x}, \theta)}{q(\mathbf{z}')}. \end{aligned}$$

Indeed, adding both quantities and cancelling terms we have

$$\mathcal{L}(q, \theta) + \text{KL}(q\|p_{\mathbf{Z}|\mathbf{x}}) = \sum_{\mathbf{z}'} q(\mathbf{z}') \log \frac{p(\mathbf{x}, \mathbf{z}'|\theta)}{p(\mathbf{z}'|\mathbf{x}, \theta)} = \underbrace{\sum_{\mathbf{z}'} q(\mathbf{z}') \log p(\mathbf{x}|\theta)}_{=1} = \log p(\mathbf{x}|\theta).$$

Notice that because $\text{KL} \geq 0$, the functional $\mathcal{L}(q, \theta)$ is *always a lower bound* to $\log p(\mathbf{x}|\theta)$.

In each step of the algorithm we will optimize this bound: first wrt. q then wrt. θ .

The lower bound can be written as

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_{\mathbf{z}'} q(\mathbf{z}') \log p(\mathbf{x}, \mathbf{z}'|\theta) - \sum_{\mathbf{z}'} q(\mathbf{z}') \log q(\mathbf{z}') \\ &= \underbrace{E_q[\log p(\mathbf{x}, \mathbf{z}'|\theta)]}_{=C} - E_q[\log q(\mathbf{z}'|\theta)],\end{aligned}\tag{4}$$

where the constant C is the entropy of q (see §8.3) and $\log p(\mathbf{x}, \mathbf{z}'|\theta) = g(\mathbf{z}')$ a function of \mathbf{z}' alone. This last line provides the intuition behind the decomposition: *if* we had the values \mathbf{z} of the latent variable $\mathbf{Z} = (Z_1, \dots, Z_N)$, we could evaluate the joint density on the complete data set $p(\mathbf{x}, \mathbf{z}|\theta)$. Under certain assumptions it would be an easy matter to maximize its log likelihood instead of the marginal $\log p(\mathbf{x}|\theta)$ and we wouldn't need the previous decomposition (e.g. assuming independence and a Gaussian mixture, the joint is a double product which the logarithm transforms into a double sum, see (6)). But we don't have the values \mathbf{z} and all we actually know about them is contained in the posterior $p_{\mathbf{Z}|\mathbf{x},\theta}$. It is therefore sensible to compute the expected value of $\log p(\mathbf{x}, \mathbf{z}|\theta)$ under the marginal $p_{\mathbf{Z}|\mathbf{x},\theta}$, which is exactly the first term in the last equation.

Algorithm EM (Expectation-Maximization)

1. Initialize.
2. **E-step:** Fix $\theta = \theta^t$ in (3). Maximize the lower bound \mathcal{L} wrt. the distribution q . This is achieved by minimizing KL because

$$\mathcal{L}(q, \theta) = \log p(\mathbf{x}|\theta) - \text{KL}(q||p_{\mathbf{Z}|\mathbf{x}}),$$

and $\log p(\mathbf{x}|\theta)$ is independent of q . This maximum is realized for the choice $q^t = p_{\mathbf{Z}|\mathbf{x},\theta^t}$ because $\text{KL}(q||p_{\mathbf{Z}|\mathbf{x},\theta^t}) \geq 0$ with equality iff $q = p_{\mathbf{Z}|\mathbf{x},\theta^t}$, see (12). For this q^t , $\mathcal{L}(q^t, \theta^t) = \log p(\mathbf{x}|\theta^t)$ by construction, it is a lower bound and the gradients of both functions are parallel, so we have a situation like in Figure 2

Finally, to prepare for the next step compute explicitly \mathcal{L} as given by (4) and $g(\mathbf{z}) := \log p(\mathbf{x}, \mathbf{z}|\theta)$:

$$\mathcal{L}(p_{\mathbf{Z}|\mathbf{x},\theta^t}, \theta) = E_{\mathbf{Z}|\mathbf{x},\theta^t}[g(\mathbf{Z})] + C.\tag{5}$$

