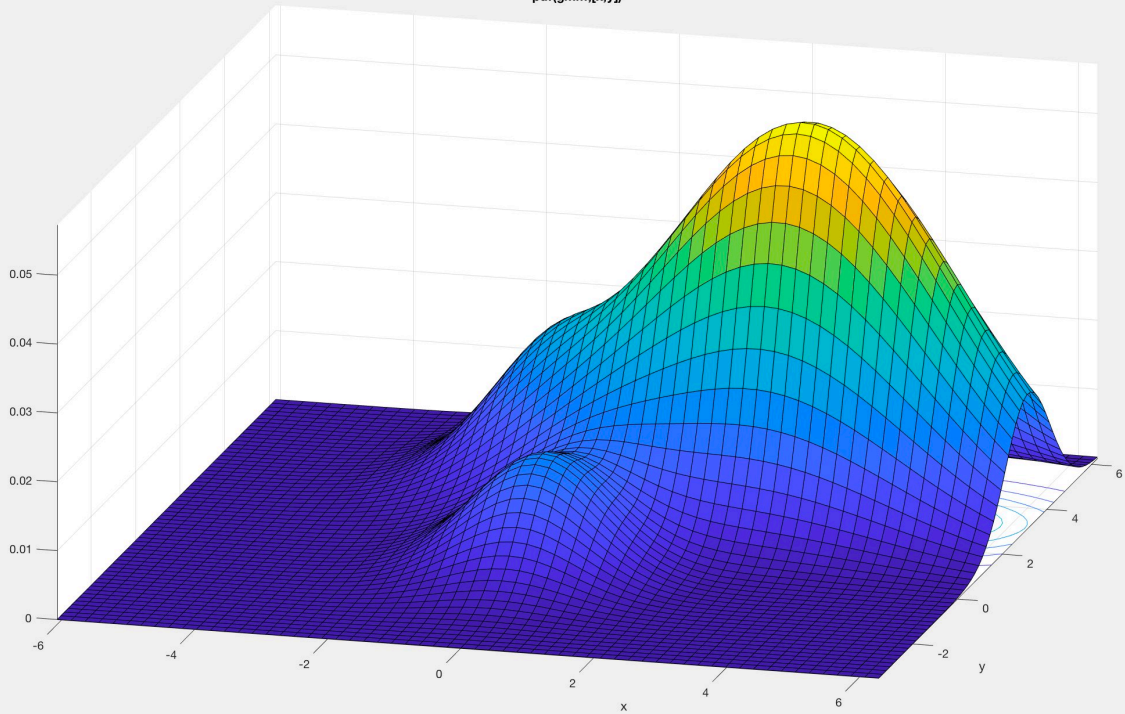
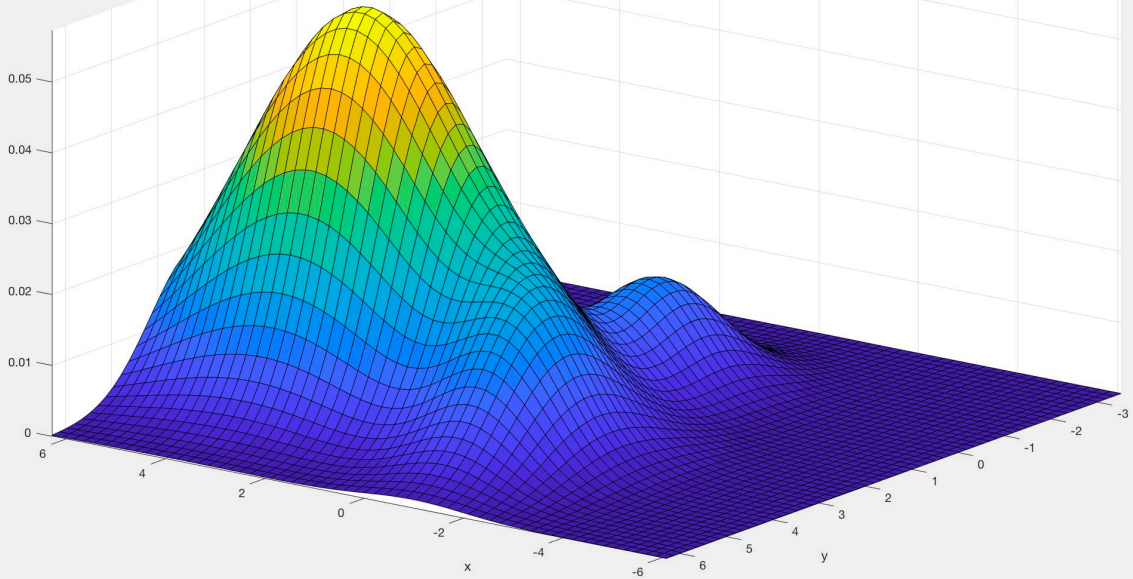


pdf(gmm,[x,y])



pdf(gmm,[x,y])



Using EM To Estimate A Probability Density With A Mixture Of Gaussians

Aaron A. D'Souza

adsouza@usc.edu

1 Introduction

The problem we are trying to address in this note is simple. Given a set of data points $\mathbf{X} = \{\mathbf{x}_i\}_1^N$, we wish to determine the underlying probability distribution $p(\mathbf{x})$, that generates this data. In general, the distribution $p(\mathbf{x})$ could be any real-valued, scalar function with the following constraints:

$$p(\mathbf{x}) \geq 0, \forall \mathbf{x} \text{ and } \int_{-\infty}^{\infty} p(\mathbf{x}) d\mathbf{x} = 1$$

The framework of mixture modeling using Gaussians makes the following assumptions:

- The data was generated using a *set* of M probability distributions.
- Each of the individual probability distributions is a Gaussian:

$$\mathbf{x}; \theta_m \sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)^*$$

i.e., the probability of generating a data point \mathbf{x} under the m^{th} model is given according to a Gaussian distribution with mean $\boldsymbol{\mu}_m$ and covariance $\boldsymbol{\Sigma}_m$.

