

Parameter Estimation Methods

Slides compiled by Sanghamitra De

Introduction

- In pattern recognition applications complete knowledge about the probabilistic structure of the problem is rarely present.
- Typically there is some vague, general knowledge about the situation, together with a number of design samples or training data — particular representatives of the patterns we want to classify.
- The problem, is to find some way to use this information to design or train the classifier.
- One approach to this problem is to use the samples to estimate the unknown probabilities and probability densities, and to use the resulting estimates as if they were true values.
- In typical supervised pattern classification problems, the estimation of the prior probabilities presents no difficulties.

Introduction

- However, in estimation of the class-conditional densities, the number of available samples always seems too small, and serious problems arise when the dimensionality of the feature vector \mathbf{x} is large.
- If we know the number of parameters in advance and our general knowledge about the problem permits us to parameterize the conditional densities, then the severity of these problems can be reduced significantly.
- Assume that $p(\mathbf{x}|\omega_i)$ is a normal density with mean μ_i and covariance matrix Σ_i , although we don't know the exact value of these quantities.
- So now we need to estimate values of μ_i and Σ_i
- Problem of parameter estimation is a classical one in statistics, and it can be approached in several ways.
- Two common and reasonable procedures are *maximum likelihood estimation* (MLE) and *Bayesian estimation*.
- The two methods give identical results but are conceptually different.

Introduction

- Maximum likelihood and several other methods view the parameters as quantities whose values are fixed but unknown.
- The best estimate of their value is defined to be the one that maximizes the probability of obtaining the samples actually observed.
- Bayesian methods view the parameters as random variables having some known a priori distribution.
- Observation of the samples converts this to a posterior density, thereby revising our opinion about the true values of the parameters.
- In the Bayesian case, we shall see that a typical effect of observing additional samples is to sharpen the a posteriori density function, causing it to peak near the true values of the parameters. This phenomenon is known as *Bayesian learning*.

Maximum Likelihood Estimation

➤ Advantages:

- ✓ they nearly always have good convergence properties as the number of training samples increases.
- ✓ maximum likelihood estimation often can be simpler than alternate methods, such as Bayesian techniques or other methods

Maximum Likelihood Parameter Estimation: General Principal

- Suppose we separate a collection of samples according to class, so that there are c data sets, D_1, \dots, D_c , with the samples in D_j having been drawn independently as per the probability law $p(x|\omega_j)$
 - Such samples are i.i.d. – independent and identically distributed random variables.
 - Assume that $p(x|\omega_j)$ has a known parametric form, and so is determined uniquely by the value of a parameter Θ_j
 - Find parameters that maximize probability of observations
-
- Assume c classes and
 - ✓ $p(x | \omega_j) \sim N(\mu_j, \Sigma_j)$
 - ✓ $p(x | \omega_j) \equiv p(x | \omega_j, \theta_j)$ where:

$$\theta = (\mu_j, \Sigma_j)$$

Maximum Likelihood Estimation

- Use n training samples in a class to estimate θ
- If \mathcal{D} contains n independently drawn samples, x_1, x_2, \dots, x_n

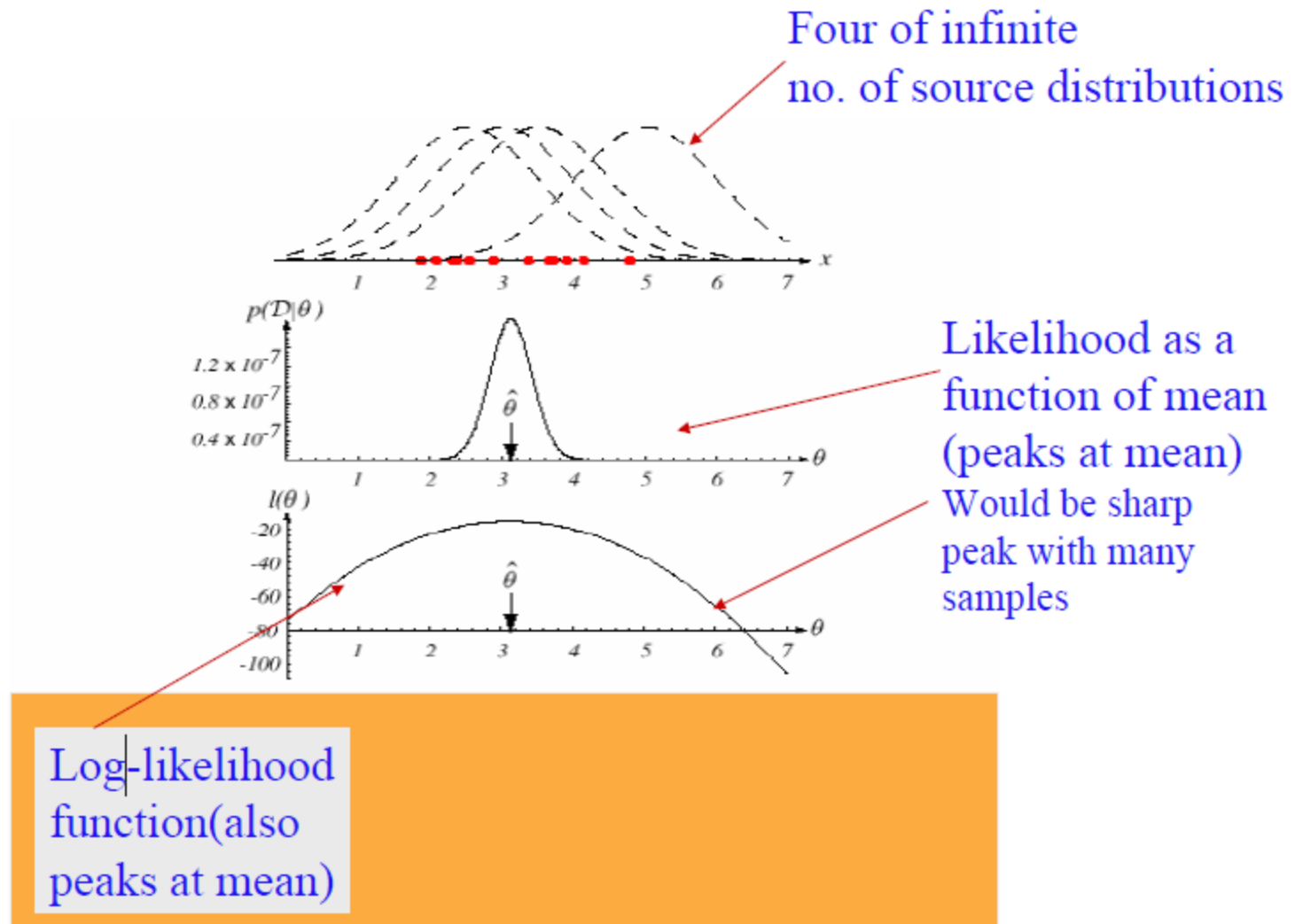
$$p(\mathcal{D} | \theta) = \prod_{k=1}^n p(x_k | \theta)$$

$p(\mathcal{D} | \theta)$ is called the likelihood of θ w.r.t. the set of samples

$\ell(\theta)$ is the log-likelihood of θ

- $l(\Theta) \equiv \ln p(\mathcal{D} | \theta)$
- ML estimate of θ is, by definition the value that maximizes $p(\mathcal{D} | \theta)$
- “It is the value of θ that best agrees with the actually observed training samples”

One-Dimensional Example



Maximizing the log-likelihood function

- Let $\theta = (\theta_1, \theta_2, \dots, \theta_p)^t$ and let ∇_θ be the gradient operator

$$\nabla_\theta = \left[\frac{\partial}{\partial \theta_1}, \frac{\partial}{\partial \theta_2}, \dots, \frac{\partial}{\partial \theta_p} \right]^t$$

- We define $l(\theta)$ as the log-likelihood function

$$l(\theta) = \ln p(D | \theta)$$

- Determine θ that maximizes the log-likelihood $\hat{\theta} = \arg \max_{\theta} \ell(\theta)$

- Set of necessary conditions for an optimum is:

$$\nabla_\theta l = \sum_{k=1}^n \nabla_\theta \ln p(x_k | \theta) = 0$$

MLE: The Gaussian Case: unknown μ

$$p(x_i | \mu) \sim N(\mu, \Sigma)$$

$$\ln p(x_k | \mu) = -\frac{1}{2} \ln[(2\pi)^d |\Sigma|] - \frac{1}{2} (x_k - \mu)^t \Sigma^{-1} (x_k - \mu)$$

and $\nabla_{\mu} \ln P(x_k | \mu) = \Sigma^{-1} (x_k - \mu)$

$\theta = \mu$ therefore:

The ML estimate for μ must satisfy:

$$\sum_{k=1}^n \Sigma^{-1} (x_k - \hat{\mu}) = 0$$

Multiplying by Σ :

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$$

Just the sample mean

MLE: Gaussian Case - *unknown* μ and Σ

$$\theta = (\theta_1, \theta_2) = (\mu, \sigma^2)$$

Consider the univariate case where:

$$\Theta_1 = \mu \text{ and } \Theta_2 = \sigma^2$$

$$\ln p(\mathbf{x}_k | \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2}(\mathbf{x}_k - \theta_1)^2$$

$$\nabla_{\theta} \mathbf{l} = \begin{pmatrix} \frac{\partial}{\partial \theta_1} (\ln p(\mathbf{x}_k | \theta)) \\ \frac{\partial}{\partial \theta_2} (\ln p(\mathbf{x}_k | \theta)) \end{pmatrix} = 0$$

$$\begin{cases} \frac{1}{\theta_2}(\mathbf{x}_k - \theta_1) = 0 \\ -\frac{1}{2\theta_2} + \frac{(\mathbf{x}_k - \theta_1)^2}{2\theta_2^2} = 0 \end{cases}$$

$$\begin{cases} \sum_{k=1}^n \frac{1}{\hat{\theta}_2} (x_k - \theta_1) = 0 & (1) \\ -\sum_{k=1}^n \frac{1}{\hat{\theta}_2} + \sum_{k=1}^n \frac{(x_k - \hat{\theta}_1)^2}{\hat{\theta}_2^2} = 0 & (2) \end{cases}$$

$$\mu = \sum_{k=1}^n \frac{x_k}{n} \quad ; \quad \sigma^2 = \frac{\sum_{k=1}^n (x_k - \mu)^2}{n}$$

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \hat{\boldsymbol{\mu}})(\mathbf{x}_k - \hat{\boldsymbol{\mu}})^t.$$

For the multivariate case: MLE for the mean vector is the sample mean. MLE for the covariance matrix is the arithmetic average of the n matrices $(\mathbf{x}_k - \hat{\boldsymbol{\mu}})(\mathbf{x}_k - \hat{\boldsymbol{\mu}})^t$

MLE Bias

- ML estimate for σ^2 is biased (i.e. expected value of the sample variance is not equal to the true variance)

$$\mathcal{E} \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right] = \frac{n-1}{n} \sigma^2 \neq \sigma^2.$$

- An elementary unbiased estimator for Σ is:

$$\mathbf{C} = \frac{1}{n-1} \sum_{k=1}^n (\mathbf{x}_k - \hat{\boldsymbol{\mu}})(\mathbf{x}_k - \hat{\boldsymbol{\mu}})^t,$$

MLE Bias

- Here \mathbf{C} is the sample covariance matrix.
- If an estimator is unbiased for all distributions, then it is called *absolutely unbiased*.
- If the estimator tends to become unbiased as the number of samples become very large, then the estimator is *asymptotically unbiased*.

Bayesian estimation



Bayesian Estimation

- Although the answers we get by Bayesian estimation method will generally be nearly identical to those obtained by maximum likelihood, there is a conceptual difference:
 - ✓ In MLE θ , was supposed fix
 - ✓ In BE θ , is a random variable
and training data allows us to convert a distribution on this variable into a posterior probability density.

Bayesian Estimation: the class-conditional densities

- The computation of posterior probabilities $P(\omega_i | \mathbf{x})$ lies at the heart of Bayesian classification
- Bayes' formula allows us to compute these probabilities from the Prior probabilities $P(\omega_i)$ and the class-conditional densities $p(\mathbf{x}|\omega_i)$
- The best we can do is to compute $P(\omega_i|\mathbf{x})$ using all of the information at our disposal.
- Part of this information might be prior knowledge, such as knowledge of the functional forms for unknown densities and ranges for the values of unknown parameters.
- Part of this information might reside in a set of training samples.

Bayesian Estimation : the class-conditional densities

- If we again let D denote the set of samples, then we can emphasize the role of the samples by saying that our Goal is:
compute $P(\omega_i | \mathbf{x}, D)$
- Given the sample D , Bayes formula can be written

$$P(\omega_i | \mathbf{x}, D) = \frac{p(\mathbf{x} | \omega_i, D)P(\omega_i | D)}{\sum_{j=1}^c p(\mathbf{x} | \omega_j, D)P(\omega_j | D)}$$

Bayesian Estimation : the class-conditional densities

- Now, assume that true values of prior probabilities are known.
- So, we substitute $P(\omega_i) = P(\omega_i | D)$.
- Also, as we are treating the supervised case, we can separate the training samples by class into c subsets D_1, \dots, D_c , with the samples in D_i belonging to ω_i .
- Now, the samples in D_i have no influence on $p(x | \omega_j, D)$ if $i \neq j$.
- So, we can work with each class separately, using only the samples in D_i to determine $p(x | \omega_i, D)$
- In conjunction to the assumption that prior probabilities are known, this gives us the equation:

Bayesian Estimation : the class-conditional densities

$$P(\omega_i|\mathbf{x}, \mathcal{D}) = \frac{p(\mathbf{x}|\omega_i, \mathcal{D}_i)P(\omega_i)}{\sum_{j=1}^c p(\mathbf{x}|\omega_j, \mathcal{D}_j)P(\omega_j)}.$$

So, we now have c separate problems of the following form: use a set \mathbf{D} of samples drawn independently according to the fixed but unknown probability distribution $p(\mathbf{x})$ to determine $\mathbf{p}(\mathbf{x}|\mathbf{D})$. This is the central problem of Bayesian learning.

Bayesian Estimation : the parameter distribution

- Although the desired probability density $\mathbf{p}(\mathbf{x})$ is unknown, we assume that it has a known parametric form. The only thing assumed unknown is the value of a parameter vector $\boldsymbol{\theta}$.
- Also we have the fact that $\mathbf{p}(\mathbf{x})$ is unknown but has known parametric form by saying that the function $\mathbf{p}(\mathbf{x}|\boldsymbol{\theta})$ is completely known.
- Any information that is there about $\boldsymbol{\theta}$ prior to observing the samples is assumed to be contained in a known prior density $\mathbf{p}(\boldsymbol{\theta})$.
- Observation of the samples converts this to a posterior density $\mathbf{p}(\boldsymbol{\theta}|\mathbf{D})$, which, is sharply peaked about the true value of $\boldsymbol{\theta}$.

Bayesian Estimation : the parameter distribution

- So, our supervised learning problem is converted to an unsupervised density estimation problem.
- To this end, our goal is to compute $\mathbf{p}(\mathbf{x}|\mathbf{D})$, which is as close as we can come to obtaining the unknown $\mathbf{p}(\mathbf{x})$.
- For this, we integrate the joint density $\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}|\mathbf{D})$ over $\boldsymbol{\theta}$.

$$p(\mathbf{x}|\mathcal{D}) = \int p(\mathbf{x}, \boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta},$$

where the integration extends over the entire parameter space.

- We can write $\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}|\mathbf{D})$ as the product $\mathbf{p}(\mathbf{x}|\boldsymbol{\theta}, \mathbf{D})\mathbf{p}(\boldsymbol{\theta}|\mathbf{D})$.
- Since the selection of \mathbf{x} and that of the training samples in \mathbf{D} is done independently, the first factor is merely $\mathbf{p}(\mathbf{x}|\boldsymbol{\theta})$. That is, the distribution of \mathbf{x} is known completely once we know the value of the parameter vector.
- So, above equation can be rewritten as:

$$p(\mathbf{x}|\mathcal{D}) = \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta}.$$

Bayesian Estimation : the parameter distribution

- This equation links the desired class-conditional density $\mathbf{p}(\mathbf{x}|\mathbf{D})$ to the posterior density $\mathbf{p}(\boldsymbol{\theta}|\mathbf{D})$ for the unknown parameter vector.
- If $\mathbf{p}(\boldsymbol{\theta}|\mathbf{D})$ peaks very sharply about some value $\hat{\boldsymbol{\theta}}$, we obtain $\mathbf{p}(\mathbf{x}|\mathbf{D}) \approx \mathbf{p}(\mathbf{x}|\hat{\boldsymbol{\theta}})$, i.e., the result we would obtain by substituting the estimate $\hat{\boldsymbol{\theta}}$ for the true parameter vector.
- This result rests on the assumption that $\mathbf{p}(\mathbf{x}|\boldsymbol{\theta})$ is smooth, and that the tails of the integral are not important.
- *These conditions are typically but not invariably the case.*
- In general, if we are less certain about the exact value of $\boldsymbol{\theta}$, this equation directs us to average $\mathbf{p}(\mathbf{x}|\boldsymbol{\theta})$ over the possible values of $\boldsymbol{\theta}$.
- Thus, when the unknown densities have a known parametric form, the samples exert their influence on $\mathbf{p}(\mathbf{x}|\mathbf{D})$ through the posterior density $\mathbf{p}(\boldsymbol{\theta}|\mathbf{D})$.