3. **M-step:** Fix now $q^t = p_{\mathbf{Z}|\mathbf{x},\theta^t}$ coming from the E-step and maximize $\mathcal{L}(q^t, \theta)$ in (5) wrt. θ to obtain θ^{t+1} . This new choice of parameters yields a new lower bound $\mathcal{L}(q^t, \theta^{t+1}) \geq \mathcal{L}(q^t, \theta^t)$ and a yet greater increase in the target function:

$$\begin{aligned}\log p(\mathbf{x}|\theta^{t+1}) &= \mathcal{L}(q^t, \theta^{t+1}) + \underbrace{\text{KL}(q^t||p_{\mathbf{Z}|\mathbf{x},\theta^{t+1}})}_{>0} \\ &> \mathcal{L}(q^t, \theta^t) + \underbrace{\text{KL}(q^t||p_{\mathbf{Z}|\mathbf{x},\theta^t})}_{=0} \\ &= \log p(\mathbf{x}|\theta^t),\end{aligned}$$

where the first KL divergence is strictly positive because $q^t \neq p_{\mathbf{Z}|\mathbf{x},\theta^{t+1}}$. This means that this step always increases the log likelihood unless we were already at a maximum and $\theta^t = \theta^{t+1}$.

4. Check whether we have finished.

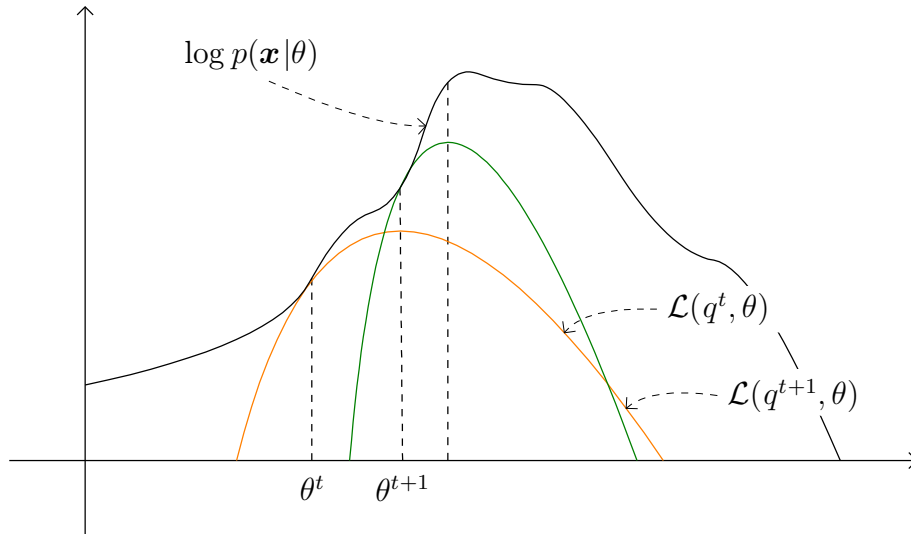


Figure 2. Example for a mixture in the exponential family, where the lower bound is always a concave function. Computing θ^{t+1} in the M-step increases the log likelihood. The E-step then computes the posterior q^{t+1} over the latent variables for the new parameters θ^{t+1} , yielding a new lower bound functional $\mathcal{L}(q^{t+1}, \theta)$, tangent to the log likelihood at θ^{t+1} .

