

Expectation Maximization

Benjamin Bray

Monday, November 14, 2016*

Abstract

These notes provide a theoretical treatment of **Expectation-Maximization**, an iterative parameter estimation algorithm used to find local maxima of the likelihood function in the presence of hidden variables. Introductory textbooks (Murphy, 2012; Bishop, 2006) typically state the algorithm without explanation and expect students to work blindly through derivations. We find this approach to be unsatisfying, and instead choose to tackle the theory head-on, followed by plenty of examples. Following (Neal und Hinton, 1998), we view expectation-maximization as coordinate ascent on the **Evidence Lower Bound**. This perspective takes much of the “mystery” out of the algorithm and allows us to easily derive variants like **Hard EM** and **Variational EM**.

Note: I actually wrote these notes for myself when learning about expectation maximization for the first time. Accordingly, they have only been lightly proofread and some parts are incomplete. Send an email to benrbray@umich.edu if you spot any mistakes!

Part 1

Introduction

Expectation-maximization is an iterative algorithm for performing maximum likelihood estimation in latent variable models. While special cases of the algorithm have existed for decades, the algorithm was not appreciated in its full generality until (Dempster u. a., 1977) and later (Neal und Hinton, 1998). Instances of Expectation-Maximization have since been derived for common statistical models, most commonly Hidden Markov Models and Gaussian Mixtures, and presented as canned algorithms to be memorized, not understood. Sadly, many students and even researchers are unaware of the rich mathematical theory that lends Expectation-Maximization its robustness.

1.1 Probabilistic Modeling and Maximum Likelihood

A *probabilistic model* P is a set X of random variables along with parameters θ which describe how the data is generated. Usually, P is represented by a graphical model or list

**Last modified:* Monday, November 14, 2016

of conditional distributions for each variable. The model is often interpreted as a *recipe* for generating new data from scratch. In this paper, assume we make N i.i.d. observations $\mathcal{X} = (X_1, \dots, X_N)$ of the random variables X . The goal of **maximum likelihood estimation** is to compute a point estimate of the parameters θ that maximizes the probability of the observed data:

$$\hat{\theta}_{ML} = \arg \max_{\theta} \ell(\theta | \mathcal{X}) = \arg \max_{\theta} \log p(\mathcal{X}|\theta) \quad (1)$$

The likelihood function is typically assumed to be an exponential family distribution, making this optimization problem convex. Therefore, generic gradient-based optimization techniques like gradient descent, Newton’s method, conjugate gradient, and others are guaranteed to find a single global maximum to the likelihood.

1.2 Latent Variable Models

Probabilistic inference works best when the observed variables are *independent* or *uncorrelated*. In applications, this is rarely the case¹, and many models make independence assumptions that do not necessarily reflect reality. When working with graphical models, this is equivalent to deleting edges in the graph. Although these additional assumptions make inference tractable, they weaken our model and could lead to poor accuracy.

An alternative approach is to assume that observed variables are correlated because they are influenced by **latent variables** which are not observed but nonetheless impact observations. Latent variable models are capable of modeling complex interdependencies between observed variables without a huge increase in computational cost. Suppose we observe data \mathcal{X} generated from a model p with true parameters θ^* in the presence of hidden variables Z . The goal of maximum likelihood estimation is the same as before:

$$\hat{\theta}_{ML} = \arg \max_{\theta} \ell(\theta|\mathcal{X}) = \arg \max_{\theta} \log p(\mathcal{X}|\theta)$$

In some cases, we also seek to *infer* the values \mathcal{Z} of the hidden variables Z . In the Bayesian spirit, we will treat the parameter θ^* as a realization of some random variable Θ . The **observed data log-likelihood** $\ell(\theta|\mathcal{X}) = \log p(\mathcal{X}|\theta)$ of the parameters given the observed data is useful for both inference and parameter estimation, in which we must grapple with uncertainty about the hidden variables. Working directly with this quantity is often difficult in latent variable models because the inner sum cannot be brought out of the logarithm when we marginalize over the latent variables:

$$\ell(\theta|\mathcal{X}) = \log p(\mathcal{X}|\theta) = \log \sum_z p(\mathcal{X}, z|\theta) \quad (2)$$

