ordering of variable dependencies in the actual network. The program succeeded in reconstructing the correct Bayesian network structure almost exactly, with the exception of one incorrectly deleted arc and one incorrectly added arc.

Constraint-based approaches to learning Bayesian network structure have also been developed (e.g., Spirtes et al. 1993). These approaches infer indepen- dence and dependence relationships from the data, and then use these relation- ships to construct Bayesian networks. Surveys of current approaches to learning Bayesian networks are provided by Heckerman (1995) and Buntine (1994).

## 6.12 THE EM ALGORITHM

In many practical learning settings, only a subset of the relevant instance features might be observable. For example, in training or using the Bayesian belief network of Figure 6.3, we might have data where only a subset of the network variables *Storm, Lightning, Thunder, ForestFire, Campfire*, and *BusTourGroup* have been observed. Many approaches have been proposed to handle the problem of learning in the presence of unobserved variables. As we saw in Chapter 3, if some variable is sometimes observed and sometimes not, then we can use the cases for which it has been observed to learn to predict its values when it is not. In this section we describe the EM algorithm (Dempster et al. 1977), a widely used approach to learning in the presence of unobserved variables. The EM algorithm can be used even for variables whose value is never directly observed, provided the general form of the probability distribution governing these variables is known. The EM algorithm has been used to train Bayesian belief networks (see Heckerman 1995) as well as radial basis function networks discussed in Section 8.4. The EM algorithm is also the basis for many unsupervised clustering algorithms (e.g., Cheeseman et al. 1988), and it is the basis for the widely used Baum-Welch forward-backward algorithm for learning Partially Observable Markov Models (Rabiner 1989).

### 6.12.1 Estimating Means of $k$ Gaussians

The easiest way to introduce the EM algorithm is via an example. Consider a problem in which the data $D$ is a set of instances generated by a probability distribution that is a mixture of $k$ distinct Normal distributions. This problem setting is illustrated in Figure 6.4 for the case where $k = 2$ and where the instances are the points shown along the $x$ axis. Each instance is generated using a two-step process. First, one of the $k$ Normal distributions is selected at random. Second, a single random instance $x_i$ is generated according to this selected distribution. This process is repeated to generate a set of data points as shown in the figure. To simplify our discussion, we consider the special case where the selection of the single Normal distribution at each step is based on choosing each with uniform probability, where each of the $k$ Normal distributions has the same variance $\sigma^2$, and where $\sigma^2$ is known. The learning task is to output a hypothesis $h = \langle \mu_1, \ldots \mu_k \rangle$ that describes the means of each of the $k$ distributions. We would like to find
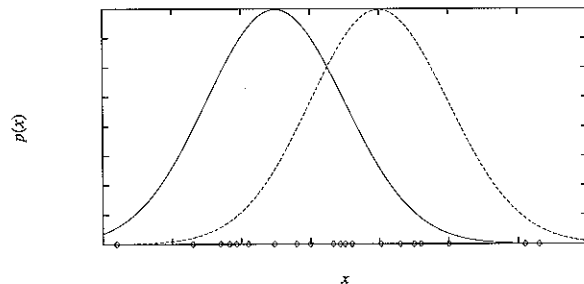
**FIGURE 6.4**
Instances generated by a mixture of two Normal distributions with identical variance $\sigma$. The instances are shown by the points along the $x$ axis. If the means of the Normal distributions are unknown, the EM algorithm can be used to search for their maximum likelihood estimates.

a maximum likelihood hypothesis for these means; that is, a hypothesis $h$ that maximizes $p(D|h)$.

Note it is easy to calculate the maximum likelihood hypothesis for the mean of a single Normal distribution given the observed data instances $x_1, x_2, \ldots, x_m$ drawn from this single distribution. This problem of finding the mean of a single distribution is just a special case of the problem discussed in Section 6.4, Equation (6.6), where we showed that the maximum likelihood hypothesis is the one that minimizes the sum of squared errors over the $m$ training instances. Restating Equation (6.6) using our current notation, we have

$$\mu_{ML} = \underset{\mu}{\operatorname{argmin}} \sum_{i=1}^{m} (x_i - \mu)^2 \tag{6.27}$$

