

(note: the notation is slightly different from lecture slides)

Mixture models

There are many problems in machine learning that are not simple classification problems but rather modeling problems (e.g., clustering, diagnosis, combining multiple information sources for sequence annotation, and so on). Moreover, even within classification problems, we often have unobserved variables that would make a difference in terms of classification. For example, if we are interested in classifying tissue samples into specific categories (e.g., tumor type), it would be useful to know the composition of the tissue sample in terms of cells that are present and in what proportions. While such variables are not typically observed, we can still model them and make use of their presence in prediction.

Mixture models are simple probability models that try to capture ambiguities in the available data. They are simple, widely used and useful. As the name suggests, a mixture model “mixes” different predictions on the probability scale. The mixing is based on alternative ways of *generating* the observed data. Let \underline{x} be a vector of observations. A mixture model over vectors \underline{x} is defined as follows. We assume each \underline{x} could be of m possible types. If we knew the type, j , we would model \underline{x} with a conditional distribution $P(\underline{x}|j)$ (e.g., Gaussian with a specific mean). If the overall frequency of type j in the data is $P(j)$, then the “mixture distribution” over \underline{x} is given by

$$P(\underline{x}) = \sum_{j=1}^m P(\underline{x}|j)P(j) \quad (1)$$

In other words, \underline{x} could be generated in m possible ways. We imagine the generative process to be as follows: sample j from the frequencies $P(j)$, then \underline{x} from the corresponding conditional distribution $P(\underline{x}|j)$. Since we do not observe the particular way the example was generated (assuming the model is correct), we sum over the possibilities, weighted by the overall frequencies.

We have already encountered mixture models. Take, for example, the Naive Bayes model $P(\underline{x}|y)P(y)$ over the feature vector \underline{x} and label y . If we pool together examples labeled $+1$ and those labeled -1 , and throw away the label information, then the Naive Bayes model predicts feature vectors \underline{x} according to

$$P(\underline{x}) = \sum_{y=\pm 1} P(\underline{x}|y)P(y) = \sum_{y=\pm 1} \left[\prod_{i=1}^d P(x_i|y) \right] P(y) \quad (2)$$

In other words, the distribution $P(\underline{x})$ assumes that the examples come in two different varieties corresponding to their label. This type of unobserved label information is precisely what the mixtures aim to capture.

Let's start with a simple two component mixture of Gaussians model (in two dimensions):

$$P(\underline{x}|\theta) = P(1)N(\underline{x}; \mu_1, \sigma_1^2 I) + P(2)N(\underline{x}; \mu_2, \sigma_2^2 I) \quad (3)$$

The parameters θ defining the mixture model are $P(1)$, $P(2)$, μ_1 , μ_2 , σ_1^2 , and σ_2^2 . Figure ?? shows data generated from such a model. Note that the frequencies $P(j)$ (a.k.a. *mixing proportions*) control the size of the resulting clusters in the data in terms of how many examples they involve, μ_j 's specify the location of cluster centers, and σ_j^2 's control how spread out the clusters are. Note that each example in the figure could in principle have been generated in two possible ways (which mixture component it was sampled from).

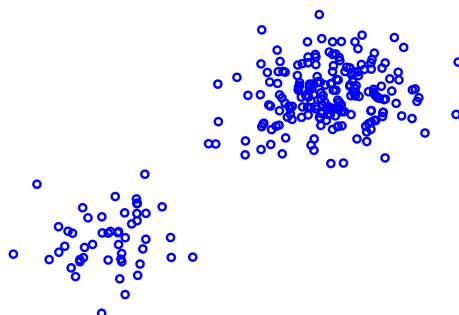


Figure 1: Samples from a mixture of two Gaussians

There are many ways of using mixtures. Consider, for example, the problem of predicting final exam score vectors for students in machine learning. Each observation is a vector \underline{x} with components that specify the points the student received in a particular question. We would expect that different types of students succeed in different types of questions. This “student type” information is not available in the exam score vectors, however, but we can model it.

Student exam model: 1-year

Each \underline{x} is a vector of scores from a particular student and samples correspond to students. We expect that the population of students in a particular year consists of m different types

(e.g., due to differences in background). If we expect each type to be present with an overall probability $P(j)$, then each student score is modeled as a mixture

$$P(\underline{x}|\theta) = \sum_{j=1}^m P(\underline{x}|\theta_j)P(j) \quad (4)$$

where we sum over the types (weighted by $P(j)$) since we don't know what type of student t is prior to seeing their exam score \underline{x}_t .