In general, this likelihood is non-convex with many local maxima. Finding the global optimum can be NP-Hard even for simple models like Gaussian mixtures (Aloise u. a., 2009) for which the number of local optima is exponential in the cluster count. For this reason, the gradient-based methods mentioned in the previous section for complete-data maximum likelihood estimation fail in the latent variable case. Further, when using these algorithms, much care must be taken to constrain model parameters to valid ranges, as specified by the model.

¹For example, the weather today in Berlin is loosely dependent on the weather last week in Detroit, even though this dependence is small and nearly impossible to model directly.

Part 2

Expectation Maximization

The **Expectation-Maximization** algorithm exploits the fact that learning is easy when the values of all variables are known. It is shown in (Murphy, 2012) that when $p(x_n, z_n|\theta)$ are exponential family distributions, the likelihood is convex, and gradient-based methods are again feasible. Therefore, it makes sense to alternate between the following two steps:

- (1) Infer the values of the latent variables.
- (2) Re-estimate the model parameters, assuming we have complete data.

2.1 Evidence Lower Bound

Our general approach will be to reason about the hidden variables through a proxy distribution q , which we use to compute a lower-bound on the log-likelihood. This section is devoted to deriving one such bound, called the **Evidence Lower Bound (ELBO)**. We can expand the data log-likelihood by marginalizing over the hidden variables:

$$\ell(\theta|\mathcal{X}) = \log p(\mathcal{X}|\theta) = \log \sum_z p(\mathcal{X}, z|\theta) \quad (3)$$

Through Jensen's inequality, we obtain the following bound, valid for any distribution q :

$$\ell(\theta|\mathcal{X}) = \log \sum_z p(\mathcal{X}, z|\theta) \quad (4)$$

$$= \log \sum_z q(z) \frac{p(\mathcal{X}, z|\theta)}{q(z)} \quad (5)$$

$$\geq \sum_z q(z) \log \frac{p(\mathcal{X}, z|\theta)}{q(z)} \equiv \mathcal{L}(q, \theta) \quad (6)$$

The lower bound $\mathcal{L}(q, \theta)$ can be rewritten as follows:

$$\ell(\theta|\mathcal{X}) \geq \mathcal{L}(q, \theta) = \sum_z q(z) \log \frac{p(\mathcal{X}, z|\theta)}{q(z)} \quad (7)$$

$$= \sum_z q(z) \log p(\mathcal{X}, z|\theta) - \sum_z q(z) \log q(z) \quad (8)$$

$$= E_q[\log p(\mathcal{X}, Z|\theta)] - E_q[\log q(z)] \quad (9)$$

$$= E_q[\log p(\mathcal{X}, Z|\theta)] + H(q) \quad (10)$$

Relationship to Relative Entropy

The first term in the last line above closely resembles the cross entropy between $q(Z)$ and the joint distribution $p(X, Z)$ of the observed and hidden variables. However, the variables X are fixed to our observations $X = \mathcal{X}$ and so $p(\mathcal{X}, Z)$ is an *unnormalized*² distribution

²In this case, $\int p(\mathcal{X}, z) dz \neq 1$.

over Z . It is easy to see that this does not set us back too far; in fact, the lower bound $\mathcal{L}(q, \theta)$ differs from a Kullback-Liebler divergence only by a constant with respect to Z :

$$D_{KL}(q||p(Z|\mathcal{X}, \theta)) = H(q, p(Z|\mathcal{X}, \theta)) - H(q) \quad (11)$$

$$= E_q[-\log p(Z|\mathcal{X}, \theta)] - H(q) \quad (12)$$

$$= E_q[-\log p(\mathcal{X}, Z|\theta)] - E_q[-\log p(\mathcal{X}|\theta)] - H(q) \quad (13)$$

$$= E_q[-\log p(\mathcal{X}, Z|\theta)] + \log p(\mathcal{X}|\theta) - H(q) \quad (14)$$

$$= -\mathcal{L}(q, \theta) + \text{const.} \quad (15)$$

This yields a second proof of the evidence lower bound, following from the nonnegativity of relative entropy. In fact, this is the proof given in Tzikas u. a. (2008) and Murphy (2012).