5. BACK TO GAUSSIAN MIXTURES

Recall from §3 the expressions for p_Z and $p_{X|z}$. The joint density of (X, Z) is

$$p_{X,Z}(x, z) = \prod_{k=1}^K [\pi_k \mathcal{N}(x | \mu_k, \Sigma_k)]^{z_k}.$$

Suppose we had all the values $\mathbf{z} = (z_1, \dots, z_N)$ in addition to the data $\mathbf{x} = (x_1, \dots, x_N)$ (one talks of the **complete dataset**) and we wanted to maximize the log likelihood:

$$\begin{aligned} \log p(\mathbf{x}, \mathbf{z} | \theta) &= \log \prod_{n=1}^N \prod_{k=1}^K [\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)]^{z_{nk}} \\ &= \sum_{n=1}^N \sum_{k=1}^K z_{nk} [\log \pi_k + \log \mathcal{N}(x_n | \mu_k, \Sigma_k)] \end{aligned} \quad (6)$$

where we now had to assume $\pi_k \in (0, 1)$ to be able to take the logarithm. The problem is we don't have the $z_n!$ But EM saves the day:

E-step: Fix $\theta = \theta^t$. We know that the optimal q^t is the posterior over the latent variables, and we immediately see that it factorizes:

$$p_{\mathbf{Z}|\mathbf{x}, \theta^t}(\mathbf{z}) = \frac{p(\mathbf{x}|\mathbf{z}, \theta^t) p(\mathbf{z}|\theta^t)}{p(\mathbf{x}|\theta^t)} = \prod_{n=1}^N \underbrace{\frac{1}{p(x|\theta^t)} \prod_{k=1}^K [\pi_k^t \mathcal{N}(x | \mu_k^t, \Sigma_k^t)]^{z_k}}_{=: p_{Z|x}(z)}.$$

With this choice for q^t the new lower bound is the function of θ :

$$\mathcal{L}(p_{\mathbf{Z}|\mathbf{x}, \theta^t}, \theta) = E_{\mathbf{Z}|\mathbf{x}, \theta^t}[g(\mathbf{Z})] + C,$$

where $g(\boldsymbol{\zeta}) := \log p_{\mathbf{X}, \mathbf{Z}|\theta}(\mathbf{x}, \boldsymbol{\zeta}|\theta)$. Setting $z_{nk} = \mathbb{1}_{\{e_k\}}(Z_n)$, the expectation is

$$\begin{aligned} E_{\mathbf{Z}|\mathbf{x}, \theta^t}[g(\mathbf{Z})] &= E_{\mathbf{Z}|\mathbf{x}, \theta^t} \left[\sum_{n=1}^N \sum_{k=1}^K z_{nk} [\log \pi_k + \log \mathcal{N}(x_n | \mu_k, \Sigma_k)] \right] \\ &= \sum_{n=1}^N \sum_{k=1}^K \underbrace{E_{\mathbf{Z}|\mathbf{x}, \theta^t}[\mathbb{1}_{\{e_k\}}(Z_n)]}_{(*)} [\log \pi_k + \log \mathcal{N}(x_n | \mu_k, \Sigma_k)]. \end{aligned}$$

In order to compute $(*)$ we use that the posterior factorizes over n :¹

$$\begin{aligned} E_{\mathbf{Z}|\mathbf{x}, \theta^t}[\mathbb{1}_{\{e_k\}}(Z_n)] &= \sum_{\mathbf{z} \in \{e_1, \dots, e_K\}^N} \mathbb{1}_{\{e_k\}}(z_n) p_{\mathbf{Z}|\mathbf{x}, \theta^t}(\mathbf{z}) \\ &= \sum_{\mathbf{z} \in \{e_1, \dots, e_K\}^N} \mathbb{1}_{\{e_k\}}(z_n) \left(\prod_{j=1}^N p_{Z|x, \theta^t}(z_j) \right) \\ &= \underbrace{\sum_{z \in \{e_1, \dots, e_K\}} \mathbb{1}_{\{e_k\}}(z) p_{Z|x, \theta^t}(z)}_{=E_{Z|x, \theta^t}[\mathbb{1}_{\{e_k\}}(Z_n)]} \underbrace{\left(\prod_{j=1, j \neq n}^N \sum_{z \in \{e_1, \dots, e_K\}} p_{Z|x, \theta^t}(z) \right)}_{=1} \\ &= p_{Z|x, \theta^t}(e_k) \\ &= \frac{1}{p(x|\theta^t)} \pi_k^t \mathcal{N}(x | \mu_k^t, \Sigma_k^t) \\ &= \frac{\pi_k^t \mathcal{N}(x | \mu_k^t, \Sigma_k^t)}{\sum_{j=1}^K \pi_j^t \mathcal{N}(x | \mu_j^t, \Sigma_j^t)} =: \gamma_{nk}^t. \end{aligned}$$