If there are n students taking the exam a particular year, then the likelihood of all the student score vectors, $D_1 = \{\underline{x}_1, \dots, \underline{x}_n\}$, would be

$$L(D_1; \theta) = \prod_{t=1}^n P(\underline{x}_t|\theta) = \prod_{t=1}^n \left(\sum_{j=1}^m P(\underline{x}_t|\theta_j)P(j) \right) \quad (5)$$

Student exam model: K-years

Suppose now that we have student data from K years of offering the course. In year k we have n_k students who took the course. Let $\underline{x}_{k,t}$ denote the score vector for a student t in year k . Note that t is just an index to identify samples each year and the same index does not imply that the same student took the course multiple years. We can now assume that the number of student types as well as $P(\underline{x}|\theta_j)$ remain the same from year to year (the parameters θ_j are the same for all years). However, the population of students may easily change from year to year, and thus the prior probabilities over the types have to be set differently. Let $P(j|k)$ denote the prior probabilities over the types in year k (all of these would have to be estimated of course). Now, according to our mixture distribution, we expect example scores for students in year k be sampled from

$$P(\underline{x}|k, \theta) = \sum_{j=1}^m P(\underline{x}|\theta_j)P(j|k) \quad (6)$$

The likelihood of all the data, across K years, $D = \{D_1, \dots, D_K\}$, is given by

$$L(D; \theta) = \prod_{k=1}^K \prod_{t=1}^{n_k} P(\underline{x}_{k,t}|k, \theta) = \prod_{k=1}^K \prod_{t=1}^{n_k} \left(\sum_{j=1}^m P(\underline{x}_{k,t}|\theta_j)P(j|k) \right) \quad (7)$$

The parameters θ here include the mixing portions $\{P(j|k)\}$ that change from year to year in addition to $\{\theta_j\}$.

Collaborative filtering

		3	
		5	2
	1		5
2		2	

movies j

users i

r_{ij}

Figure 2: Partially observed rating matrix for a collaborative filtering task.

Mixture models are useful also in recommender systems. Suppose we have n users and m movies and our task is to recommend movies for users. The users have each rated a small fraction of movies and our task is to fill-in the rating matrix, i.e., provide a predicted rating for each user across all m movies. Such a prediction task is known as a collaborative filtering problem (see Figure ??).

Say the possible ratings are $r_{ij} \in \{1, \dots, 5\}$ (i.e., how many stars you assign to each movie). We will use i to index users and j for movies; a rating r_{ij} , if provided, specifies how user i rated movie j . Since only a small fraction of movies are rated by each user, we need a way to index these elements of the user/movie matrix: we say $(i, j) \in I_D$ if rating r_{ij} is available (observed). D denotes all the observed ratings.

We can build on the previous discussion on mixture models. We can represent each movie as a distribution over “movie types” $z_m \in \{1, \dots, K_m\}$. Similarly, a user is represented by a distribution over “user types” $z_u \in \{1, \dots, K_u\}$. We do not assume that each movie corresponds to a single movie type across all users. Instead, we interpret the distribution over movie types as a bag of features corresponding to the movie and we resample from this bag in the context of each user. This is analogous to predicting exam scores for students in a particular year (we didn’t assume that all the students had the same type). We also make the same assumption about user types, i.e., that the type is sampled from the “bag” for each rating, resulting potentially in different types for different movies. Since all the unobserved quantities are summed over, we do not explicitly assign any fixed movie/user type to a rating.