Bayesian Parameter Estimation: Gaussian Case

Goal: Estimate θ using the a-posteriori density $p(\theta|D)$ & the desired probability density $p(x|D)$ for the case where $p(\mathbf{x}|\boldsymbol{\mu}) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

➤ **The univariate case: $p(\mu | D)$**

μ is the only unknown parameter

$$p(x | \mu) \sim N(\mu, \sigma^2)$$

$$p(\mu) \sim N(\mu_0, \sigma_0^2)$$

(μ_0 and σ_0 are known!)

Bayesian Parameter Estimation: Gaussian Case

$$p(\mu | \mathbf{D}) = \frac{p(\mathbf{D} | \mu) p(\mu)}{\int p(\mathbf{D} | \mu) p(\mu) d\mu} \quad \text{Bayes Formula}$$
$$= \alpha \prod_{k=1}^{k=n} p(x_k | \mu) p(\mu)$$

- α is the normalization factor dependent on \mathbf{D} but independent of μ
- We know: $p(x_k | \mu) \sim N(\mu, \sigma^2)$ and $p(\mu) \sim N(\mu_0, \sigma_0^2)$
- Plugging in their gaussian expressions, extracting out factors which do not depend on μ (and getting them absorbed into constant α^n) yields:

$$p(\mu | D) = \alpha^n \exp \left(-\frac{1}{2} \left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left(\frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right) \mu \right] \right)$$

(from eq. 29 page 93)

Bayesian Parameter Estimation: Gaussian Case

- Observation: $p(\mu|D)$ is an exponential of a quadratic

$$p(\mu | D) = \alpha^n \exp\left(-\frac{1}{2}\left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right)\mu^2 - 2\left(\frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2}\right)\mu\right]\right)$$

- It is normal! It is called a *reproducing density*

$$p(\mu | D) \sim N(\mu_n, \sigma_n^2)$$

$p(\mu)$ is said to be a *conjugate prior*

Bayesian Parameter Estimation: Gaussian Case

$$p(\mu | D) = \alpha^n \exp \left(-\frac{1}{2} \left[\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left(\frac{1}{\sigma^2} \sum_{k=1}^n x_k + \frac{\mu_0}{\sigma_0^2} \right) \mu \right] \right)$$

Identifying coefficients in the top equation with that of the generic Gaussian

$$p(\mu | D) = \frac{1}{\sqrt{2\pi}\sigma_n} \exp \left(-\frac{1}{2} \left[\frac{\mu - \mu_n}{\sigma_n} \right]^2 \right)$$

Yields expressions for μ_n and σ_n^2

$$\frac{1}{\sigma_n^2} = \frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \quad \text{and} \quad \frac{\mu_n}{\sigma_n^2} = \frac{n}{\sigma^2} \hat{\mu}_n + \frac{\mu_0}{\sigma_0^2}$$

Bayesian Parameter Estimation: Gaussian Case

Solving for μ_n and σ_n^2 yields:

$$\mu_n = \left(\frac{n\sigma_0^2}{n_0\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \mu_0$$
$$\text{and } \sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

From these equations we see as n increases:

- the variance decreases monotonically
- the estimate of $p(\mu|D)$ becomes more peaked, approaching a Dirac delta function as n approaches infinity. This behavior is called *Bayesian Learning*.

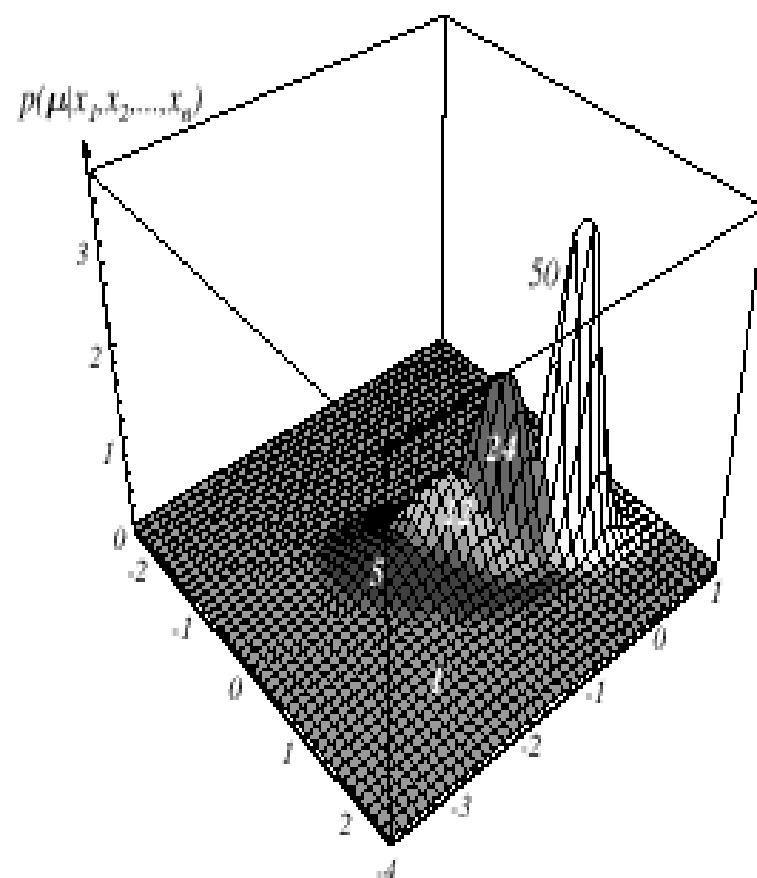
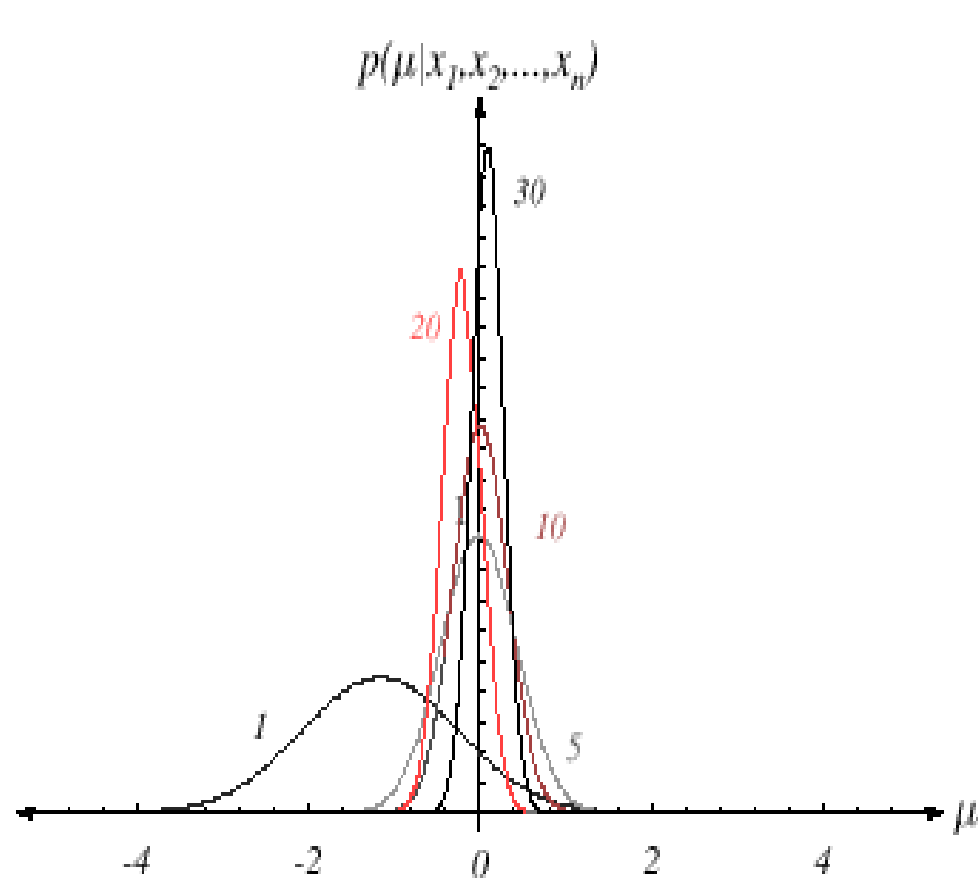