In this case, the sum of squared errors is minimized by the sample mean

$$\mu_{ML} = \frac{1}{m} \sum_{i=1}^{m} x_i \tag{6.28}$$

Our problem here, however, involves a mixture of $k$ different Normal distributions, and we cannot observe which instances were generated by which distribution. Thus, we have a prototypical example of a problem involving hidden variables. In the example of Figure 6.4, we can think of the full description of each instance as the triple $\langle x_i, z_{i1}, z_{i2} \rangle$, where $x_i$ is the observed value of the $i$th instance and where $z_{i1}$ and $z_{i2}$ indicate which of the two Normal distributions was used to generate the value $x_i$. In particular, $z_{ij}$ has the value 1 if $x_i$ was created by the $j$th Normal distribution and 0 otherwise. Here $x_i$ is the observed variable in the description of the instance, and $z_{i1}$ and $z_{i2}$ are hidden variables. If the values of $z_{i1}$ and $z_{i2}$ were observed, we could use Equation (6.27) to solve for the means $\mu_1$ and $\mu_2$. Because they are not, we will instead use the EM algorithm.

Applied to our $k$-means problem the EM algorithm searches for a maximum likelihood hypothesis by repeatedly re-estimating the expected values of the hidden variables $z_{ij}$ given its current hypothesis $\langle \mu_1 \ldots \mu_k \rangle$, then recalculating the

maximum likelihood hypothesis using these expected values for the hidden variables. We will first describe this instance of the EM algorithm, and later state the EM algorithm in its general form.

Applied to the problem of estimating the two means for Figure 6.4, the EM algorithm first initializes the hypothesis to $h = \langle \mu_1, \mu_2 \rangle$, where $\mu_1$ and $\mu_2$ are arbitrary initial values. It then iteratively re-estimates $h$ by repeating the following two steps until the procedure converges to a stationary value for $h$.

**Step 1:** Calculate the expected value $E[z_{ij}]$ of each hidden variable $z_{ij}$, assuming the current hypothesis $h = \langle \mu_1, \mu_2 \rangle$ holds.

**Step 2:** Calculate a new maximum likelihood hypothesis $h' = \langle \mu_1', \mu_2' \rangle$, assuming the value taken on by each hidden variable $z_{ij}$ is its expected value $E[z_{ij}]$ calculated in Step 1. Then replace the hypothesis $h = \langle \mu_1, \mu_2 \rangle$ by the new hypothesis $h' = \langle \mu_1', \mu_2' \rangle$ and iterate.

Let us examine how both of these steps can be implemented in practice. Step 1 must calculate the expected value of each $z_{ij}$. This $E[z_{ij}]$ is just the probability that instance $x_i$ was generated by the $j$th Normal distribution

$$
\begin{aligned}
E[z_{ij}] &= \frac{p(x = x_i | \mu = \mu_j)}{\sum_{n=1}^{2} p(x = x_i | \mu = \mu_n)} \\
&= \frac{e^{-\frac{1}{2\sigma^2}(x_i - \mu_j)^2}}{\sum_{n=1}^{2} e^{-\frac{1}{2\sigma^2}(x_i - \mu_n)^2}}
\end{aligned}
$$

Thus the first step is implemented by substituting the current values $\langle \mu_1, \mu_2 \rangle$ and the observed $x_i$ into the above expression.

In the second step we use the $E[z_{ij}]$ calculated during Step 1 to derive a new maximum likelihood hypothesis $h' = \langle \mu_1', \mu_2' \rangle$. As we will discuss later, the maximum likelihood hypothesis in this case is given by

$$\mu_j \leftarrow \frac{\sum_{i=1}^{m} E[z_{ij}] \, x_i}{\sum_{i=1}^{m} E[z_{ij}]}$$

Note this expression is similar to the sample mean from Equation (6.28) that is used to estimate $\mu$ for a single Normal distribution. Our new expression is just the weighted sample mean for $\mu_j$, with each instance weighted by the expectation $E[z_{ij}]$ that it was generated by the $j$th Normal distribution.

The above algorithm for estimating the means of a mixture of $k$ Normal distributions illustrates the essence of the EM approach: The current hypothesis is used to estimate the unobserved variables, and the expected values of these variables are then used to calculate an improved hypothesis. It can be proved that on each iteration through this loop, the EM algorithm increases the likelihood $P(D|h)$ unless it is at a local maximum. The algorithm thus converges to a local maximum likelihood hypothesis for $\langle \mu_1, \mu_2 \rangle$.