- Each data point is generated according to the following algorithm:

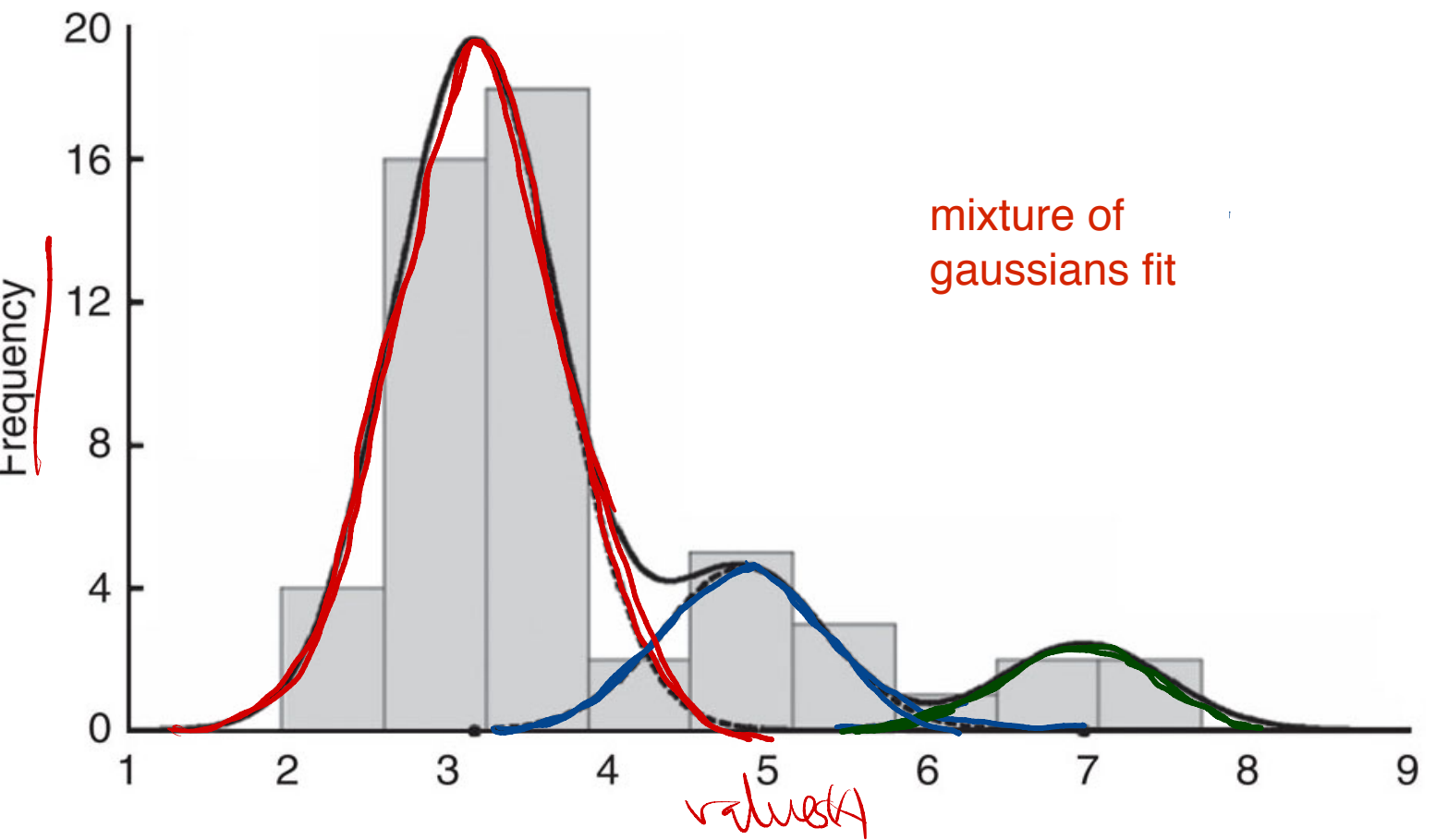
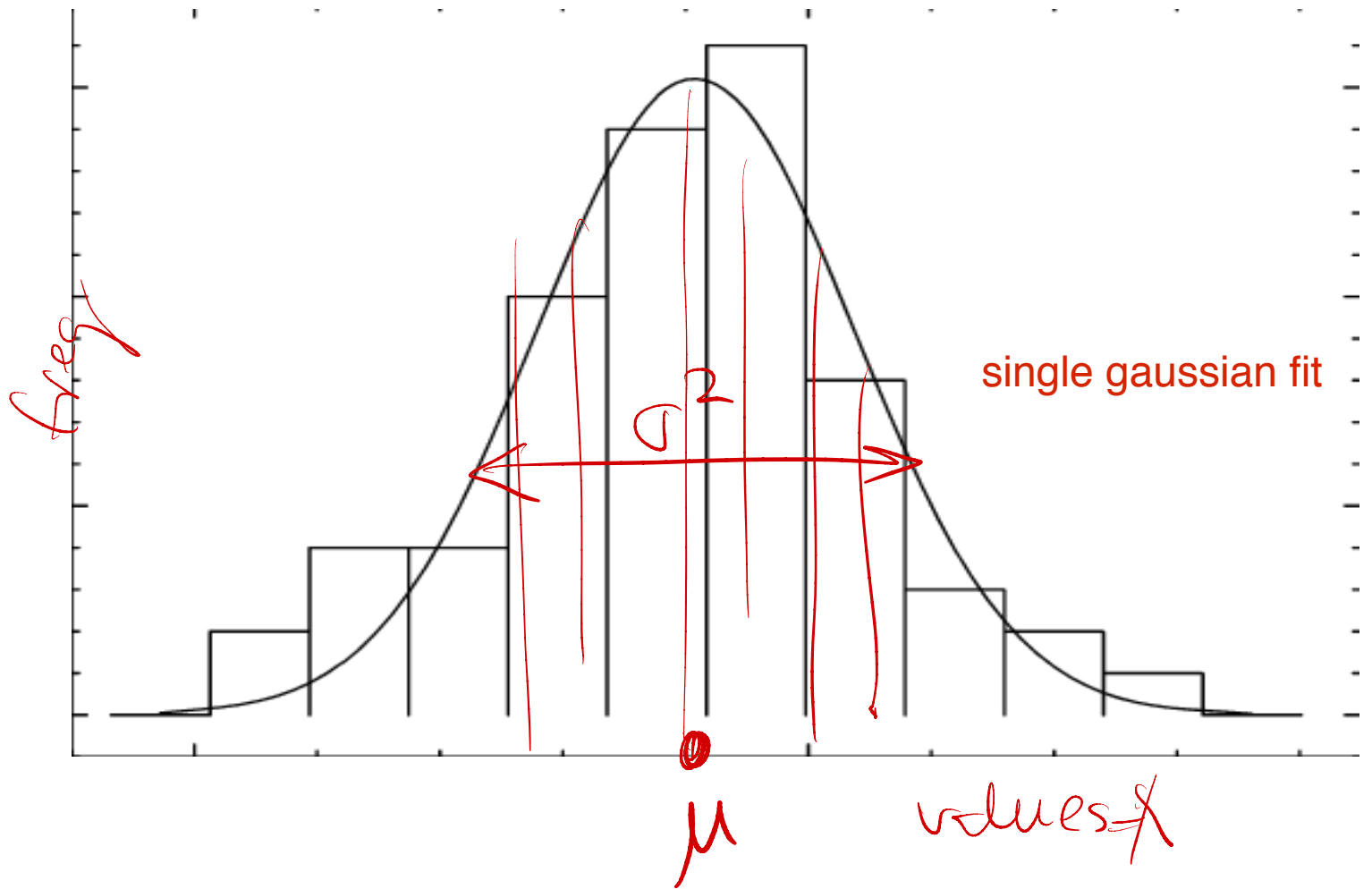
- 1: **for** $i = 1$ to N **do**
- 2: $m \leftarrow$ index of one of the M models randomly selected according to the prior probability vector $\boldsymbol{\pi}$ → mixture coef
- 3: Randomly generate \mathbf{x}_i according to the distribution $\mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$
- 4: **end for**