$$\log p(\mathcal{X}|\theta) = D_{KL}(q||p(Z|\mathcal{X}, \theta)) + \mathcal{L}(q, \theta) \geq \mathcal{L}(q, \theta) \quad (16)$$

Selecting a Proxy Distribution

The quality of our lower bound $\mathcal{L}(q, \theta)$ depends heavily on the choice of proxy distribution $q(Z)$. We now show that the evidence lower bound is *tight* in the sense that equality holds when the proxy distribution $q(Z)$ is chosen to be the hidden posterior $p(Z|\mathcal{X}, \theta)$. This will be useful later for proving that the Expectation Maximization algorithm converges.

Remark 1. *Maximizing $\mathcal{L}(q, \theta)$ with respect to q is equivalent to minimizing the relative entropy between q and the hidden posterior $p(Z|\mathcal{X}, \theta)$. Hence, the optimal choice for q is exactly the hidden posterior, for which $D_{KL}(q||p(Z|\mathcal{X}, \theta)) = 0$, and*

$$\log p(\mathcal{X}|\theta) = E_q[\log p(\mathcal{X}, Z|\theta)] + H(q) = \mathcal{L}(q, \theta)$$

2.2 Expectation Maximization

Recall that the maximum likelihood estimate of the parameters θ given observed data \mathcal{X} in the presence of hidden variables Z is

$$\hat{\theta}_{ML} = \arg \max_{\theta} \ell(\theta|\mathcal{X}) = \arg \max_{\theta} \log p(\mathcal{X}|\theta) \quad (17)$$

Unfortunately, when reasoning about hidden variables, finding a global maximum is difficult. Instead, the **Expectation Maximization** algorithm is an iterative procedure for computing a local maximum of the likelihood function, under the assumption that the hidden posterior $p(Z|\mathcal{X}, \theta)$ is tractable. We will take advantage of the evidence lower bound

$$\ell(\theta|\mathcal{X}) \geq \mathcal{L}(q, \theta) \quad (18)$$

on the data likelihood. Consider only proxy distributions of the form $q_{\vartheta}(Z) = p(Z|\mathcal{X}, \vartheta)$, where ϑ is some fixed configuration of the variables Θ , possibly different from our estimate θ . The optimal value for ϑ , in the sense that $\mathcal{L}(q_{\vartheta}, \theta)$ is maximum, depends on the particular choice of θ . Similarly, the optimal value for θ depends on the choice of ϑ . This suggests an iterative scheme in which we alternate between maximizing with respect to ϑ and with respect to θ , gradually improving the log-likelihood.

Iterative Procedure

Suppose at time t we have an estimate θ_t of the parameters. To improve our estimate, we perform two steps of coordinate ascent on $\mathcal{L}(\vartheta, \theta) \equiv \mathcal{L}(q_\vartheta, \theta)$, as described in Neal und Hinton (1998),

E-Step Compute a new lower bound on the observed log-likelihood, with

$$\vartheta_{t+1} = \arg \max_{\vartheta} \mathcal{L}(\vartheta, \theta_t) = \theta_t$$

M-Step Estimate new parameters by optimizing over the lower bound,

$$\theta_{t+1} = \arg \max_{\theta} \mathcal{L}(\vartheta_{t+1}, \theta) = \arg \max_{\theta} E_q[\log p(\mathcal{X}, Z|\theta)]$$

In the M-Step, the expectation is taken with respect to $q_{\vartheta_{t+1}}$.