### 6.12.2 General Statement of EM Algorithm

Above we described an EM algorithm for the problem of estimating means of a mixture of Normal distributions. More generally, the EM algorithm can be applied in many settings where we wish to estimate some set of parameters $\theta$ that describe an underlying probability distribution, given only the observed portion of the full data produced by this distribution. In the above two-means example the parameters of interest were $\theta = \langle \mu_1, \mu_2 \rangle$, and the full data were the triples $\langle x_i, z_{i1}, z_{i2} \rangle$ of which only the $x_i$ were observed. In general let $X = \{x_1, \ldots, x_m\}$ denote the observed data in a set of $m$ independently drawn instances, let $Z = \{z_1, \ldots, z_m\}$ denote the unobserved data in these same instances, and let $Y = X \cup Z$ denote the full data. Note the unobserved $Z$ can be treated as a random variable whose probability distribution depends on the unknown parameters $\theta$ and on the observed data $X$. Similarly, $Y$ is a random variable because it is defined in terms of the random variable $Z$. In the remainder of this section we describe the general form of the EM algorithm. We use $h$ to denote the current hypothesized values of the parameters $\theta$, and $h'$ to denote the revised hypothesis that is estimated on each iteration of the EM algorithm.

The EM algorithm searches for the maximum likelihood hypothesis $h'$ by seeking the $h'$ that maximizes $E[\ln P(Y|h')]$. This expected value is taken over the probability distribution governing $Y$, which is determined by the unknown parameters $\theta$. Let us consider exactly what this expression signifies. First, $P(Y|h')$ is the likelihood of the full data $Y$ given hypothesis $h'$. It is reasonable that we wish to find a $h'$ that maximizes some function of this quantity. Second, maximizing the logarithm of this quantity $\ln P(Y|h')$ also maximizes $P(Y|h')$, as we have discussed on several occasions already. Third, we introduce the expected value $E[\ln P(Y|h')]$ because the full data $Y$ is itself a random variable. Given that the full data $Y$ is a combination of the observed data $X$ and unobserved data $Z$, we must average over the possible values of the unobserved $Z$, weighting each according to its probability. In other words we take the expected value $E[\ln P(Y|h')]$ over the probability distribution governing the random variable $Y$. The distribution governing $Y$ is determined by the completely known values for $X$, plus the distribution governing $Z$.

What is the probability distribution governing $Y$? In general we will not know this distribution because it is determined by the parameters $\theta$ that we are trying to estimate. Therefore, the EM algorithm uses its current hypothesis $h$ in place of the actual parameters $\theta$ to estimate the distribution governing $Y$. Let us define a function $Q(h'|h)$ that gives $E[\ln P(Y|h')]$ as a function of $h'$, under the assumption that $\theta = h$ and given the observed portion $X$ of the full data $Y$.

$$Q(h'|h) = E[\ln p(Y|h')|h, X]$$

We write this function $Q$ in the form $Q(h'|h)$ to indicate that it is defined in part by the assumption that the current hypothesis $h$ is equal to $\theta$. In its general form, the EM algorithm repeats the following two steps until convergence:

**Step 1:** *Estimation (E) step:* Calculate $Q(h'|h)$ using the current hypothesis $h$ and the observed data $X$ to estimate the probability distribution over $Y$.

$$Q(h'|h) \leftarrow E[\ln P(Y|h')|h, X]$$

**Step 2:** *Maximization (M) step:* Replace hypothesis $h$ by the hypothesis $h'$ that maximizes this $Q$ function.

$$h \leftarrow \underset{h'}{\operatorname{argmax}} \, Q(h'|h)$$

When the function $Q$ is continuous, the EM algorithm converges to a stationary point of the likelihood function $P(Y|h')$. When this likelihood function has a single maximum, EM will converge to this global maximum likelihood estimate for $h'$. Otherwise, it is guaranteed only to converge to a local maximum. In this respect, EM shares some of the same limitations as other optimization methods such as gradient descent, line search, and conjugate gradient discussed in Chapter 4.