Representing probabilistic systems as graphical models is rapidly becoming a useful tool in Bayesian analysis. The graphical model corresponding to our formulation of the data generation process is shown in fig. 1. We adopt the convention that circular nodes correspond to random variables in the model, while rectangular nodes correspond to variables that parameterize the distributions of these random variables. Circular nodes with double borders indicate observed variables in our model (the data).

*We use the notation $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ to denote the multivariate Gaussian distribution which is mathematically defined as:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

where d is the dimensionality of \mathbf{x} .



Mixtures of Gaussians (2)

$$Z_{ik} = \begin{cases} 1 & \text{from cluster } k \\ 0 & \text{otherwise} \end{cases}$$

$$\text{membership } \langle Z_{ik} \rangle = \text{prob}(K | x_i)$$

Combine simple models into a complex model:

$$n_k = E[\# \text{ datapoints from cluster } k] = \sum_i Z_{ik}$$

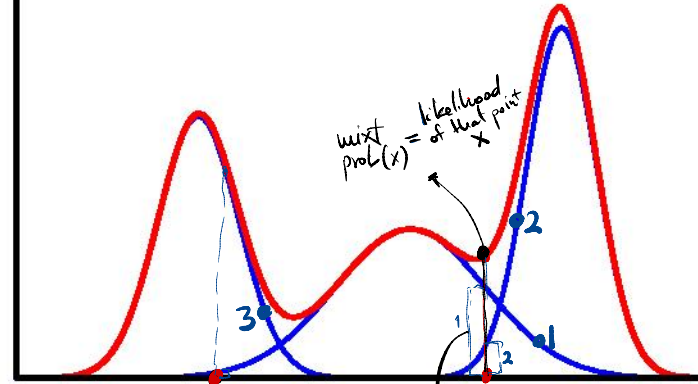
mixture

$$p(\mathbf{x}) = \sum_{k=1}^M \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)$$

generator source
 $\text{prob}(x | \text{cluster } k)$
 $p(x)$
 Component
 Mixing coefficient = $\frac{n_k}{n}$
 weight of cluster k

$$\forall k : \pi_k \geq 0$$

$$\sum_{k=1}^M \pi_k = 1$$



$K=3$
 $\text{prob}(Z_{ik}=1 | x_i)$
 = likelihood of x_i coming from source k

cluster k size $\pi_k = \text{Probab}(Z_{ik}=1)$ for any datapoint i

membership for x_i $\langle Z_{ik} \rangle = \text{Prob}(Z_{ik}=1 | x_i)$

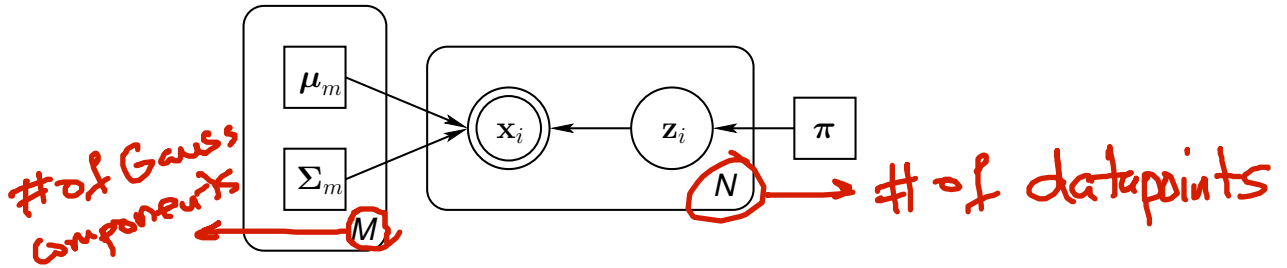


Figure 1: Graphical model for maximum likelihood density estimation using a mixture of Gaussians

In this model we have introduced an additional variable \mathbf{z}_i associated with each \mathbf{x}_i . The \mathbf{z}_i variables are *indicator* variables that are multinomially distributed according to the parameter vector $\boldsymbol{\pi}$, and indicate which component generates the corresponding \mathbf{x}_i . It is easiest to think of each \mathbf{z}_i as an M dimensional vector with a 1 in the element corresponding to the selected mixture component, and 0's in all other elements. The probability of the m^{th} element being 1, is π_m .

$$\mathbf{z}_i = \underbrace{[0 \ 0 \ \dots \ 0 \ 1 \ 0 \ \dots \ 0 \ 0]^T}_{M \text{ elements}}$$

Since we do not know the corresponding \mathbf{z}_i for each \mathbf{x}_i (if we did, then we simply group the \mathbf{x}_i according to their \mathbf{z}_i , and fit a single Gaussian to each group), these variables are called *hidden* variables.

Our problem is now reduced to finding the values of the model parameters $\boldsymbol{\mu}_m$ and $\boldsymbol{\Sigma}_m$ for each of the M models, as well as the prior probability vector $\boldsymbol{\pi}$, which when plugged into the generative model, is most likely to generate the observed data distribution. In other words we are interested in maximizing the *likelihood* $\mathcal{L}(\boldsymbol{\theta}) = p(\mathbf{X}; \boldsymbol{\theta})$ of generating the observed data given the model parameters $\boldsymbol{\theta} = \{\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m, \pi_m\}_1^M$. This approach is called the Maximum Likelihood (ML) framework since it finds the parameter settings that maximize the likelihood of observing the data.

Although the ML approach is an intuitively appealing solution, we often find that maximizing the expressions for likelihood w.r.t. the parameters $\boldsymbol{\theta}$ are often analytically intractable. The Expectation Maximization (EM) algorithm can be used to simplify the math considerably.