Notice that this quantity is precisely the **responsibility** of cluster k for the point x_n as defined in (10) when maximizing the log likelihood of the marginal $p(\mathbf{x}|\theta)$. Substituting above, we have:

$$E_{\mathbf{Z}|\mathbf{x}}[\log p(\mathbf{x}, \mathbf{z}|\theta^t)] = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk}^t [\log \pi_k^t + \log \mathcal{N}(x_n | \mu_k^t, \Sigma_k^t)], \quad (7)$$

and this is the quantity we maximize next.

M-step: We maximize (7) wrt. each of μ_k, Σ_k, π_k obtaining *closed formulas* for μ_k^{t+1} and we are in business for the next step.

1. Differentiate (7) wrt. μ_k and equate to zero to obtain: $0 = \sum_{n=1}^N \gamma_{nk}^t (\Sigma_k^t)^{-1} (x_n - \mu_k)$. Now solve for μ_k and multiply on the left by Σ_k^t :

$$\mu_k^{t+1} = \frac{1}{N_k^t} \sum_{n=1}^N \gamma_{nk}^t x_n, \quad \text{where } N_k^t = \sum_{n=1}^N \gamma_{nk}^t. \quad (8)$$

1. This is just an application of ‘‘Fubini’s’’ theorem, and might be easier to see if we write integrals to make notation easier.

2. Differentiate (7) wrt. Σ_k and equate to zero. After some computations and solving for Σ_k (here the constraint that Σ_k be symmetric positive definite complicates matters considerably, see e.g. [AO85]) we find:

$$\Sigma_k^{t+1} = \frac{1}{N_k^t} \sum_{n=1}^N \gamma_{nk}^t (x_n - \mu_k^t) \otimes (x_n - \mu_k^t).$$

3. Differentiate (7) wrt. π_k , with the constraint $\sum_{k=1}^K \pi_k = 1$ using a Lagrange multiplier (recall that now $\pi_k \in (0, 1)$):

$$0 = \partial_{\pi_k} E_{\mathbf{Z}|\mathbf{x}}[\log p(\mathbf{x}, \mathbf{z}|\theta^t)] + \partial_{\pi_k} \lambda \left(\sum_k \pi_k - 1 \right) = \sum_{n=1}^N \frac{\gamma_{nk}^t}{\pi_k} + \lambda = \frac{N_k^t}{\pi_k} + \lambda.$$

Multiply by π_k and sum over K using the constraint on the π_k and the fact that $\sum_{k=1}^K \gamma_{nk}^t = 1$:

$$0 = \sum_{k=1}^K N_k^t + \lambda \sum_{k=1}^K \pi_k \Rightarrow \lambda = -N$$

Substituting above this yields

$$\pi_k^{t+1} = \frac{N_k^t}{N}.$$

Evaluate: We use the estimators obtained in the M-step to evaluate the log likelihood and check if we should stop.

6. AN EXAMPLE WITH MIXTURES OF BERNOULLIS

This model is called *latent class analysis* and is an exercise for the reader! We will use it for Hidden Markov Models over discrete variables, if we ever do those.

We now have D i.i.d. binary variables $B = (B_1, \dots, B_D)$, with $B_i \sim \mathcal{B}(\mu_i)$, i.e.

$$p(b|\mu) = \prod_{i=1}^D \mu_i^{b_i} (1 - \mu_i)^{1-b_i}.$$

Consider a mixture

$$X \sim \sum_{k=1}^K \pi_k B(\mu^k).$$

Compute mean, covariance and log likelihood for some data $\mathbf{x} = (x_1, \dots, x_N)$.