We imagine generating the rating r_{ij} associated with (i, j) element of the rating matrix as

follows. Sample a movie type from $P(z_m|j)$, sample a user type from $P(z_u|i)$, then sample a rating r_{ij} with probability $P(r_{ij}|z_u, z_m)$. All the probabilities involved have to be estimated from the available data. Note that we will resample movie and user types for each rating. The resulting mixture model for rating r_{ij} is given by

$$P(r_{ij}|i, j, \theta) = \sum_{z_u=1}^{K_u} \sum_{z_m=1}^{K_m} P(r_{ij}|z_u, z_m)P(z_u|i)P(z_m|j) \quad (8)$$

where the parameters θ refer to the mapping from types to ratings $\{P(r|z_u, z_m)\}$ and the user and movie specific distributions $\{P(z_u|i)\}$ and $\{P(z_m|j)\}$, respectively. The likelihood of the observed data D is

$$L(D; \theta) = \prod_{(i,j) \in I_D} P(r_{ij}|i, j, \theta) \quad (9)$$

Users rate movies differently. For example, some users may only use a part of the scale (e.g., 3, 4 or 5) while others may be bi-modal, rating movies either very bad 1 or very good 5. We can improve the model by assuming that each user has a rating style $s \in \{1, \dots, K_s\}$. The styles are assumed to be the same for all users, we just don't know how to assign each user to a particular style. The prior probability that any randomly chosen user would have style s is specified by $P(s)$. These parameters are common to all users. We also assume that the rating predictions $P(r_{ij}|z_u, z_m)$ associated with user/movie types now depend on the style s as well: $P(r_{ij}|z_u, z_m, s)$.

We have to be a bit careful in writing the likelihood of the data. All the ratings of one user have to come from one rating style s but we can sum over the possibilities. As a result, the likelihood of the observed ratings is modified to be

$$L'(D; \theta) = \prod_{i=1}^n \left[\sum_{s=1}^{K_s} P(s) \overbrace{\prod_{j:(i,j) \in I_D} \left(\sum_{z_u=1}^{K_u} \sum_{z_m=1}^{K_m} P(r_{ij}|z_u, z_m, s)P(z_u|i)P(z_m|j) \right)}^{\text{likelihood of user } i\text{'s ratings with style } s} \right] \quad (10)$$

The model does not actually involve that many parameters to estimate. There are exactly

$$\underbrace{\{P(s)\}}_{(K_s - 1)} + \underbrace{\{P(r|z_u, z_m, s)\}}_{(5 - 1)K_u K_m K_s} + \underbrace{\{P(z_u|i)\}}_{n(K_u - 1)} + \underbrace{\{P(z_m|j)\}}_{m(K_m - 1)} \quad (11)$$

independent parameters in the model.

A realistic model would include, in addition, a prediction of the “missing elements” in the rating matrix, i.e., a model of why the entry was missing (a user failed to rate a movie they had seen, not seen but could, chosen not to see, etc.).

Estimating mixtures: the EM-algorithm

We have seen a number of different types of mixture models. The advantage of mixture models lies in their ability to incorporate and resolve ambiguities in the data, especially in terms of unidentified sub-groups. However, we can make use of them only if we can estimate such models easily from the available data.

Complete data. The simplest way to understand how to estimate mixture models is to start by pretending that we knew all the sub-typing (component) assignments for each available data point. This is analogous to knowing the label for each example in a classification context. We don't actually know these (they are unobserved in the data) but solving the estimation problem in this context will help us later on.

Let's begin with the simple Gaussian mixture model in Eq.(??),

$$P(\underline{x}; \theta) = \sum_{j=1}^m P(j)N(\underline{x}; \underline{\mu}_j, \Sigma_j) \quad (12)$$

and pretend that each observation $\underline{x}_1, \dots, \underline{x}_n$ also had information about the component that was responsible for generating it, i.e., we also observed j_1, \dots, j_n . This additional component information is convenient to include in the form of 0-1 assignments $\delta(j|t)$ where $\delta(j_t|t) = 1$ and $\delta(j|t) = 0$ for all $j \neq j_t$. The log-likelihood of this *complete data* is

$$l(\underline{x}_1, \dots, \underline{x}_n, j_1, \dots, j_n; \theta) = \sum_{t=1}^n \log \left[P(j_t)N(\underline{x}_t; \underline{\mu}_{j_t}, \Sigma_{j_t}) \right] \quad (13)$$

$$= \sum_{t=1}^n \sum_{j=1}^m \delta(j|t) \log \left[P(j)N(\underline{x}_t; \underline{\mu}_j, \Sigma_j) \right] \quad (14)$$

$$= \sum_{j=1}^m \left(\sum_{t=1}^n \delta(j|t) \right) \log P(j) + \sum_{j=1}^m \left(\sum_{t=1}^n \delta(j|t) \log N(\underline{x}_t; \underline{\mu}_j, \Sigma_j) \right) \quad (15)$$