FIGURE 3.2. Bayesian learning of the mean of normal distributions in one and two dimensions. The posterior distribution estimates are labeled by the number of training samples used in the estimation. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Bayesian Parameter Estimation: Gaussian Case

➤ The univariate case $P(\mathbf{x} \mid \mathbf{D})$

- ✓ $p(\mu \mid \mathbf{D})$ computed (in preceding discussion)
- ✓ $p(\mathbf{x} \mid \mathbf{D})$ remains to be computed!

$$p(x \mid \mathbf{D}) = \int p(x \mid \mu) p(\mu \mid \mathbf{D}) d\mu \text{ is Gaussian}$$

It provides: $p(x \mid \mathbf{D}) \sim N(\mu_n, \sigma^2 + \sigma_n^2)$

We know σ^2 and how to compute μ_n and σ_n^2

Density $p(x \mid D)$ is the desired class-conditional density

$p(x \mid \omega_j, D_j)$ and together with prior probabilities $P(\omega_j)$ we obtain the Bayesian classification rule to design the classifier

$$\underset{\omega_j}{\text{Max}} [P(\omega_j \mid x, \mathbf{D})] \equiv \underset{\omega_j}{\text{Max}} [P(x \mid \omega_j, \mathbf{D}_j) P(\omega_j)]$$

Bayesian Parameter Estimation: Gaussian Case

- **The multivariate case:** (Σ is known but μ is not a direct generalization of the univariate case)
- As before assume: $p(x_k | \mu) \sim N(\mu, \Sigma)$ and $p(\mu) \sim N(\mu_0, \Sigma_0)$

Where Σ , Σ_0 and μ_0 are assumed to be known.

After observing a set \mathcal{D} of n independent samples x_1, \dots, x_n , Bayes formula is used to get:

$$\begin{aligned} p(\mu | \mathcal{D}) &= \alpha \prod_{k=1}^n p(x_k | \mu) p(\mu) \\ &= \alpha' \exp \left[-\frac{1}{2} \left(\mu^t (n\Sigma^{-1} + \Sigma_0^{-1}) \mu - 2\mu^t \left(\Sigma^{-1} \sum_{k=1}^n x_k + \Sigma_0^{-1} \mu_0 \right) \right) \right], \end{aligned} \quad (40)$$

Which has the form:

$$p(\mu | \mathcal{D}) = \alpha'' \exp \left[-\frac{1}{2} (\mu - \mu_n)^t \Sigma_n^{-1} (\mu - \mu_n) \right]$$

Bayesian Parameter Estimation: Gaussian Case

- So, $p(\mu|D) \sim N(\mu_n, \Sigma_n)$ and we have a reproducing density.
- Equating the coefficients, we get analogs of eqns. in slide 28 as:

$$\Sigma_n^{-1} = n\Sigma^{-1} + \Sigma_0^{-1} \quad \& \quad \Sigma_n^{-1}\mu_n = n\Sigma^{-1}\hat{\mu}_n + \Sigma_0^{-1}\mu_0,$$

where $\hat{\mu}_n$ is the sample mean

$$\hat{\mu}_n = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k.$$

Bayesian Parameter Estimation: Gaussian Case

➤ With little manipulation, we get:

$$\mu_n = \Sigma_0 \left(\Sigma_0 + \frac{1}{n} \Sigma \right)^{-1} \hat{\mu}_n + \frac{1}{n} \Sigma \left(\Sigma_0 + \frac{1}{n} \Sigma \right)^{-1} \mu_0$$

(which, as in the univariate case, is a linear combination of $\hat{\mu}_n$ and μ_0) and

$$\Sigma_n = \Sigma_0 \left(\Sigma_0 + \frac{1}{n} \Sigma \right)^{-1} \frac{1}{n} \Sigma.$$

The proof that $p(\mathbf{x}|\mathcal{D}) \sim N(\mu_n, \Sigma + \Sigma_n)$ can be obtained as before by performing the integration

$$p(\mathbf{x}|\mathcal{D}) = \int p(\mathbf{x}|\mu)p(\mu|\mathcal{D}) d\mu. \quad (48)$$

Bayesian Parameter Estimation: Gaussian Case

- Since the sum of two independent, normally distributed vectors is again a normally distributed vector whose mean is the sum of the means and whose covariance matrix is the sum of the covariance matrices, we have:

$$p(\mathbf{x}|\mathcal{D}) \sim N(\mu_n, \Sigma + \Sigma_n),$$

and the generalization is complete

Bayesian Parameter Estimation: General Theory

- $P(x \mid D)$ computation can be applied to any situation in which the unknown density can be parameterized: the basic assumptions are:
 - ✓ The form of $P(x \mid \theta)$ is assumed known, but the value of θ is not known exactly
 - ✓ Our knowledge about θ is assumed to be contained in a known prior density $P(\theta)$
 - ✓ The rest of our knowledge θ is contained in a set D of n random variables x_1, x_2, \dots, x_n that follows $P(x)$

Bayesian Parameter Estimation: General Theory

The basic problem is:

“Compute the posterior density $p(\boldsymbol{\theta} \mid D)$ ” then “Derive $p(\mathbf{x} \mid D)$ ”,

where
$$p(\mathbf{x} \mid D) = \int p(\mathbf{x} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid D) d\boldsymbol{\theta}$$

Using Bayes formula, we have:
$$p(\boldsymbol{\theta} \mid D) = \frac{p(D \mid \boldsymbol{\theta}) p(\boldsymbol{\theta})}{\int p(D \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}},$$

And by the independence assumption:

$$p(D \mid \boldsymbol{\theta}) = \prod_{k=1}^n p(\mathbf{x}_k \mid \boldsymbol{\theta})$$

➤ This constitutes the formal solution to the problem.

Problems of Dimensionality

- Problems involving 50 or 100 features are common (usually binary valued)
- Note: microarray data might entail ~20000 real-valued features
- Classification accuracy dependant on
 - ✓ Dimensionality (amount of training data)
 - ✓ discrete vs continuous

Problems of Dimensionality

- Case of two class multivariate normal with the same covariance
 - ✓ $p(\mathbf{x}|\omega_j) \sim N(\mu_j, \Sigma), j=1,2$
 - ✓ Statistically independent features
 - ✓ If the priors are equal then:

$$P(error) = \frac{1}{\sqrt{2\pi}} \int_{r/2}^{\infty} e^{-u^2/2} du \quad (\text{Bayes error})$$

$$\text{where : } r^2 = (\mu_1 - \mu_2)^t \Sigma^{-1} (\mu_1 - \mu_2)$$

r^2 is the squared Mahalanobi s distance

$$\lim_{r \rightarrow \infty} P(error) = 0$$

Problems of Dimensionality

- If features are *conditionally* independent then:

$$\Sigma = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2)$$

$$r^2 = \sum_{i=1}^{i=d} \left(\frac{\mu_{i1} - \mu_{i2}}{\sigma_i} \right)^2$$

- ✓ Do we remember what conditional independence is?