Apply this to the MNIST data set of handwritten digits, see the Kaggle at [Kag12] or the original source at [LCB12].

7. WHERE TO GO FROM HERE

Some ideas:

- Implement **maximum a posteriori** estimation using EM. The idea is to maximize the log likelihood of the posterior over the parameters $\log p(\theta|\mathbf{x})$ to solve the problem of “collapsing variances”. Using Bayes’ rule and the decomposition (3) we have

$$\begin{aligned} \log p(\theta|\mathbf{x}) &= \log p(\theta, \mathbf{x}) - \log p(\mathbf{x}) \\ &= \log p(\mathbf{x}|\theta) + \log p(\theta) - \underbrace{\log p(\mathbf{x})}_{=C} \\ &= \mathcal{L}(q, \theta) + \text{KL}(q\|p_{\mathbf{Z}|\mathbf{x},\theta}) + \log p(\theta) - C \\ &= \tilde{\mathcal{L}}(q, \theta) + \text{KL}(q\|p_{\mathbf{Z}|\mathbf{x},\theta}), \end{aligned}$$

and we now apply EM to this expression. The only necessary modification is to the estimators in the M step, due to the appearance of the prior in $\log p(\theta)$.

- Fully Bayesian use of EM including model selection.
- Expectation Propagation.

8. APPENDIX

8.1. A word about conditional densities.

If (X, Z) is a RV with real values, then

$$p(z|x) = p_{Z|X=x}(z) := \frac{p_{X,Z}(x, z)}{p_X(x)}$$

defines a probability density at every x such that $p(x) \neq 0$, since it is non-negative, measurable and

$$\int p(z|x) dz = \frac{1}{p(x)} \int p(x, z) dz = 1.$$

For discrete variables this agrees with the usual product rule, but for continuous variables needs interpreting, since $\{X = x\}$ is a null set, see e.g. [JP04, §12].

8.2. Problems related to the log likelihood approach.

Collapsing variance: Suppose that at some point in our algorithm for maximizing (2), whatever it is, we set $\mu_k = x_n$ for some n, k , then the k -th summand in (2) will be

$$\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) = C \frac{1}{\sqrt{\det \Sigma_k}},$$

and we can obviously send this quantity to $+\infty$ by adjusting the covariance matrix. If the algorithm goes on trying to maximize the whole expression, it will do exactly that and will almost surely fail with a division by zero.

We talk about collapsing variance because Σ being symmetric positive definite, it may be diagonalized and the new diagonal entries are the variances of the X_i in the new coordinates.

We have at least two solutions to this problem, the hacky and the not-so-hacky:

1. Hack hack: use some heuristics in the code to avoid the situation $\mu_k = x_n$ for any k, n . For instance, randomly reset μ_k if there is any $n \in \{1, \dots, N\}$ such that $|x_n - \mu_k| < \varepsilon$ for some fixed $\varepsilon > 0$. Ugly!
2. Penalize: postulate a prior on the parameters p_Θ and add a term $\log p_\Theta(\theta)$ (called **capacity**) to penalize lower probabilities of θ . Some good choice of p_Θ should fix the issue. We can justify this idea within the framework presented in §4, as sketched in §7.

Non-closed estimators: Maximizing (2) is a matter of differentiating wrt. each parameter in θ , setting equal to 0 and solving for the parameter, then checking the second derivatives. Consider for example the mean μ_k :

$$\partial_{\mu_k} \log p(\mathbf{x}|\theta) = \dots \text{compute compute} \dots = \sum_{n=1}^N \gamma_{nk} \Sigma_k^{-1} (x_n - \mu_k).$$

Solving for μ_k yields:

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} x_n, \quad \text{where } N_k = \sum_{n=1}^N \gamma_{nk}, \quad (9)$$

and we defined the **responsibility**

$$\gamma_{nk} = P(Z_n = e_k | X_n = x_n) \quad (10)$$

to be the *posterior probability* of having z_n given the data x_n . You can see that (9) is not a closed formula for μ_k since on the rhs. γ_{nk} depends itself on μ_k . Indeed, recalling our previous computations for the joint and marginal and using the product rule twice we have explicitly

$$p(z|x) = \frac{p(x, z)}{p(x)} = \frac{p(x|z) p(z)}{p(x)} = \frac{\prod_{k=1}^K \mathcal{N}(x|\mu_k, \Sigma_k)^{z_k} \prod_{k=1}^K \pi_k^{z_k}}{\sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \Sigma_k)},$$

and consequently:

$$\gamma_{nk} = p(z_n = e_k | x_n) = \frac{\mathcal{N}(x_n|\mu_k, \Sigma_k) \pi_k}{\sum_{k=1}^K \pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)}.$$

Analogously, the other two estimators are also in non-closed form (they are obtained in the same way as in §5):

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk}(x_n - \mu_k) \otimes (x_n - \mu_k) \text{ and } \pi_k = \frac{N_k}{N}. \quad (11)$$

We already saw a solution with Philipp to this problem: the expressions (9), (11) for the estimators for μ_k, Σ_k, π_k feature the variables themselves at the rhs., which suggests that we use a Newton-like iterative scheme. Our different modelling guides us to Expectation-Maximization.

Model order: This is actually not a problem of the log likelihood estimation *per se*, but of the mixture model. We have always assumed that we know some value for the number of components K , but for many applications this will not be the case.

8.3. A few sketchy ideas from information theory.

We want to know how much information about a random variable X is transmitted when we send a specific value x . This should be a monotonically decreasing function of the probability of the event because of the following simple ideas:

- News of event with probability one are no information.
- News of event with low probability are “high” information.

Information content: $h(x) := -\log p(x)$.

Entropy of a RV:

$$E_X[h(X)] = - \int_{\mathbb{R}} \log(p(x)) p(x) dx,$$

where we set $p \log p = 0$ if $p = 0$ by continuity $p \ln p \rightarrow 0$ for $p \rightarrow 0$. This integral converges (!).

Noiseless decoding theorem. (Shanon 1948) *Entropy is a lower bound on the number of bits (resp. nats) needed to transmit the state of a random variable.*

Kullback-Leibler divergence: We have approximated an unknown distribution p of some random variable X (assumed to be $p > 0$) with a distribution q . How well did we do? Define the KL-*divergence* as

$$\begin{aligned} \text{KL}(p||q) &:= - \int \log q(x) p(x) dx - \underbrace{\left(- \int \log p(x) p(x) dx \right)}_{=H[X]} \\ &= - \int \log \frac{q(x)}{p(x)} p(x) dx. \end{aligned}$$

Intuition: KL is the *additional* amount of information needing to be transmitted to specify X by using q instead of p . $\text{KL}(\cdot\|\cdot)$ is not a distance since it is obviously not symmetric. We have however:

Non-negativity of KL.

$$\text{KL}(p\|q) \geq 0 \text{ with equality iff } p = q. \quad (12)$$

Proof. This is a direct application of Jensen's inequality (recall that for φ convex and μ finite this is: $\varphi(\int f d\mu) \leq \int \varphi(f) d\mu$, when things make sense, blah blah...). By the convexity of the $-\log$ (recall that we assume $p > 0$):

$$\text{KL}(p\|q) = \int -\log \frac{q(x)}{p(x)} p(x) dx \geq -\log \int \frac{q(x)}{p(x)} p(x) dx = -\log \int q(x) dx = 0,$$

where we used that q is a probability distribution. □

An example using KL: Suppose we stick to one family of distributions $q = q(x|\theta)$. We would like to estimate the best θ by minimizing KL:

$$\text{KL}(p\|q) = - \int \log \frac{q(x)}{p(x)} p(x) dx = -E_p[\log q] + C,$$

but this is not possible since p is unknown. However, the expectation wrt. p may be approximated (!) by a finite sum over the observed data x_1, \dots, x_N :

$$\begin{aligned} \text{KL}(p\|q) &\simeq \frac{1}{N} \sum_{n=1}^N [-\log q(x_n|\theta) + \log p(x_n)] \\ &= \frac{1}{N} \sum_{n=1}^N -\log q(x_n|\theta) + C, \end{aligned}$$

and this quantity may be now minimized. Notice that the first term is the negative log likelihood of $q(\cdot|\theta)$ under the observed data $\mathbf{x} = (x_1, \dots, x_n)$.