It's important to note that in trying to maximize this log-likelihood, all the Gaussians can be estimated separately from each other. Put another way, because our pretend observations are "complete", we can estimate each component from only data pertaining to it; the problem of resolving which component should be responsible for which data points is not

present. As a result, the maximizing solution can be written in closed form:

$$\hat{P}(j) = \frac{\hat{n}(j)}{n}, \text{ where } \hat{n}(j) = \sum_{t=1}^n \delta(j|t) \quad (16)$$

$$\hat{\mu}_j = \frac{1}{\hat{n}(j)} \sum_{t=1}^n \delta(j|t) \underline{x}_t \quad (17)$$

$$\hat{\Sigma}_j = \frac{1}{\hat{n}(j)} \sum_{t=1}^n \delta(j|t) (\underline{x}_t - \hat{\mu}_j)(\underline{x}_t - \hat{\mu}_j)^T \quad (18)$$

In other words, the prior probabilities simply recover the empirical fractions from the “observed” j_1, \dots, j_n , and each Gaussian component is estimated in the usual way (evaluating the empirical mean and the covariance) based on data points explicitly assigned to that component. So, the estimation of mixture models would be very easy if we knew the assignments j_1, \dots, j_n .

Incomplete data. What changes if we don’t know the assignments? We can always guess what the assignments should be based on the current setting of the parameters. Let $\theta^{(l)}$ denote the current (initial) parameters of the mixture distribution. Using these parameters, we can evaluate for each data point x_t the posterior probability that it was generated from component j :

$$P(j|\underline{x}_t, \theta^{(l)}) = \frac{P^{(l)}(j)N(\underline{x}_t; \underline{\mu}_j^{(l)}, \Sigma_j^{(l)})}{\sum_{j'=1}^m P^{(l)}(j')N(\underline{x}_t; \underline{\mu}_{j'}^{(l)}, \Sigma_{j'}^{(l)})} = \frac{P^{(l)}(j)N(\underline{x}_t; \underline{\mu}_j^{(l)}, \Sigma_j^{(l)})}{P(\underline{x}_t; \theta^{(l)})} \quad (19)$$

Instead of using the 0-1 assignments $\delta(j|t)$ of data points to components we can use “soft assignments” $p^{(l)}(j|t) = P(j|\underline{x}_t, \theta^{(l)})$. By substituting these in the above closed form estimating equations we get an iterative algorithm for estimating Gaussian mixtures. The algorithm is iterative since the soft posterior assignments were evaluated based on the current setting of the parameters $\theta^{(l)}$ and may have to be revised later on (once we have a better idea of where the clusters are in the data).

The resulting algorithm is known as the *Expectation Maximization algorithm* (EM for short) and applies to all mixture models and beyond. For Gaussian mixtures, the EM-algorithm can be written as

(Step 0) Initialize the Gaussian mixture, i.e., specify $\theta^{(0)}$. A simple initialization consists of setting $P^{(0)}(j) = 1/m$, equating each $\underline{\mu}_j^{(0)}$ with a randomly chosen data point, and making all $\Sigma_j^{(0)}$ equal to the overall data covariance.

(E-step) Evaluate the posterior assignment probabilities $p^{(l)}(j|t) = P(j|\underline{x}_t, \theta^{(l)})$ based on the current setting of the parameters $\theta^{(l)}$.

(M-step) Update the parameters according to

$$P^{(l+1)}(j) = \frac{\hat{n}(j)}{n}, \text{ where } \hat{n}(j) = \sum_{t=1}^n p^{(l)}(j|t) \quad (20)$$

$$\mu_j^{(l+1)} = \frac{1}{\hat{n}(j)} \sum_{t=1}^n p^{(l)}(j|t) \underline{x}_t \quad (21)$$

$$\Sigma_j^{(l+1)} = \frac{1}{\hat{n}(j)} \sum_{t=1}^n p^{(l)}(j|t) (\underline{x}_t - \mu_j^{(l+1)})(\underline{x}_t - \mu_j^{(l+1)})^T \quad (22)$$