Problems of Dimensionality

- Most useful features are the ones for which the difference between the means is large relative to the standard deviation
 - ✓ Doesn't require independence
- Adding independent features helps increase $r \rightarrow$ reduce error
- Caution: adding features increases cost & complexity of feature extractor and classifier
- It has frequently been observed in practice that, beyond a certain point, the inclusion of additional features leads to worse rather than better performance:
 - ✓ we have the wrong model !
 - ✓ we don't have enough training data to support the additional dimensions

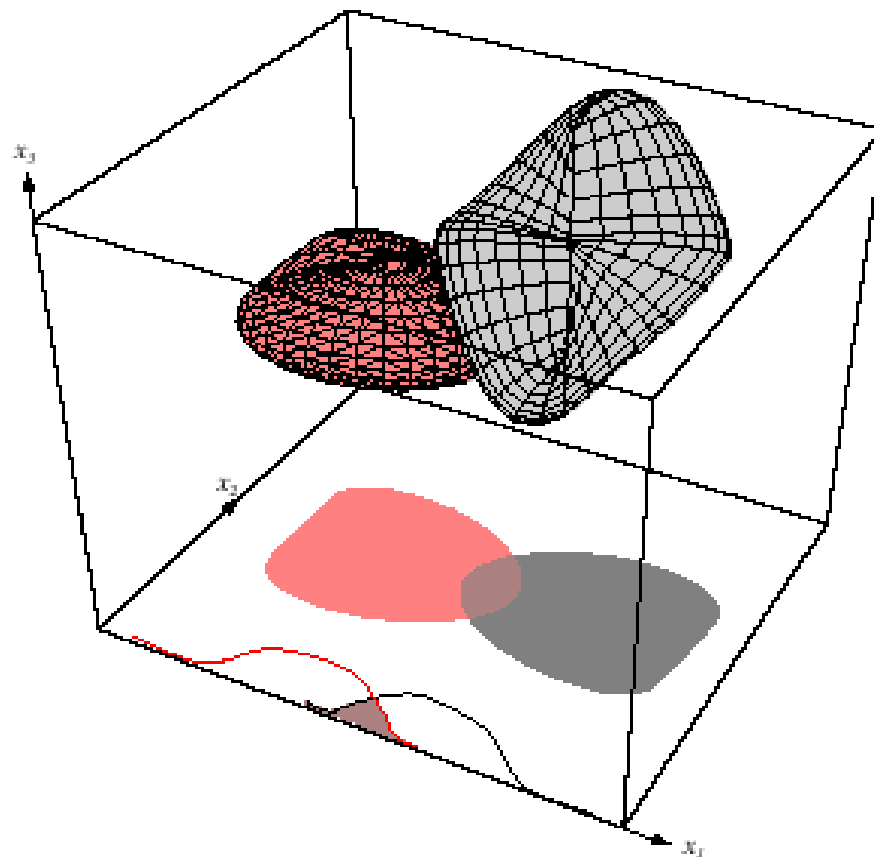


FIGURE 3.3. Two three-dimensional distributions have nonoverlapping densities, and thus in three dimensions the Bayes error vanishes. When projected to a subspace—here, the two-dimensional $x_1 - x_2$ subspace or a one-dimensional x_1 subspace—there can be greater overlap of the projected distributions, and hence greater Bayes error. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Computational Complexity

- Our design methodology is affected by the computational difficulty

- “big oh” notation

$f(x) = O(h(x))$ “big oh of $h(x)$ ”

If: there exist constants c_0 and x_0 such that

$$|f(x)| \leq c_0 |h(x)| \text{ for all } x > x_0.$$

(An upper bound on $f(x)$ grows no worse than $h(x)$ for sufficiently large x !)

Example:

$$f(x) = 2 + 3x + 4x^2$$

$$g(x) = x^2$$

$$f(x) = O(x^2)$$

Computational Complexity

➤ “big oh” is not unique!

$$f(x) = O(x^2); f(x) = O(x^3); f(x) = O(x^4)$$

➤ “big theta” notation

$$f(x) = \theta(h(x))$$

If: there are constants x_0 , c_1 and c_2 such that
for $x > x_0$,

$f(x)$ always lies between $c_1 h(x)$ and $c_2 h(x)$.

So, $f(x) = \theta(x^2)$ but $f(x) \neq \theta(x^3)$

Computational Complexity

➤ Complexity of the ML Estimation

- ✓ Gaussian priors in d dimensions classifier with n training samples for each of c classes
- ✓ For each category, we have to compute the discriminant function

$$g(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \overset{O(dn)}{\hat{\boldsymbol{\mu}}})^t \overset{O(d^2n)}{\hat{\boldsymbol{\Sigma}}^{-1}}(\mathbf{x} - \hat{\boldsymbol{\mu}}) - \overset{O(1)}{\frac{d}{2} \ln 2\pi} - \underbrace{\frac{1}{2} \ln |\hat{\boldsymbol{\Sigma}}|}_{O(d^2n)} + \underbrace{\ln P(\omega)}_{O(n)}$$

Total = $O(d^2n)$

Total for c classes = $O(cd^2n) \cong O(d^2n)$

- Cost increase when d and n are large!

Overfitting

- Dimensionality of model vs size of training data
 - ✓ **Issue: not enough data to support the model**
 - ✓ Possible solutions:
 - ❑ Reduce model dimensionality either by redesigning the feature extractor, by selecting an appropriate subset of the existing features, or by combining the existing features in some way
 - ❑ Make (possibly incorrect) assumptions to better estimate Σ

Overfitting

- Estimate better Σ
 - use data pooled from all classes
 - ✓ normalization issues
 - use pseudo-Bayesian form $\lambda \Sigma_0 + (1-\lambda) \Sigma_n$ where Σ_0 is a prior estimate
 - by thresholding sample covariance matrix Σ
 - ✓ reduces chance correlations
 - assume statistical independence
 - ✓ zero all off-diagonal elements

Shrinkage

- **Issue: not enough data to support the model**
- Shrinkage: weighted combination of common and individual covariances. If i is an index on the c categories in question, then:

$$\Sigma_i(\alpha) = \frac{(1-\alpha)n_i\Sigma_i + \alpha n\Sigma}{(1-\alpha)n_i + \alpha n} \quad \text{for } 0 < \alpha < 1$$

- We can also shrink the estimate of the common covariances toward the identity matrix

$$\Sigma(\beta) = (1-\beta)\Sigma + \beta\mathbf{I} \quad \text{for } 0 < \beta < 1$$

Gaussian Mixture Model (GMM)

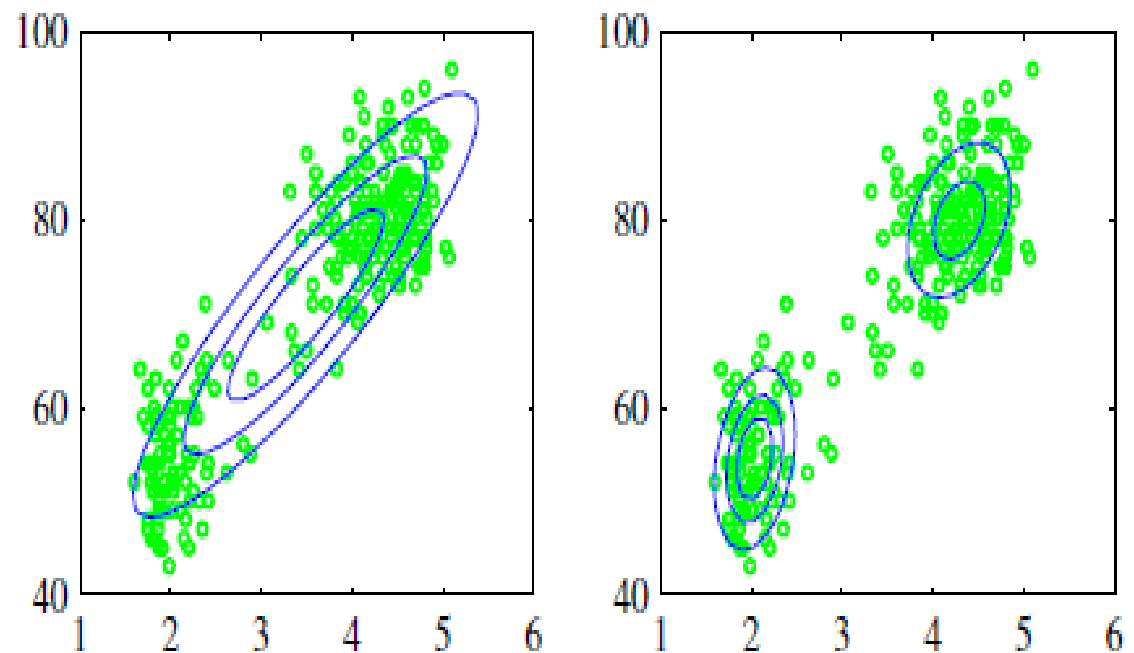
- A **Gaussian Mixture Model (GMM)** is a parametric probability density function represented as a weighted sum of Gaussian component densities.
- GMMs are commonly used as a parametric model of the probability distribution of continuous measurements or features in a biometric system, such as vocal-tract related spectral features in a speaker recognition system.
- GMM parameters are estimated from training data using the iterative **Expectation-Maximization (EM)** algorithm or **Maximum A Posteriori (MAP)** estimation from a well-trained prior model.