Alternative Formulation

In the M-Step, the entropy term of the evidence lower bound $\mathcal{L}(\vartheta_{t+1}, \theta)$ does not depend on θ . The remaining term $Q(\theta_t, \theta) = E_q[\log p(\mathcal{X}, Z|\theta)]$ is sometimes called the **auxiliary function** or **Q-function**. To us, this is the **expected complete-data log-likelihood**.

Proof of Convergence

To prove convergence of this algorithm, we show that the data likelihood $\ell(\theta|\mathcal{X})$ increases after each update.

Theorem 1. *After a single iteration of Expectation Maximization, the observed data likelihood of the estimated parameters has not decreased, that is,*

$$\ell(\theta_t|\mathcal{X}) \leq \ell(\theta_{t+1}|\mathcal{X})$$

Proof. This result is a simple consequence of all the hard work we have put in so far:

$$\begin{aligned} \ell(\theta_t|\mathcal{X}) &= \mathcal{L}(q_{\vartheta_{t+1}}, \theta_t) && \text{(Remark 1)} \\ &\leq \mathcal{L}(q_{\vartheta_{t+1}}, \theta_{t+1}) && \text{(M-Step)} \\ &\leq \ell(\theta_{t+1}|\mathcal{X}) && \text{(ELBO)} \end{aligned}$$

□

It is also possible to show that Expectation-Maximization converges to something *useful*.

Theorem 2. *(Neal & Hinton 1998, Thm. 2) Every local maximum of the evidence lower bound $\mathcal{L}(q, \theta)$ is a local maximum of the data likelihood $\ell(\theta|\mathcal{X})$.*

Starting from an initial guess θ_0 , We run this procedure until some stopping criterion is met and obtain a sequence $\{(\vartheta_t, \theta_t)\}_{t=1}^T$ of parameter estimates.

2.3 Example: Coin Flips

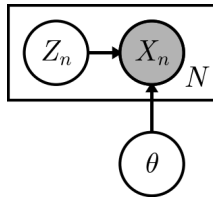
Now that we have a good grasp on the theory behind Expectation Maximization, let's get some intuition by means of a simple example. As usual, the simplest possible example involves coin flips!

Probabilistic Model

Suppose we have two coins, each with a different probability of heads, θ_A and θ_B , unknown to us. We collect data from a series of N trials in order to estimate the bias of each coin. Each trial k consists of flipping the same random coin Z_k a total of M times and recording only the total number X_k of heads.

This situation is best described by the following **generative probabilistic model**, which precisely describes our assumptions about how the data was generated. The corresponding graphical model and a set of sample data are shown in Figure 1.

$$\begin{aligned}
 \theta &= (\theta_A, \theta_B) && \text{fixed coin biases} \\
 Z_n &\sim \text{Uniform}\{A, B\} \quad \forall n = 1, \dots, N && \text{coin indicators} \\
 X_n | Z_n, \theta &\sim \text{Bin}[\theta_{Z_n}, M] \quad \forall n = 1, \dots, N && \text{head count}
 \end{aligned} \tag{19}$$



#	Sequence	Heads
1	HTTTHHTHTH	5
2	HHHHTHHHHH	9
3	HTHHHHHTHH	8
4	HTHTTTHHTT	4
5	THHHTHHHTH	7

Figure 1: Sample data and graphical model representation for the coin flip example, with $N = 5$ trials and $M = 10$ flips per trial. Adapted from Do and Batzoglou (2008).

Complete Data Log-Likelihood

The complete data log-likelihood for a single trial (x_n, z_n) is

$$\log p(x_n, z_n | \theta) = \log p(z_n) + \log p(x_n | z_n, \theta) \tag{20}$$

In this model, $P(z_n) = \frac{1}{2}$ is uniform. The remaining term is

$$\log p(x_n | z_n, \theta) = \log \binom{M}{x_n} \theta_{z_n}^{x_n} (1 - \theta_{z_n})^{M - x_n} \tag{21}$$

$$= \log \binom{M}{x_n} + x_n \log \theta_{z_n} + (M - x_n) \log(1 - \theta_{z_n}) \tag{22}$$

Expectation Maximization

Now that we have specified the probabilistic model and worked out all relevant probabilities, we are ready to derive an Expectation Maximization algorithm.