2 Estimating the Model Parameters using EM

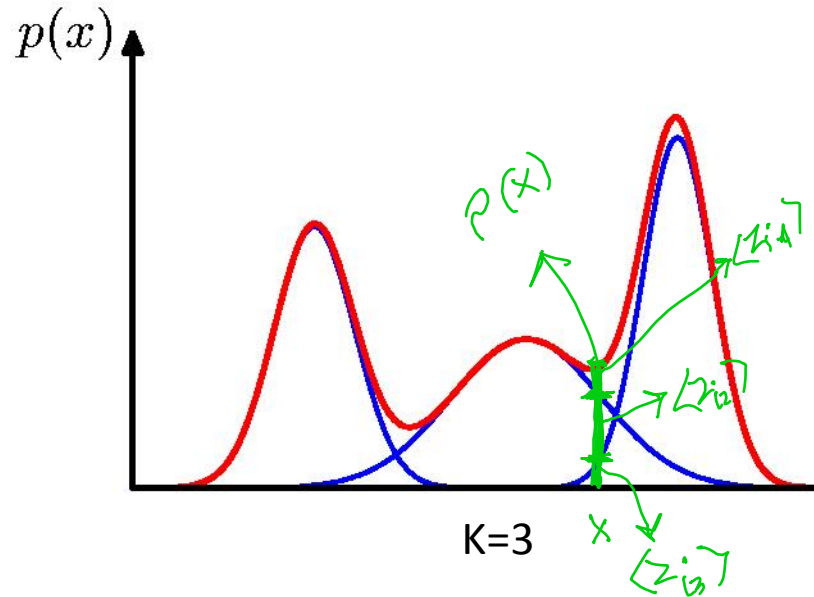
Instead of attempting to maximize the likelihood of the observed data $p(\mathbf{X}; \boldsymbol{\theta})$, we attempt instead to maximize the likelihood of the joint distribution of \mathbf{X} and $\mathbf{Z} = \{\mathbf{z}_i\}_1^N$, $p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})$. For the purposes of maximization we can also work with the logarithm of this quantity, $\log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})$. This quantity is also known as the *complete* log-likelihood. Since we cannot observe the values of the random variables \mathbf{z}_i we must work with the expectation of this quantity w.r.t. some distribution $Q(\mathbf{Z})$.

Mixtures of Gaussians (2)

Combine simple models into a complex model:

$$p(\mathbf{x}) = \sum_{k=1}^M \underbrace{\pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)}_{\text{Component}} \cdot \underbrace{\langle z_{ik} \rangle}_{\text{Mixing coefficient}}$$

$$\forall k : \pi_k \geq 0 \quad \sum_{k=1}^M \pi_k = 1$$



$\pi_k = \sum_{i=1}^n z_{ik}$ Probab ($z_{ik}=1$) for any datapoint i in the whole comp/gauss

Expectation

$\langle z_{ik} \rangle =$ mixture weight of k^{th} component for datapoint x_i

$$\pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k) / \sum_{t=1}^M \pi_t \mathcal{N}(x_i | \mu_t, \Sigma_t)$$

task compute param (μ_k, Σ_k, π_k) and membership (z_{ik})

to max $\rightarrow \log [P(\text{data observed} | \text{models})]$

The log of the complete data likelihood can be written as follows:

group/source cluster $m=k$

$$l_c(\theta) = \log p(\mathbf{X}, \mathbf{Z}; \theta) = \log \left(\prod_{i=1}^N p(\mathbf{x}_i, \mathbf{z}_i; \theta) \right)$$

datapoints = indep

$$= \log \prod_{i=1}^N \prod_{m=1}^M [p(\mathbf{x}_i | z_{im} = 1; \theta) p(z_{im} = 1)]^{z_{im}}$$

$P = \sum \pi_k p_k$

$$= \sum_{i=1}^N \sum_{m=1}^M z_{im} \log p(\mathbf{x}_i | z_{im} = 1; \theta) + z_{im} \log \pi_m$$

Since we have assumed that each of the individual models is a Gaussian, the quantity $p(\mathbf{x}_i | m, \theta)$ is simply the conditional probability of generating \mathbf{x}_i given that the m^{th} model is chosen:

Gauss \rightarrow five dist of points from sources

$$\log p(\mathbf{x}_i | z_{im} = 1; \theta) = \frac{1}{(2\pi)^{d/2} |\Sigma_m|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \mu_m)^T \Sigma_m^{-1} (\mathbf{x}_i - \mu_m) \right\} \quad (1)$$

Taking expectations w.r.t. $Q(\mathbf{Z})$ we get:

$E[\text{likelihood}]$ over all possible \mathbf{Z} membership

$$\langle l_c(\theta) \rangle = \sum_{i=1}^N \sum_{m=1}^M \langle z_{im} \rangle \log p(\mathbf{x}_i | z_{im} = 1; \theta) + \langle z_{im} \rangle \log \pi_m \quad (2)$$

2.1 The M step

The "M" step in EM takes the expected complete log-likelihood as defined in eq. (2) and maximizes it w.r.t. the parameters that are to be estimated; in this case $\pi_m, \mu_m,$ and Σ_m .

Differentiating eq. (2) w.r.t. μ_m we get: $k = \text{fixed}$

diff(likelihood) w.r.t μ

$$\frac{\partial \langle l_c(\theta) \rangle}{\partial \mu_m} = \sum_{i=1}^N \langle z_{im} \rangle \frac{\partial}{\partial \mu_m} \log p(\mathbf{x}_i | z_{im} = 1; \theta) = 0 \quad (3)$$

We can compute $\frac{\partial}{\partial \mu_m} \log p(\mathbf{x}_i | z_{im} = 1; \theta)$ using eq. (1) as follows:

fixed $m=k$

$$\frac{\partial}{\partial \mu_m} \log p(\mathbf{x}_i | z_{im} = 1; \theta) = \frac{\partial}{\partial \mu_m} \log \left\{ \frac{1}{(2\pi)^{d/2} |\Sigma_m|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \mu_m)^T \Sigma_m^{-1} (\mathbf{x}_i - \mu_m) \right\} \right\}$$

$$= -\frac{1}{2} \frac{\partial}{\partial \mu_m} (\mathbf{x}_i - \mu_m)^T \Sigma_m^{-1} (\mathbf{x}_i - \mu_m)$$

$$= (\mathbf{x}_i - \mu_m)^T \Sigma_m^{-1}$$

$\Sigma_m \rightarrow (A + A^T)$
 $A = \Sigma_m^{-1} = A^T$

sum $\sum_{i=1}^N$

Substituting this result into eq. (3), we get:

$$\sum_{i=1}^N \langle z_{im} \rangle (\mathbf{x}_i - \mu_m)^T \Sigma_m^{-1} = 0$$

$d = 1 \text{ dim}$

[†]Where we have used the relation

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$$

diff w.r.t vector $\mathbf{x} = (\mathbf{x}_i - \mu_m)$

$$\mathbf{A} = \Sigma_m^{-1} = \mathbf{A}^T$$

$$\sum \langle z_{im} \rangle \mu_m = \dots = \sum \langle z_{im} \rangle x_i$$

giving us the update equation:

$$\mu_m = \frac{\sum_{i=1}^N \langle z_{im} \rangle \mathbf{x}_i}{\sum_{i=1}^N \langle z_{im} \rangle} \quad (4)$$

weighted avg.