GMM

- Gaussian distribution has significant limitations when it comes to modelling real data sets.
- Consider the example shown in next slide:
 - ✓ This is known as the ‘Old Faithful’ data set, and comprises 272 measurements of the eruption of the Old Faithful geyser at Yellowstone National Park in the USA.
 - ✓ Each measurement comprises the duration of the eruption in minutes (horizontal axis) and the time in minutes to the next eruption (vertical axis).
 - ✓ The data set forms two dominant clumps, and that a simple Gaussian distribution is unable to capture this structure, whereas a linear superposition of two Gaussians gives a better characterization of the data set.

GMM

Figure 2.21 Plots of the 'old faithful' data in which the blue curves show contours of constant probability density. On the left is a single Gaussian distribution which has been fitted to the data using maximum likelihood. Note that this distribution fails to capture the two clumps in the data and indeed places much of its probability mass in the central region between the clumps where the data are relatively sparse. On the right the distribution is given by a linear combination of two Gaussians which has been fitted to the data by maximum likelihood using techniques discussed Chapter 9, and which gives a better representation of the data.

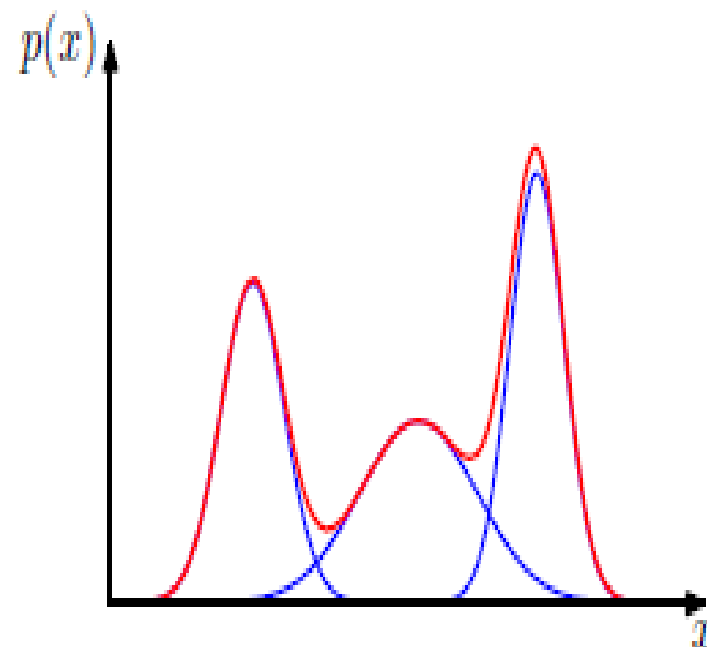


GMM

- Such superpositions, formed by taking linear combinations of more basic distributions such as Gaussians, can be formulated as probabilistic models known as *mixture distributions*.
- In the next figure it can be seen that a linear combination of Gaussians can give rise to very complex densities.
- By using a sufficient number of Gaussians, and by adjusting their means and covariances as well as the coefficients in the linear combination, almost any continuous density can be approximated to arbitrary accuracy.

GMM

Figure 2.22 Example of a Gaussian mixture distribution in one dimension showing three Gaussians (each scaled by a coefficient) in blue and their sum in red.



GMM

- Therefore consider a superposition of K Gaussian densities of the form:

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

- The above is called a *mixture of Gaussians*.
- Each Gaussian density $\mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)$ is called a *component* of the mixture and has its own mean μ_k and covariance Σ_k .
- Contour and surface plots for a Gaussian mixture having 3 components are shown in next figure
- Generally speaking, mixture models can comprise linear combinations of other distributions.

GMM

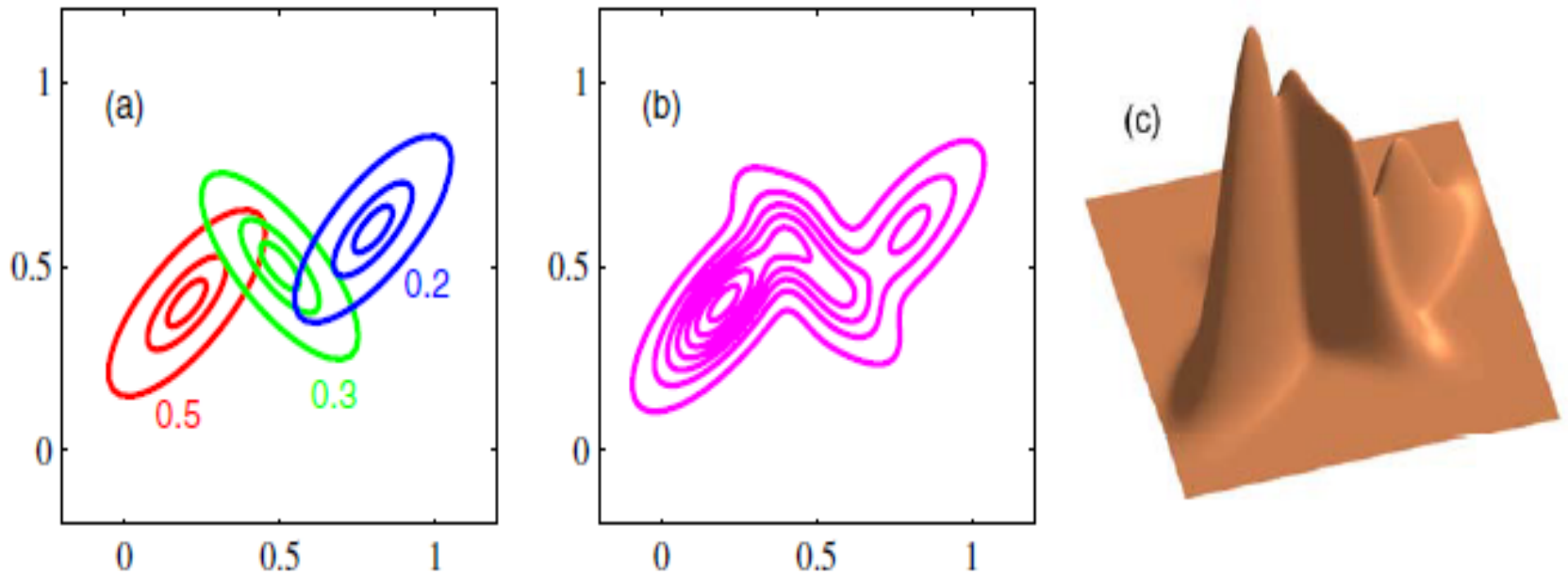


Figure 2.23 Illustration of a mixture of 3 Gaussians in a two-dimensional space. (a) Contours of constant density for each of the mixture components, in which the 3 components are denoted red, blue and green, and the values of the mixing coefficients are shown below each component. (b) Contours of the marginal probability density $p(\mathbf{x})$ of the mixture distribution. (c) A surface plot of the distribution $p(\mathbf{x})$.

GMM

- The parameters π_k in the equation (see previous slide) are called *mixing coefficients*.
- If we integrate both sides of the equation with respect to x , and note that both $p(x)$ and the individual Gaussian components are normalized, we obtain:

$$\sum_{k=1}^K \pi_k = 1.$$

- Also, the requirement that $p(x) \geq 0$, together with $N(x|\mu_k, \Sigma_k) \geq 0$, implies $\pi_k \geq 0$ for all k . Combining this with the condition above, we obtain:

$$0 \leq \pi_k \leq 1.$$

GMM

- Therefore, the mixing coefficients satisfy the requirements to be probabilities.
- From the sum and product rules, the marginal density is given by:

$$p(\mathbf{x}) = \sum_{k=1}^K p(k)p(\mathbf{x}|k)$$

- This is equivalent to eqn. in slide 53 in which we can view $\pi_k = p(k)$ as the prior probability of picking the k th component, and the density $N(\mathbf{x}|\mu_k, \Sigma_k) = p(\mathbf{x}|k)$ as the probability of \mathbf{x} conditioned on k .
- NOTE: an important role is played by the posterior probabilities $p(k|\mathbf{x})$, which are also known as *responsibilities*.

GMM

- From Bayes' theorem these are given by:

$$\begin{aligned}\gamma_k(\mathbf{x}) &\equiv p(k|\mathbf{x}) \\ &= \frac{p(k)p(\mathbf{x}|k)}{\sum_l p(l)p(\mathbf{x}|l)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)}{\sum_l \pi_l \mathcal{N}(\mathbf{x}|\mu_l, \Sigma_l)}.\end{aligned}$$

- The form of the Gaussian mixture distribution is governed by the parameters π , μ and Σ , where we have used the notation $\pi \equiv \{\pi_1, \dots, \pi_K\}$, $\mu \equiv \{\mu_1, \dots, \mu_K\}$ and $\Sigma \equiv \{\Sigma_1, \dots, \Sigma_K\}$.
- One way to set the values of these parameters is to use maximum likelihood.