The **E-Step** is straightforward. The **M-Step** computes a new parameter estimate θ_{t+1} by optimizing over the lower bound found in the E-Step. Let $\vartheta = \vartheta_{t+1} = \theta_t$. Then,

$$\theta_{t+1} = \arg \max_{\theta} \mathcal{L}(\theta, q_{\vartheta}) = \arg \max_{\theta} E_q[\log p(\mathcal{X}, Z|\theta)] \quad (23)$$

$$= \arg \max_{\theta} E_q[\log p(\mathcal{X}|Z, \theta)p(Z)] \quad (24)$$

$$= \arg \max_{\theta} E_q[\log p(\mathcal{X}|Z, \theta)] + \log p(Z) \quad (25)$$

$$= \arg \max_{\theta} E_q[\log p(\mathcal{X}|Z, \theta)] \quad (26)$$

Now, because each trial is conditionally independent of the others, given the parameters,

$$E_q[\log p(\mathcal{X}|Z, \theta)] = E_q \left[\log \prod_{n=1}^N p(x_n|Z_n, \theta) \right] = \sum_{n=1}^N E_q[\log p(x_n|Z_n, \theta)] \quad (27)$$

$$= \sum_{n=1}^N E_q \left[x_n \log \theta_{z_n} + (M - x_n) \log(1 - \theta_{z_n}) \right] + \sum_{n=1}^N \log \binom{M}{x_n} \quad (28)$$

$$= \sum_{n=1}^N E_q \left[x_n \log \theta_{z_n} + (M - x_n) \log(1 - \theta_{z_n}) \right] + \text{const. w.r.t. } \theta \quad (29)$$

$$= \sum_{n=1}^N q_{\vartheta}(z_n = A) \left[x_n \log \theta_A + (M - x_n) \log \theta_A \right] \quad (30)$$

$$+ \sum_{n=1}^N q_{\vartheta}(z_n = B) \left[x_n \log \theta_B + (M - x_n) \log \theta_B \right] + \text{const. w.r.t. } \theta \quad (31)$$

Let $a_k = q(z_k = A)$ and $b_k = q(z_k = B)$. Note $\sum_{k=1}^N a_k = \sum_{k=1}^N b_k = 1$. To maximize the above expression with respect to the parameters, we take derivatives with respect to θ_A and θ_B and set to zero:

$$\frac{\partial}{\partial \theta_A} \left[E_q[\log p(\mathcal{X}|Z, \theta)] \right] = \frac{1}{\theta_A} \sum_{n=1}^N a_n x_n + \frac{1}{1 - \theta_A} \sum_{n=1}^N a_n (M - x_n) = 0 \quad (32)$$

$$\frac{\partial}{\partial \theta_B} \left[E_q[\log p(\mathcal{X}|Z, \theta)] \right] = \frac{1}{\theta_B} \sum_{n=1}^N b_n x_n + \frac{1}{1 - \theta_B} \sum_{n=1}^N b_n (M - x_n) = 0 \quad (33)$$

$$(34)$$

Solving for θ_A and θ_B , we obtain

$$\theta_A = \frac{\sum_{n=1}^N a_n x_n}{\sum_{n=1}^N a_n M} \quad \theta_B = \frac{\sum_{n=1}^N b_n x_n}{\sum_{n=1}^N b_n M} \quad (35)$$

2.4 Example: Gaussian Mixture Model

Probabilistic Model

In a Gaussian Mixture Model, samples are drawn from a random *cluster*, each normally distributed with its own mean and variance. Our goal will be to estimate the following parameters:

$$\begin{aligned}\boldsymbol{\pi} &= (\pi_1, \dots, \pi_K) && \text{mixing weights} \\ \boldsymbol{\mu} &= (\mu_1, \dots, \mu_K) && \text{cluster centers} \\ \boldsymbol{\Sigma} &= (\Sigma_1, \dots, \Sigma_K) && \text{cluster variance}\end{aligned}\tag{36}$$