8.4. Maximization of the log likelihood for the complete dataset.

Assuming we have the complete dataset, i.e. all of \mathbf{x} and \mathbf{z} (which we don't) finding maxima wrt. to the parameters is easy. For the means we compute

$$\partial_{\mu_k} \log p(\mathbf{x}, \mathbf{z}|\theta) = \sum_{n \in C_k} \partial_{\mu_k} \log \mathcal{N}(x_n|\mu_k, \Sigma_k) = \sum_{n \in C_k} \Sigma_k^{-1} (x_n - \mu_k).$$

Equating to zero and solving for μ_k we find that

$$\bar{\mu}_k = \frac{1}{\#C_k} \sum_{n \in C_k} x_n$$

is a critical point, and it is a maximum by the properties of Σ_k . Analogously, differentiation wrt. Σ_k yields (with quite some effort since we have the constraint that Σ_k be symmetric positive definite, see e.g. [AO85])

$$\bar{\Sigma}_k = \frac{1}{\#C_k} \sum_{n \in C_k} (x_n - \bar{\mu}_k) \otimes (x_n - \bar{\mu}_k).$$

Finally, for the weights π_k we have to include the constraint $\sum_k \pi_k = 1$, which we do with a Lagrange multiplier (recall that now $\pi_k \in (0, 1)$, so this is no issue):

$$0 = \partial_{\pi_k} \log p(\mathbf{x}, \mathbf{z} | \theta) + \partial_{\pi_k} \lambda \left(\sum_k \pi_k - 1 \right) = \sum_{n \in C_k} \frac{1}{\pi_k} + \lambda = \frac{\#C_k}{\pi_k} + \lambda.$$

We multiply by π_k and sum over K using the constraint on π_k : $\lambda \pi_k = -\#C_k \Rightarrow \lambda \sum_{k=1}^K \pi_k = -\sum_{k=1}^K \#C_k = -N$, so we have $\lambda = -N$. Substituting above this yields

$$\pi_k = \frac{\#C_k}{N}.$$

9. REFERENCES

You can read a much better version of all this stuff in [Bis06, Chapter 9]. Another reference, explaining less but with more applications is [Mur12, Chapter 11]. For more on probabilistic graphical models, read [Mur12, Chapter 10] or [KF09, Chapter 3]. A $\text{T}_{\text{E}}\text{X}_{\text{MACS}}$ file with an implementation of EM for Gaussians is also available.

-
- [AO85] T. W. Anderson and I. Olkin. Maximum-likelihood estimation of the parameters of a multivariate normal distribution. *Linear Algebra and its Applications*, 70:147–171, oct 1985.
 - [Bis06] Christopher M. Bishop. *Pattern recognition and machine learning*. Information science and statistics. Springer, 1 edition, aug 2006.
 - [JP04] Jean Jacod and Philip Protter. *Probability essentials*. Universitext. Springer Berlin Heidelberg, 2 edition, jan 2004.
 - [Kag12] Kaggle.com. Classify handwritten digits using the famous MNIST data. jul 2012.
 - [KF09] Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. Adaptive Computation and Machine Learning. MIT Press, 2009.
 - [LCB12] Yan LeCun, Corinna Cortes and Christopher J.C. Burges. The MNIST database of handwritten digits. 2012.
 - [Mur12] Kevin P. Murphy. *Machine learning: a probabilistic perspective*. MIT Press, aug 2012.