GMM

- From eqn. in slide 53 the log of the likelihood function is given by

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

where $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.

- This situation is now more complex than with a single Gaussian, due to the presence of the summation over k inside the logarithm.
- As a result, the maximum likelihood solution for the parameters no longer has a closed-form analytical solution.
- One approach to maximizing the likelihood function is to use iterative numerical optimization techniques.
- Other alternative is to use *expectation maximization* (EM) algorithm.

Mixtures of Gaussians

- Recall that the Gaussian mixture distribution can be written as a linear superposition of Gaussians in the form

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

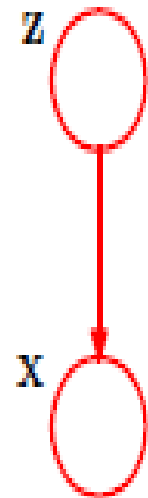
- We now turn to a formulation of Gaussian mixtures in terms of discrete latent variables.
- This will provide us with a deeper insight into this important distribution, and will serve to motivate the expectation-maximization algorithm.

Mixtures of Gaussians

- Let us introduce a K -dimensional binary random variable z having a 1-of- K representation in which a particular element z_k is equal to 1 and all other elements are equal to 0.
- The values of z_k therefore satisfy $z_k \in \{0, 1\}$ and $\sum_k z_k = 1$
- It can be seen that there are K possible states for the vector z according to which element is nonzero.
- We shall define the joint distribution $p(x, z)$ in terms of a marginal distribution $p(z)$ and a conditional distribution $p(x|z)$, corresponding to the graphical model in next figure.

Mixtures of Gaussians

Graphical representation of a mixture model, in which the joint distribution is expressed in the form $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$.



Mixtures of Gaussians

- The marginal distribution over \mathbf{z} is specified in terms of the mixing coefficients π_k , such that: $p(z_k = 1) = \pi_k$

where the parameters $\{\pi_k\}$ must satisfy $0 \leq \pi_k \leq 1$ together with:

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

in order to be valid probabilities.

Since \mathbf{z} uses a 1-of- K representation, we can also write this distribution in the form:

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

Mixtures of Gaussians

- Similarly, the conditional distribution of \mathbf{x} given a particular value for \mathbf{z} is a Gaussian:

$$p(\mathbf{x}|\mathbf{z}_k = 1) = \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)$$

- This can also be written in the form:
$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)^{z_k}.$$

- The joint distribution is given by $p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$, and the marginal distribution of \mathbf{x} is then obtained by summing the joint distribution over all possible states of \mathbf{z} to give:

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

where we have made use of previous two equations.

Mixtures of Gaussians

- Thus the marginal distribution of \mathbf{x} is a Gaussian mixture of the form as in slide 59.
- If we have several observations $\mathbf{x}_1, \dots, \mathbf{x}_N$, then, because we have represented the marginal distribution in the form:

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z)$$

it follows that for every observed data point \mathbf{x}_n there is a corresponding latent variable z_n .

- Thus, an equivalent formulation of the Gaussian mixture involving an explicit latent variable is found.
- Now, we can work with the joint distribution $p(\mathbf{x}, \mathbf{z})$ instead of the marginal distribution $p(\mathbf{x})$.
- This will lead to significant simplifications, specially through the introduction of the expectation-maximization (EM) algorithm.

Mixtures of Gaussians

- Another quantity that will play an important role is the conditional probability of z given \mathbf{x} .
- We shall use $\gamma(z_k)$ to denote $p(z_k = 1|\mathbf{x})$, whose value can be found using Bayes' theorem:

$$\begin{aligned}\gamma(z_k) \equiv p(z_k = 1|\mathbf{x}) &= \frac{p(z_k = 1)p(\mathbf{x}|z_k = 1)}{\sum_{j=1}^K p(z_j = 1)p(\mathbf{x}|z_j = 1)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\mu_j, \Sigma_j)}.\end{aligned}$$

Mixtures of Gaussians

- Here, π_k can be viewed as the prior probability of $z_k = 1$, and the quantity $\gamma(z_k)$ as the corresponding posterior probability once we have observed x .
- $\gamma(z_k)$ can also be viewed as the *responsibility* that component k takes for ‘explaining’ the observation x .

Expectation Maximization

➤ The EM algorithm

- ✓ An excellent way of doing our unsupervised learning problem.
- ✓ Many, many other uses, including inference of Hidden Markov Models
- ✓ The expectation maximization (EM) algorithm, is a general technique for finding maximum likelihood solutions for probabilistic models having latent variables.

EM Algorithm

- Consider a probabilistic model in which we collectively denote all of the observed variables by X and all of the hidden variables by Z .
- The joint distribution $p(X, Z | \theta)$ is governed by a set of parameters denoted by θ .
- Our goal is to maximize the likelihood function that is given by:

$$p(X | \theta) = \sum_Z p(X, Z | \theta).$$

EM Algorithm

- Next we introduce a distribution $q(\mathbf{Z})$ defined over the latent variables.
- It is observed that, for any choice of $q(\mathbf{Z})$, the following decomposition holds:

$$\ln p(\mathbf{X}|\boldsymbol{\theta}) = \mathcal{L}(q, \boldsymbol{\theta}) + \text{KL}(q\|p)$$

where we have defined:

$$\mathcal{L}(q, \boldsymbol{\theta}) = \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z})} \right\}$$

$$\text{KL}(q\|p) = - \sum_{\mathbf{Z}} q(\mathbf{Z}) \ln \left\{ \frac{p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})}{q(\mathbf{Z})} \right\}.$$

EM Algorithm

- Note that $\mathcal{L}(q, \theta)$ is a functional of the distribution $q(Z)$, and a function of the parameters θ .
- The last two equations differ in sign and also that $\mathcal{L}(q, \theta)$ contains the joint distribution of X and Z while $\text{KL}(q||p)$ contains the conditional distribution of Z given X .
- To verify the 1st equation of slide 69, we first make use of the product rule of probability to give:

$$\ln p(\mathbf{X}, \mathbf{Z}|\theta) = \ln p(\mathbf{Z}|\mathbf{X}, \theta) + \ln p(\mathbf{X}|\theta)$$

which we then substitute into the expression for $\mathcal{L}(q, \theta)$.

EM Algorithm

- This gives rise to two terms, one of which cancels $KL(q||p)$ while the other gives the required log likelihood $\ln p(X|\theta)$ after noting that $q(Z)$ is a normalized distribution that sums to 1.

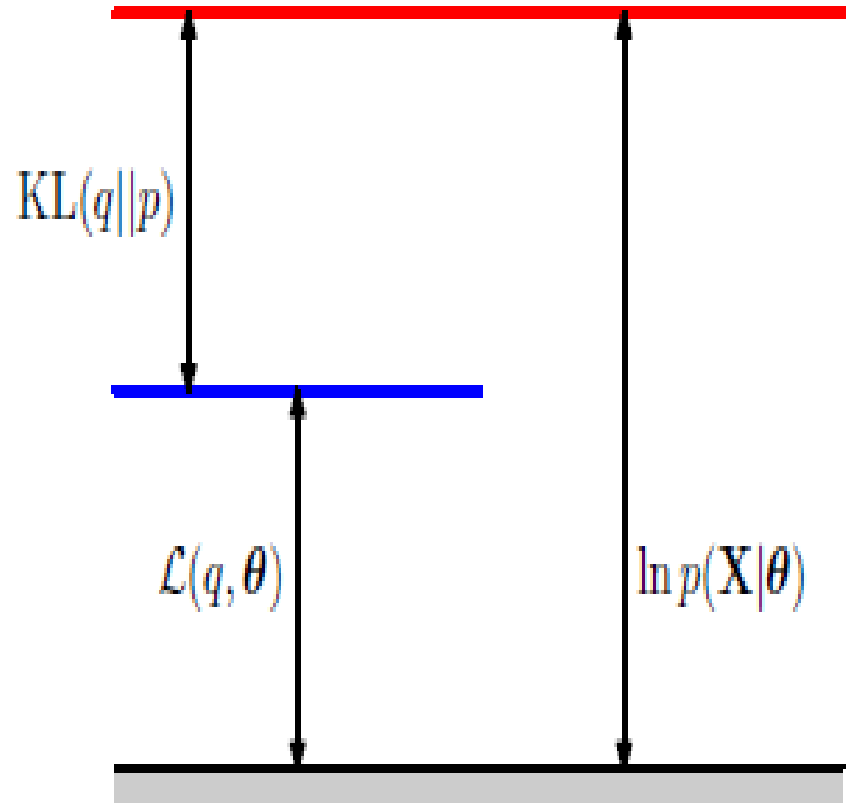
NOTE:

$KL(q||p)$ is the Kullback-Leibler divergence between $q(Z)$ and the posterior distribution $p(Z|X, \theta)$. Recall that the Kullback-Leibler divergence satisfies $KL(q||p) \geq 0$, with equality if, and only if, $q(Z) = p(Z|X, \theta)$. So, $\mathcal{L}(q, \theta) \leq \ln p(X|\theta)$, in other words $\mathcal{L}(q, \theta)$ is a lower bound on $\ln p(X|\theta)$

- The EM algorithm is a two-stage iterative optimization technique for finding maximum likelihood solutions.
- The 1st equation of slide 69 can be used to define the EM algorithm and to demonstrate that it does indeed maximize the log likelihood.