The full model specification is below. A graphical model is shown in Figure 2.

$$\begin{aligned}\boldsymbol{\theta} &= (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) && \text{model parameters} \\ z_n &\sim \text{Cat}[\boldsymbol{\pi}] && \text{cluster indicators} \\ x_n | z_n, \boldsymbol{\theta} &\sim \mathcal{N}(\boldsymbol{\mu}_{z_n}, \boldsymbol{\Sigma}_{z_n}) && \text{base distribution}\end{aligned}\tag{37}$$

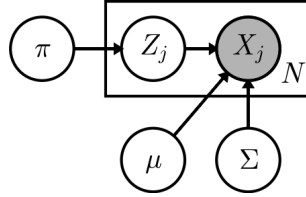


Figure 2: Gaussian Mixture Model.

Complete Data Log-Likelihood

The complete data log-likelihood for a single datapoint (x_n, z_n) is

$$\log p(x_n, z_n | \boldsymbol{\theta}) = \log \prod_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)^{\mathbb{I}(z_n=k)}\tag{38}$$

$$= \sum_{k=1}^K \mathbb{I}(z_n = k) \log \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)\tag{39}$$

Similarly, the complete data log-likelihood over all points $\{(x_n, z_n)\}_{n=1}^N$ is

$$\log p(X, Z | \boldsymbol{\theta}) = \sum_{n=1}^N \log p(x_n, z_n | \boldsymbol{\theta}) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{I}(z_n = k) \log \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)\tag{40}$$

Hidden Posterior

The hidden posterior for a single point (x_n, z_n) can be found using Bayes' rule:

$$p(z_n = k | x_n, \boldsymbol{\theta}) = \frac{P(z_n = k | \boldsymbol{\theta}) p(x_n | z_n = k, \boldsymbol{\theta})}{p(x_n | \boldsymbol{\theta})}\tag{41}$$

$$= \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{k'=1}^K \pi_{k'} \mathcal{N}(x_n | \mu_{k'}, \Sigma_{k'})}\tag{42}$$

Expectation Maximization

Our derivation will follow that of Murphy (2012), adapted to our notation.

E-Step

Before the E-step, we have an estimate θ_t of the parameters, and seek to compute a new lower bound on the observed log-likelihood. Earlier, we showed that the optimal lower bound is

$$\mathcal{L}(q_{\theta_t}, \theta) = E_q[\log p(\mathcal{X}, Z|\theta)] + \text{const.} \quad (43)$$

where $q_{\theta_t}(z) \equiv p(z|\mathcal{X}, \theta_t)$ and the second term is constant with respect to θ . The E-Step requires us to derive an expression for the first term. Using Equation 40, the expected complete data log-likelihood is given by

$$Q(\theta_t, \theta) = E_q[\log p(\mathcal{X}, Z|\theta)] = \sum_{n=1}^N \sum_{k=1}^K E_q[\mathbb{I}(z_n = k) \log \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)] \quad (44)$$

$$= \sum_{n=1}^N \sum_{k=1}^K E_q[\mathbb{I}(z_n = k)] \log \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \quad (45)$$

$$= \sum_{n=1}^N \sum_{k=1}^K p(z_n = k | x_n, \theta_t) \log \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \quad (46)$$

$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log \pi_k + \sum_{n=1}^N \sum_{k=1}^K r_{nk} \log \mathcal{N}(x_n | \mu_k, \Sigma_k) \quad (47)$$

where $r_{nk} \equiv p(z_n = k | x_n, \theta_t)$ is the **responsibility** that cluster k takes for data point x_n after step t . During the E-Step, we compute these values explicitly with Equation 42.