Perhaps surprisingly, this iterative algorithm is guaranteed to converge and each iteration increases the log-likelihood of the data. In other words, if we write

$$l(D; \theta^{(l)}) = \sum_{t=1}^n P(\underline{x}_t; \theta^{(l)}) \quad (23)$$

then

$$l(D; \theta^{(0)}) < l(D; \theta^{(1)}) < l(D; \theta^{(2)}) < \dots \quad (24)$$

until convergence. The main downside of this algorithm is that we are only guaranteed to converge to a locally optimal solution where $d/d\theta l(D; \theta) = 0$. In other words, there could be a setting of the parameters for the mixture distribution that leads to a higher log-likelihood of the data¹. For this reason, the algorithm is typically run multiple times (recall the random initialization of the means) so as to ensure we find a reasonably good solution, albeit perhaps only locally optimal.

Example. Consider a simple mixture of two Gaussians. Figure ?? demonstrates how the EM-algorithm changes the Gaussian components after each iteration. The ellipsoids specify one standard deviation distances from the Gaussian mean. The mixing proportions $P(j)$ are not visible in the figure. Note that it takes many iterations for the algorithm to resolve how to properly assign the data points to mixture components. At convergence, the assignments are still soft (not 0-1) but nevertheless clearly divide the responsibilities of the two Gaussian components across the clusters in the data.

¹In fact, the highest likelihood for Gaussian mixtures is always ∞ . This happens when one of the Gaussians concentrates around a single data point. We do not look for such solutions, however, and they can be removed by constraining the covariance matrices or via regularization. The real comparison is to a non-trivial mixture that achieves the highest log-likelihood.

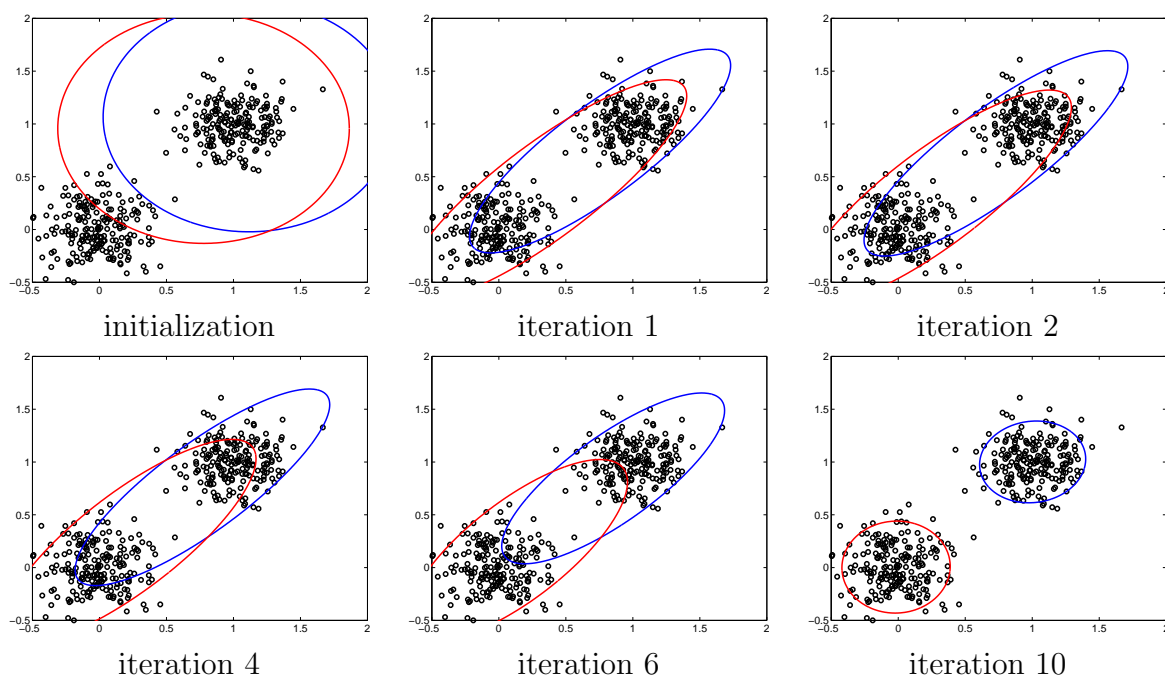


Figure 3: An example of the EM algorithm with a two-component mixture of Gaussians model.