### 6.12.3 Derivation of the $k$ Means Algorithm

To illustrate the general EM algorithm, let us use it to derive the algorithm given in Section 6.12.1 for estimating the means of a mixture of $k$ Normal distributions. As discussed above, the $k$-means problem is to estimate the parameters $\theta = \langle \mu_1 \ldots \mu_k \rangle$ that define the means of the $k$ Normal distributions. We are given the observed data $X = \{\langle x_i \rangle\}$. The hidden variables $Z = \{\langle z_{i1}, \ldots, z_{ik} \rangle\}$ in this case indicate which of the $k$ Normal distributions was used to generate $x_i$.

To apply EM we must derive an expression for $Q(h|h')$ that applies to our $k$-means problem. First, let us derive an expression for $\ln p(Y|h')$. Note the probability $p(y_i|h')$ of a single instance $y_i = \langle x_i, z_{i1}, \ldots z_{ik} \rangle$ of the full data can be written

$$p(y_i|h') = p(x_i, z_{i1}, \ldots, z_{ik}|h') = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \sum_{j=1}^{k} z_{ij}(x_i - \mu_j')^2}$$

To verify this note that only one of the $z_{ij}$ can have the value 1, and all others must be 0. Therefore, this expression gives the probability distribution for $x_i$ generated by the selected Normal distribution. Given this probability for a single instance $p(y_i|h')$, the logarithm of the probability $\ln P(Y|h')$ for all $m$ instances in the data is

$$\ln P(Y|h') = \ln \prod_{i=1}^{m} p(y_i|h')$$

$$= \sum_{i=1}^{m} \ln p(y_i|h')$$

$$= \sum_{i=1}^{m} \left( \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \sum_{j=1}^{k} z_{ij}(x_i - \mu_j')^2 \right)$$

Finally we must take the expected value of this $\ln P(Y|h')$ over the probability distribution governing $Y$ or, equivalently, over the distribution governing the unobserved components $z_{ij}$ of $Y$. Note the above expression for $\ln P(Y|h')$ is a linear function of these $z_{ij}$. In general, for any function $f(z)$ that is a *linear* function of $z$, the following equality holds

$$E[f(z)] = f(E[z])$$

This general fact about linear functions allows us to write

$$E[\ln P(Y|h')] = E\left[\sum_{i=1}^{m}\left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}\sum_{j=1}^{k} z_{ij}(x_i - \mu_j')^2\right)\right]$$

$$= \sum_{i=1}^{m}\left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}\sum_{j=1}^{k} E[z_{ij}](x_i - \mu_j')^2\right)$$

To summarize, the function $Q(h'|h)$ for the $k$ means problem is

$$Q(h'|h) = \sum_{i=1}^{m}\left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}\sum_{j=1}^{k} E[z_{ij}](x_i - \mu_j')^2\right)$$

where $h' = \langle \mu_1', \ldots, \mu_k' \rangle$ and where $E[z_{ij}]$ is calculated based on the current hypothesis $h$ and observed data $X$. As discussed earlier

$$E[z_{ij}] = \frac{e^{-\frac{1}{2\sigma^2}(x_i-\mu_j)^2}}{\sum_{n=1}^{k} e^{-\frac{1}{2\sigma^2}(x_i-\mu_n)^2}} \tag{6.29}$$

Thus, the first (estimation) step of the EM algorithm defines the $Q$ function based on the estimated $E[z_{ij}]$ terms. The second (maximization) step then finds the values $\mu_1', \ldots, \mu_k'$ that maximize this $Q$ function. In the current case

$$\underset{h'}{\operatorname{argmax}}\ Q(h'|h) = \underset{h'}{\operatorname{argmax}} \sum_{i=1}^{m}\left(\ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}\sum_{j=1}^{k} E[z_{ij}](x_i - \mu_j')^2\right)$$

$$= \underset{h'}{\operatorname{argmin}} \sum_{i=1}^{m}\sum_{j=1}^{k} E[z_{ij}](x_i - \mu_j')^2 \tag{6.30}$$

Thus, the maximum likelihood hypothesis here minimizes a weighted sum of squared errors, where the contribution of each instance $x_i$ to the error that defines $\mu_j'$ is weighted by $E[z_{ij}]$. The quantity given by Equation (6.30) is minimized by setting each $\mu_j'$ to the weighted sample mean

$$\mu_j \leftarrow \frac{\sum_{i=1}^{m} E[z_{ij}]\ x_i}{\sum_{i=1}^{m} E[z_{ij}]} \tag{6.31}$$

Note that Equations (6.29) and (6.31) define the two steps in the $k$-means algorithm described in Section 6.12.1.