M-Step

During the M-Step, we optimize our lower bound with respect to the parameters $\theta = (\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$. For the mixing weights $\boldsymbol{\pi}$, we use Lagrange multipliers to maximize the ELBO subject to the constraint $\sum_{k=1}^K \pi_k = 1$. The Lagrangian is

$$\Lambda(\boldsymbol{\pi}, \lambda) = Q(\theta_t, \theta) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (48)$$

Carrying out the optimization, we find that $\lambda = -N$. The correct update for the mixing weights is

$$\boxed{\pi_k = \frac{1}{N} \sum_{n=1}^N r_{nk} = \frac{r_k}{N}} \quad (49)$$

where $r_k \equiv \sum_{n=1}^N r_{nk}$ is the *effective* number of points assigned to cluster k . For the cluster centers $\boldsymbol{\mu}$ and variance $\boldsymbol{\Sigma}$, you should verify that the correct updates are

$$\boxed{\mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{r_k}} \quad \boxed{\Sigma_k = \frac{\sum_{n=1}^N r_{nk} x_n x_n^T}{r_k} - \mu_k \mu_k^T} \quad (50)$$

2.5 Advice for Deriving EM Algorithms

The previous two examples suggest a general approach for deriving a new algorithm.

1. **Specify the probabilistic model.** Identify the observed variables, hidden variables, and parameters. Draw the corresponding graphical model to help determine the underlying independence structure.
2. **Identify the complete-data likelihood** $P(X, Z|\theta)$. For exponential family models, the complete-data likelihood will be convex and easy to optimize. In other models, other work may be required.
3. **Identify the hidden posterior** $P(Z|X, \theta)$. If this distribution is not tractable, you may want to consider variational inference, which we will discuss later.
4. **Derive the E-Step.** Write down an expression for $E_q[\log p(\mathcal{X}|Z, \theta)]$.
5. **Derive the M-Step.** Try taking derivatives and setting to zero. If this doesn't work, you may need to resort to gradient-based methods or variational inference.

2.6 Hard Expectation-Maximization

We are now in a position to demonstrate the connection between Gaussian Mixture Models and the K-Means algorithm. As it turns out, K-Means can be interpreted as a Gaussian Mixture Model in which all covariance matrices are identical and for which we perform learning via an algorithm called **Hard EM** or **Classification Maximization**.

In our derivation of Expectation Maximization, we constructed a lower bound $\mathcal{L}(q_\vartheta, \theta)$ on the log-likelihood, where the distributions q_ϑ belong to some known parametric family $\mathcal{Q} = \{q_\vartheta \mid \vartheta \in \Theta\}$ of distributions. We showed that the choice $q_\vartheta(Z) \equiv p(Z \mid \mathcal{X}, \vartheta)$ is optimal in the sense that it leads to a *tight* bound that touches the objective when $\vartheta = \theta$. What happens if we choose a different parametric family to optimize over?

Delta Approximation

We will make use of the fact that the best approximation of an arbitrary distribution p by a point mass occurs when the point mass δ_{x_0} is centered at the mode x_0 of the distribution.

Proposition 1. *Let p be a probability distribution, and denote by $q_{x_0}(x) \equiv \delta_{x_0}(x)$ be the distribution with unit mass at x_0 . Then, the quantity $D_{KL}(q_{x_0}||p)$ is minimum when $x_0 = \arg \max_x p(x)$ is the mode of p .*

Proof. The result is obvious upon observing that

$$D_{KL}(q_{x_0}||p) = \int_{\mathcal{X}} q_{x_0}(x) \log \frac{q_{x_0}(x)}{p(x)} dx = -\log p(x_0) \quad (51)$$

□

Classification Maximization

We now examine the consequences of choosing $\mathcal{Q} = \{\delta_{z_0} \mid z_0 \in Z\}$ to be the parametric family of distributions we use to approximate $p(Z \mid \mathcal{X}, \theta)$ in the derivation of Expectation-Maximization. Recall that

$$\log p(\mathcal{X} \mid \theta) = \mathcal{L}(q_{z_0}, \theta) + D_{KL}(q_{z_0}||p(Z|\mathcal{X}, \theta)) \quad (52)$$

Just as before, we perform coordinate ascent on the $\mathcal{L}(q_{z_0}, \theta)$. Since this bound is no longer (necessarily) tight, we are not guaranteed to find a local maximum of the likelihood. However, we can still write down an algorithm.