Differentiating eq. (2) w.r.t. Σ_m^{-1} we get:

$$\frac{\partial \langle l_c(\theta) \rangle}{\partial \Sigma_m^{-1}} = \sum_{i=1}^N \langle z_{im} \rangle \frac{\partial}{\partial \Sigma_m^{-1}} \log p(\mathbf{x}_i | z_{im} = 1; \theta) = 0 \quad (5)$$

We can compute $\frac{\partial}{\partial \Sigma_m^{-1}} \log p(\mathbf{x}_i | z_{im} = 1; \theta)$ using eq. (1) as follows:

$$\begin{aligned} \frac{\partial}{\partial \Sigma_m^{-1}} \log p(\mathbf{x}_i | z_{im} = 1; \theta) &= \frac{\partial}{\partial \Sigma_m^{-1}} \log \left\{ \frac{1}{(2\pi)^{d/2} |\Sigma_m|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \mu_m)^T \Sigma_m^{-1} (\mathbf{x}_i - \mu_m) \right\} \right\} \\ &= \frac{\partial}{\partial \Sigma_m^{-1}} \left\{ \frac{1}{2} \log |\Sigma_m^{-1}| - \frac{1}{2} (\mathbf{x}_i - \mu_m)^T \Sigma_m^{-1} (\mathbf{x}_i - \mu_m) \right\} \\ &= \frac{1}{2} \Sigma_m - \frac{1}{2} (\mathbf{x}_i - \mu_m) (\mathbf{x}_i - \mu_m)^T \ddagger \end{aligned}$$

|A| = determinant
 $\partial \log x = \frac{1}{x}$
 $\partial(ax) = a$

Substituting this result into eq. (5), we get:

$$\sum_{i=1}^N \langle z_{im} \rangle \left(\frac{1}{2} \Sigma_m - \frac{1}{2} (\mathbf{x}_i - \mu_m) (\mathbf{x}_i - \mu_m)^T \right) = 0$$

giving us the update equation:

$$\Sigma_m = \frac{\sum_{i=1}^N \langle z_{im} \rangle (\mathbf{x}_i - \mu_m) (\mathbf{x}_i - \mu_m)^T}{\sum_{i=1}^N \langle z_{im} \rangle} \quad (6)$$

In order to maximize the expected log-likelihood in eq. (2) w.r.t. π_m , we have to keep in mind that the maximization has the constraint that $\sum_{m=1}^M \pi_m = 1$. In order to enforce this constraint we use the Lagrange multiplier λ , and augment eq. (2) as follows:

$$L'(\theta) = \langle l_c(\theta) \rangle_{Q(\mathbf{z})} - \lambda \left(\sum_{m=1}^M \pi_m - 1 \right) \quad (7)$$

mixt weights constraint

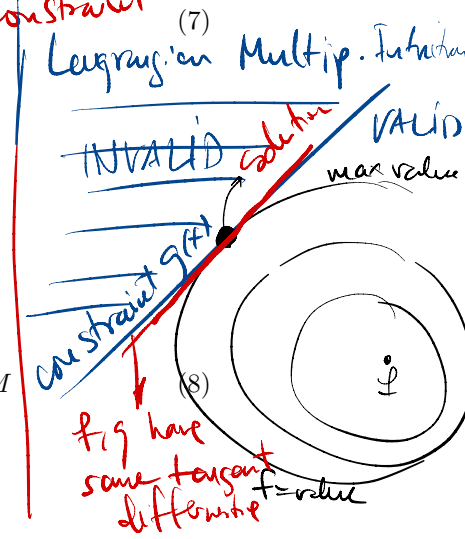
We now differentiate this new expression w.r.t. each π_m giving us:

$$\frac{\partial}{\partial \pi_m} \langle l_c(\theta) \rangle_{Q(\mathbf{z})} - \lambda = 0 \quad \text{for } 1 \leq m \leq M$$

Using eq. (2) we get:

$$\left. \begin{aligned} \frac{1}{\pi_m} \sum_{i=1}^N \langle z_{im} \rangle - \lambda &= 0 \\ \text{or equivalently } \sum_{i=1}^N \langle z_{im} \rangle - \lambda \pi_m &= 0 \end{aligned} \right\} \quad \text{for } 1 \leq m \leq M$$

\ddagger Where we have used the relation $\frac{\partial}{\partial \mathbf{X}} \log |\mathbf{X}| = (\mathbf{X}^{-1})^T$ and $\frac{\partial}{\partial \mathbf{A}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x} \mathbf{x}^T$



Summing eq. (8) over all M models we get:

$$\sum_m \sum_{i=1}^N \langle z_{im} \rangle - \lambda \sum_m \pi_m = 0$$

But since $\sum_m \pi_m = 1$ we get:

$$\lambda = \sum_m \sum_{i=1}^N \langle z_{im} \rangle = N \quad (9)$$

Substituting this result back into eq. (8) we get the following update equation:

$$\pi_m = \frac{\sum_{i=1}^N \langle z_{im} \rangle}{N} = \frac{n_k}{n} \quad (10)$$

which preserves the constraint that $\sum_m \pi_m = 1$.

2.2 The E step

compute membership / param (w/zt)
 $\langle z_{ik} \rangle = \pi_{ik}$
 μ_k, Σ_k, π_k

Now that we have derived the update equations that maximize the expected *complete* log-likelihood ($\log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})$), we wish to ensure that we are indeed also maximizing the *incomplete* log-likelihood $p(\mathbf{X}; \boldsymbol{\theta})$ (which is the quantity that we are truly interested in maximizing).