E-Step Compute a new lower bound on the observed log-likelihood, with

$$z_{t+1} = \arg \max_z \mathcal{L}(q_z, \theta_t) = \arg \max_z p(z|\mathcal{X}, \theta)$$

M-Step Estimate new parameters by optimizing over the lower bound,

$$\begin{aligned} \theta_{t+1} &= \arg \max_{\theta} \mathcal{L}(\vartheta_{t+1}, \theta) \\ &= \arg \max_{\theta} E_q[\log p(\mathcal{X}, Z|\theta)] \\ &= \arg \max_{\theta} \log p(\mathcal{X}, z_0|\theta) \end{aligned}$$

Starting from an initial guess θ_0 of the parameters, we alternate between estimating the most likely assignment z_0 of the hidden variables Z and computing a maximum likelihood estimate of θ assuming we have complete data (\mathcal{X}, z_0) . In contrast to expectation maximization, which takes into account the probabilities of different possible assignments to the hidden variables, the classification maximization reasons only about the single most likely assignment. For example, in clustering, a data point may be equidistant from two cluster centers. In Expectation Maximization, each cluster will have equal responsibility for this point, but in classification maximization, a *hard* assignment will be made arbitrarily.

Example: Connecting Gaussian Mixtures & K-Means

Let's derive a classification maximization algorithm for Gaussian Mixture Model clustering. Assume all clusters have the same covariance matrix $\Sigma = I$ and that the mixing weights π are uniform. In the E-Step, for each data point x_n , we assign

$$\begin{aligned} z_n &= \arg \max_k p(z_n = k|x_n, \theta_t) \\ &= \arg \max_k \log p(z_n = k|x_n, \theta_t) \\ &= \arg \max_k \log \mathcal{N}(x_n|\mu_k, I) \\ &= \arg \max_k \|x - \mu_k\|^2 \end{aligned}$$

In the M-Step, we re-estimate the cluster means by averaging over all data points currently assigned to that cluster. Amazingly, this is exactly the **K-Means** algorithm!

References

- [Aloise u. a. 2009] ALOISE, Daniel ; DESHPANDE, Amit ; HANSEN, Pierre ; POPAT, Preyas: NP-hardness of Euclidean sum-of-squares clustering. In: *Machine learning* 75 (2009), Nr. 2, S. 245–248
- [Bishop 2006] BISHOP, Christopher M.: *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA : Springer-Verlag New York, Inc., 2006. – ISBN 0387310738
- [Dempster u. a. 1977] DEMPSTER, Arthur P. ; LAIRD, Nan M. ; RUBIN, Donald B.: Maximum likelihood from incomplete data via the EM algorithm. In: *Journal of the royal statistical society. Series B (methodological)* (1977), S. 1–38
- [Do und Batzoglu 2008] DO, Chuong B. ; BATZOGLOU, Serafim: What is the expectation maximization algorithm? In: *Nature biotechnology* 26 (2008), Nr. 8, S. 897–899
- [Murphy 2012] MURPHY, Kevin P.: *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012. – ISBN 0262018020, 9780262018029
- [Neal und Hinton 1998] NEAL, Radford M. ; HINTON, Geoffrey E.: A view of the EM algorithm that justifies incremental, sparse, and other variants. In: *Learning in graphical models*. Springer, 1998, S. 355–368
- [Tzikas u. a. 2008] TZIKAS, Dimitris G. ; LIKAS, Aristidis C. ; GALATSANOS, Nickolaos P.: The variational approximation for Bayesian inference. In: *Signal Processing Magazine, IEEE* 25 (2008), Nr. 6, S. 131–146