As we mentioned earlier in section 2, we are guaranteed to maximize the incomplete log-likelihood only when the expectation is taken w.r.t. the posterior distribution of \mathbf{Z} , namely $p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})$. Hence each of the expectations $\langle z_{im} \rangle$ that appear in the update equations derived in the previous section (section 2.1), should be computed as follows:

$$\begin{aligned} \langle z_{im} \rangle_{p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})} &= \textcircled{1} p(z_{im} = 1 | \mathbf{x}_i; \boldsymbol{\theta}) + \textcircled{0} p(z_{im} = 0 | \mathbf{x}_i; \boldsymbol{\theta}) \\ &= p(z_{im} = 1 | \mathbf{x}_i; \boldsymbol{\theta}) \\ &= \frac{p(\mathbf{x}_i | z_{im} = 1; \boldsymbol{\theta}) p(z_{im} = 1)}{\sum_{k=1}^M p(\mathbf{x}_i | z_{ik} = 1; \boldsymbol{\theta}) p(z_{ik} = 1)} \\ &= \frac{p(\mathbf{x}_i | z_{im} = 1; \boldsymbol{\theta}) \pi_m}{\sum_{k=1}^M p(\mathbf{x}_i | z_{ik} = 1; \boldsymbol{\theta}) \pi_k} \quad \text{next comp normalized} \end{aligned}$$

3 Summary

To summarize, given a set of data points \mathbf{X} , if we wish to estimate the underlying probability distribution using EM to fit M Gaussians, we apply Algorithm 1, iterating until convergence of the model parameters.

4 Why does this work? (A brief review of EM theory)

Knowing that $\log p(\mathbf{X}; \boldsymbol{\theta})$ is difficult to maximize analytically, we (seemingly arbitrarily) chose to maximize the expected *complete* log-likelihood $\langle \log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) \rangle_{Q(\mathbf{Z})}$ in the hope that this also increases the *incomplete* log-likelihood $\log p(\mathbf{X}; \boldsymbol{\theta})$ (the quantity we are *really* interested in). This section will justify this choice and prove that we are indeed maximizing $\log p(\mathbf{X}; \boldsymbol{\theta})$

Initialize: all $\langle z_{im} \rangle$, π_m , $\boldsymbol{\mu}_m$, and $\boldsymbol{\Sigma}_m$

1: **repeat**
 2: **for** $i = 1$ to N **do** //The E step
 3: **for** $m = 1$ to M **do**
 4:

$$p(\mathbf{x}_i | z_{im} = 1; \boldsymbol{\theta}) = (2\pi)^{-d/2} |\boldsymbol{\Sigma}_m|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_m)^T \boldsymbol{\Sigma}_m^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_m) \right\}$$

$$\langle z_{im} \rangle = \frac{p(\mathbf{x}_i | z_{im} = 1; \boldsymbol{\theta}) \pi_m}{\sum_j^M p(\mathbf{x}_i | z_{ij} = 1; \boldsymbol{\theta}) \pi_j}$$

5: **end for**
 6: **end for**
 7: **for** $m = 1$ to M **do** //The M step
 8:

$$\boldsymbol{\Sigma}_m = \frac{\sum_{i=1}^N \langle z_{im} \rangle (\mathbf{x}_i - \boldsymbol{\mu}_m) (\mathbf{x}_i - \boldsymbol{\mu}_m)^T}{\sum_{i=1}^N \langle z_{im} \rangle}$$

$$\boldsymbol{\mu}_m = \frac{\sum_{i=1}^N \langle z_{im} \rangle \mathbf{x}_i}{\sum_{i=1}^N \langle z_{im} \rangle}$$

$$\pi_m = \frac{\sum_{i=1}^N \langle z_{im} \rangle}{N}$$

9: **end for**
 10: **until** model parameters converge

Algorithm 1: Estimate $\pi_m, \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m$ for $1 \leq m \leq M$

Let us rewrite $\log p(\mathbf{X} | \boldsymbol{\theta})$ as follows:

$$\begin{aligned} \log p(\mathbf{X}; \boldsymbol{\theta}) &= \log \int p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) d\mathbf{Z} \\ &= \log \int Q(\mathbf{Z}) \frac{p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})}{Q(\mathbf{Z})} d\mathbf{Z} \end{aligned} \quad (11)$$

$$\geq \int Q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})}{Q(\mathbf{Z})} d\mathbf{Z} \quad (\text{Jensen's inequality}) \quad (12)$$

$$= \int Q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) d\mathbf{Z} - \int Q(\mathbf{Z}) \log Q(\mathbf{Z}) d\mathbf{Z} \quad (13)$$

Hence we have the following lower-bound to $\log p(\mathbf{X}; \boldsymbol{\theta})$:

$$\log p(\mathbf{X}; \boldsymbol{\theta}) \geq \underbrace{\langle \log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) \rangle_{Q(\mathbf{Z})}}_{\text{exp. comp. log-lik.}} + \underbrace{\mathcal{H}[Q(\mathbf{Z})]}_{\text{entropy of } Q(\mathbf{Z})} = \mathcal{F}(Q, \boldsymbol{\theta}) \quad (14)$$

Since $Q(\mathbf{Z})$ is an arbitrary distribution, it is independent of $\boldsymbol{\theta}$. Hence in order to maximize the functional $\mathcal{F}(Q, \boldsymbol{\theta})$ w.r.t. $\boldsymbol{\theta}$, it suffices to simply maximize $\langle \log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) \rangle_{Q(\mathbf{Z})}$. (Hence the M-step).

Does this maximization achieve our aim? Eq. (14) shows that the functional $\mathcal{F}(Q, \boldsymbol{\theta})$ is a *lower bound* to the quantity we are interested in. In which case, maximizing $\mathcal{F}(Q, \boldsymbol{\theta})$ does not

guarantee that we are improving $\log p(\mathbf{X}; \boldsymbol{\theta})$ at all! If however, we set $Q(\mathbf{Z}) = p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})$ in Eq. (12), then we see that *the lower bound in fact becomes an equality*.

$$\begin{aligned}
 \int Q(\mathbf{Z}) \log \frac{p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})}{Q(\mathbf{Z})} d\mathbf{Z} &= \int p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}) \log \frac{p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta})}{p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})} d\mathbf{Z} \\
 &= \int p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}) \log \frac{p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})p(\mathbf{X}; \boldsymbol{\theta})}{p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})} d\mathbf{Z} \\
 &= \int p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}) \log p(\mathbf{X}; \boldsymbol{\theta}) d\mathbf{Z} \\
 &= \log p(\mathbf{X}; \boldsymbol{\theta}) \int p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}) d\mathbf{Z} \\
 &= \log p(\mathbf{X}; \boldsymbol{\theta})
 \end{aligned}$$

This means that when computing the expected complete log-likelihood $\langle \log p(\mathbf{X}, \mathbf{Z}; \boldsymbol{\theta}) \rangle_{Q(\mathbf{Z})}$, the expectation should be taken w.r.t. the true posterior $p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})$ of the hidden variables (the E step).

5 Examples

Figures 2, 3 and 4 show the result of fitting 2, 3, and 4 Gaussians respectively to a set of data points. In each figure the first plot shows the positions of the means and the relative covariances of each Gaussian, while the second shows the resulting estimated distribution obtained by marginalizing over the models as follows:

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_m^M p(\mathbf{x}|z_m = 1; \boldsymbol{\theta})p(z_m = 1)$$

Figure 5 demonstrates that one can model rather arbitrary non-gaussian distributions provided we have a sufficient number of mixture components. The tradeoff here is that too few components will fail to model the structure of the data, while too many will “overfit” the data. The *model selection* problem in this context is determining an appropriate compromise.

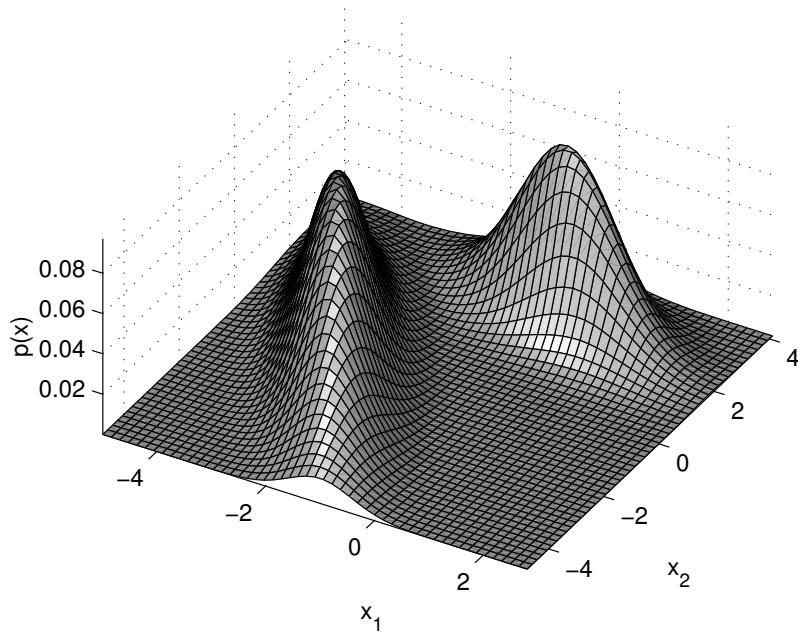
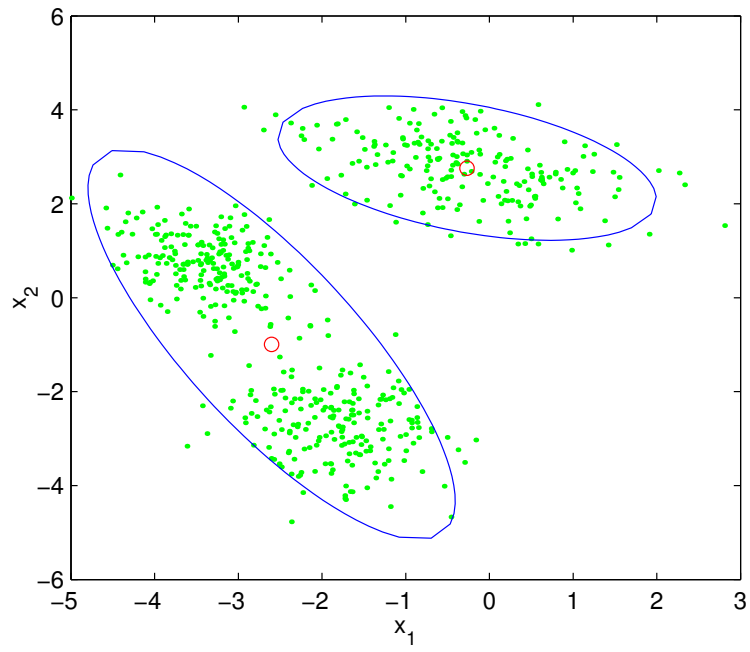


Figure 2: Fitting with 2 Gaussians

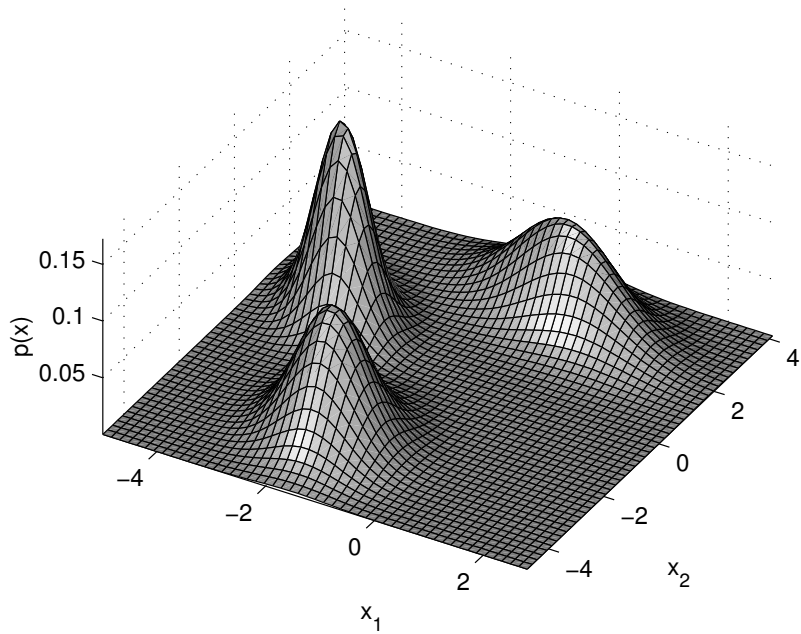
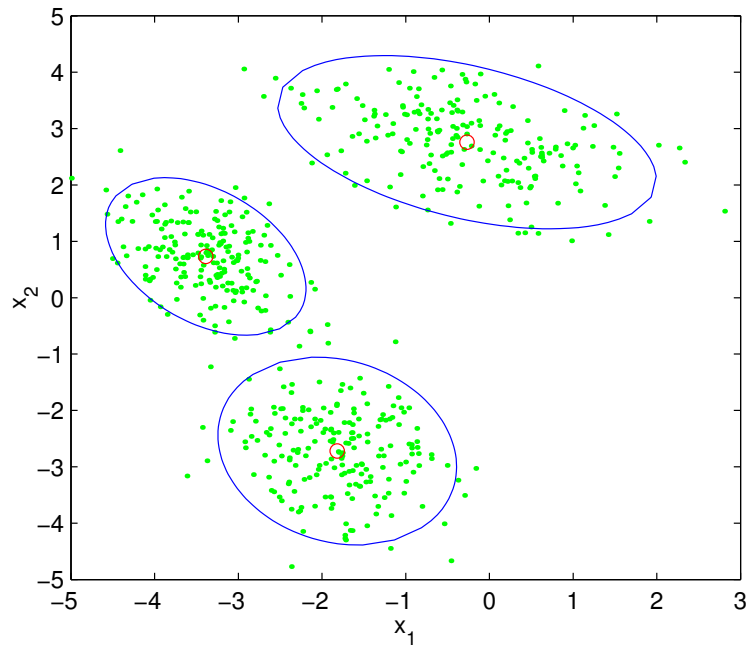


Figure 3: Fitting with 3 Gaussians

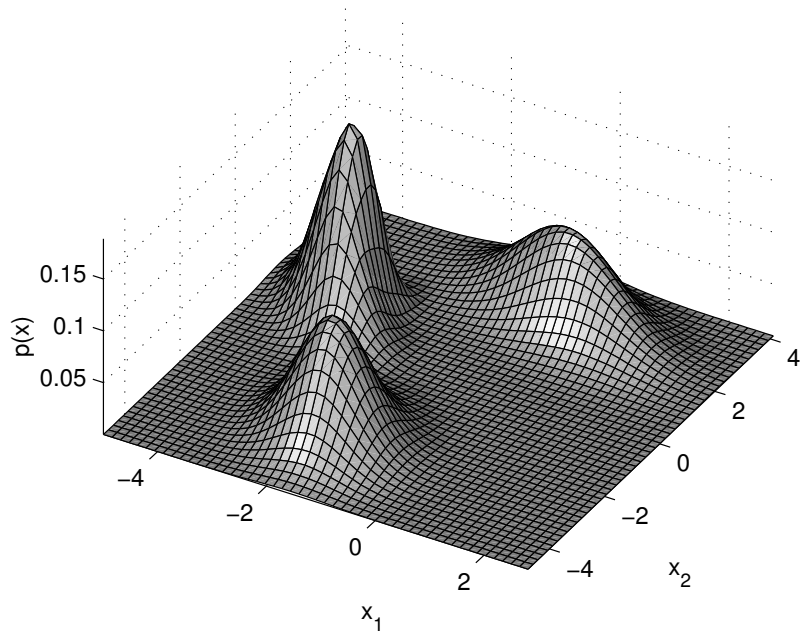
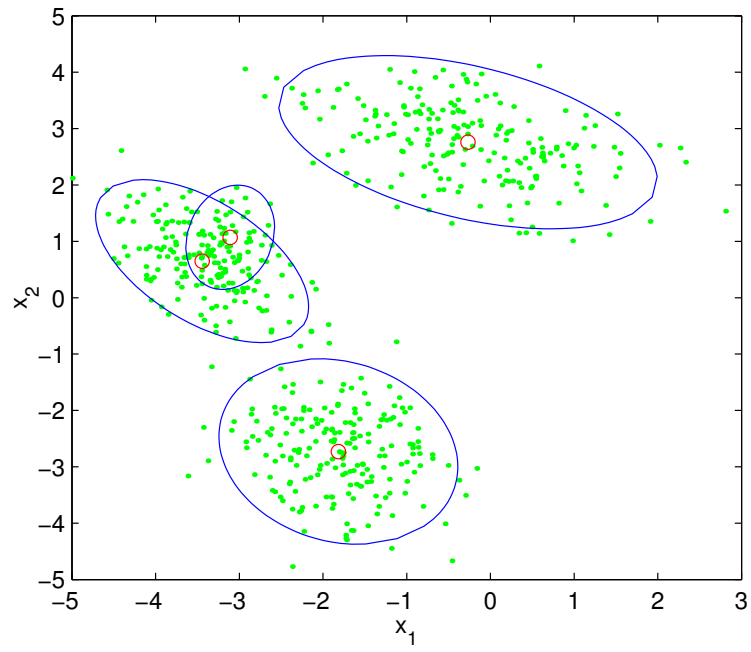


Figure 4: Fitting with 4 Gaussians

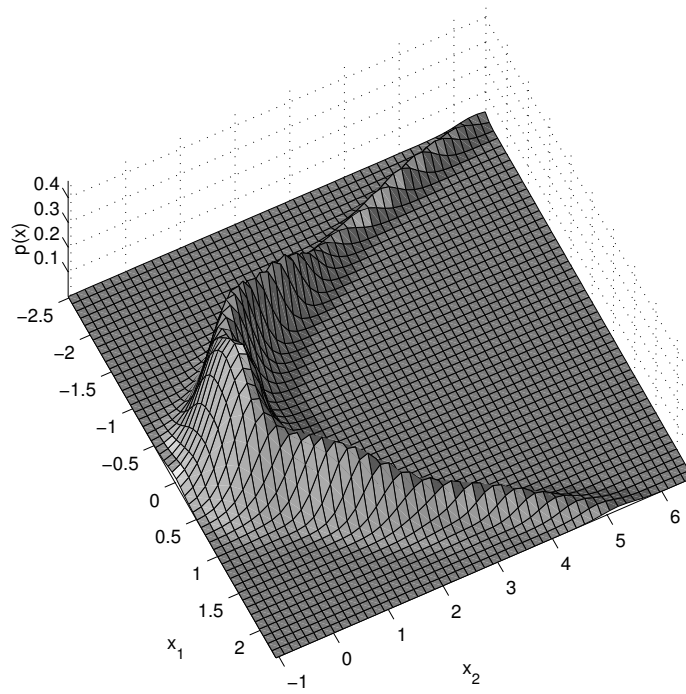
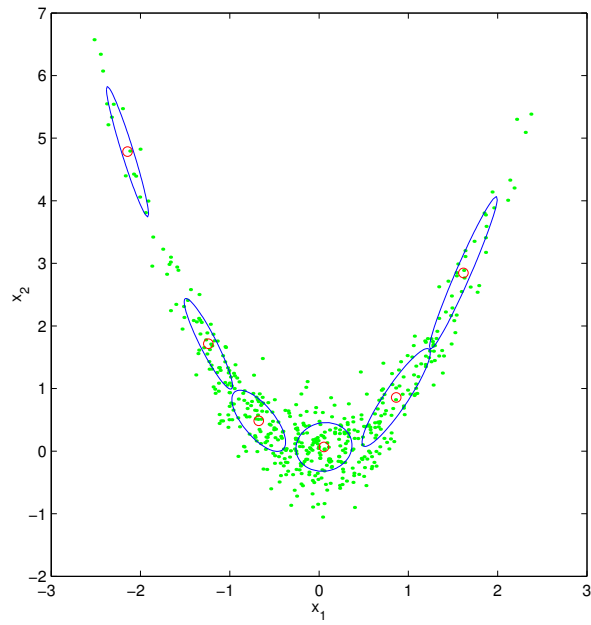


Figure 5: A sufficient number of mixture components can model arbitrary distributions