EM Algorithm

Figure 9.11 Illustration of the decomposition given by (9.70), which holds for any choice of distribution $q(\mathbf{Z})$. Because the Kullback-Leibler divergence satisfies $KL(q||p) \geq 0$, we see that the quantity $\mathcal{L}(q, \theta)$ is a lower bound on the log likelihood function $\ln p(\mathbf{X}|\theta)$.

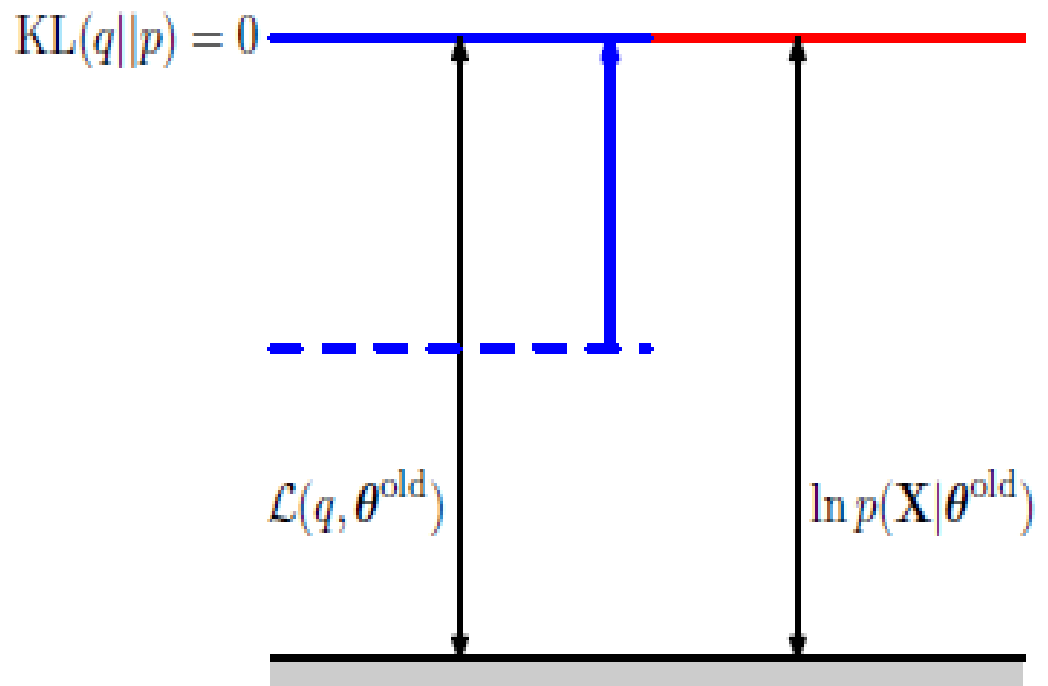


EM Algorithm

- Suppose that the current value of the parameter vector is θ^{old} .
- In the **E step**, the lower bound $\mathcal{L}(q, \theta^{\text{old}})$ is maximized with respect to $q(Z)$ while holding θ^{old} fixed.
- The solution to this maximization problem is easily seen by noting that the value of $\ln p(X|\theta^{\text{old}})$ does not depend on $q(Z)$ and so the largest value of $\mathcal{L}(q, \theta^{\text{old}})$ will occur when $q(Z)$ is equal to the posterior distribution $p(Z|X, \theta^{\text{old}})$.
- In this case, the lower bound will equal the log likelihood, as shown in the next figure.

EM Algorithm

Figure 9.12 Illustration of the E step of the EM algorithm. The q distribution is set equal to the posterior distribution for the current parameter values θ^{old} , causing the lower bound to move up to the same value as the log likelihood function, with the KL divergence vanishing.

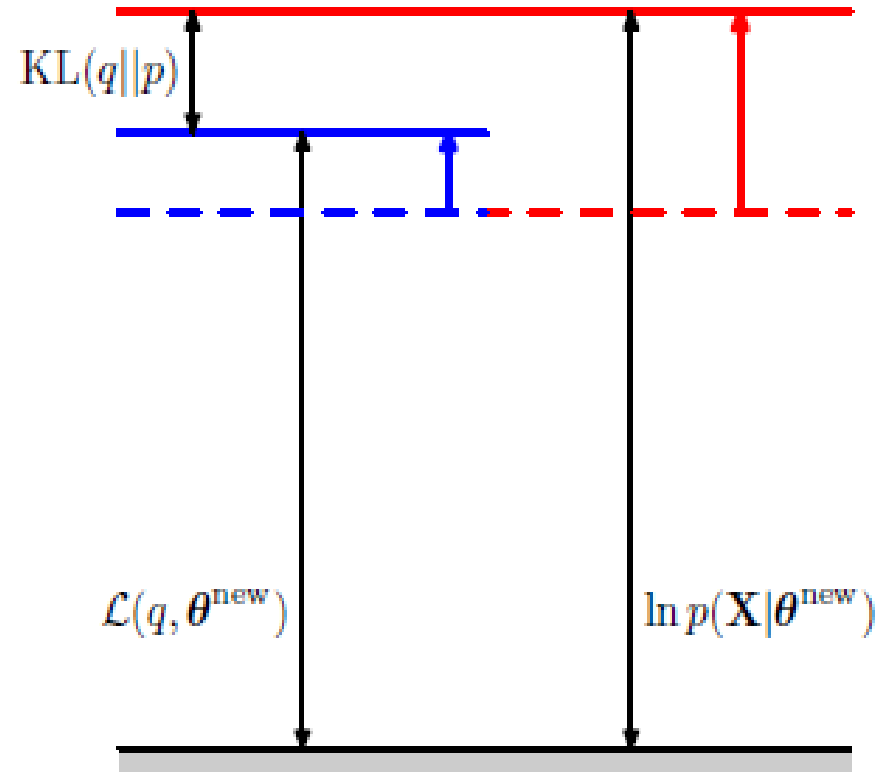


EM Algorithm

- In the subsequent **M step**, the distribution $q(Z)$ is held fixed and the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to θ to give some new value θ^{new} .
- This will cause the lower bound \mathcal{L} to increase (unless it is already at a maximum), which will cause the corresponding log likelihood function to increase.
- Since the distribution q is determined using the old parameter values rather than the new values and is held fixed during the M step, it will not equal the new posterior distribution $p(Z|X, \theta^{\text{new}})$, and hence there will be a nonzero KL divergence.
- The increase in the log likelihood function is therefore greater than the increase in the lower bound, as shown in next figure.

EM Algorithm

Figure 9.13 Illustration of the M step of the EM algorithm. The distribution $q(\mathbf{Z})$ is held fixed and the lower bound $\mathcal{L}(q, \theta)$ is maximized with respect to the parameter vector θ to give a revised value θ^{new} . Because the KL divergence is nonnegative, this causes the log likelihood $\ln p(\mathbf{X}|\theta)$ to increase by at least as much as the lower bound does.



EM Algorithm

- If we substitute $q(Z) = p(Z|X, \theta^{\text{old}})$ into 2nd equation of slide 69, we see that, after the E step, the lower bound takes the form:

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{X}, \mathbf{Z}|\theta) - \sum_{\mathbf{Z}} p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \ln p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}}) \\ &= Q(\theta, \theta^{\text{old}}) + \text{const}\end{aligned}\tag{9.74}$$

where the constant is simply the negative entropy of the q distribution and is therefore independent of θ .

EM Algorithm

- In the M step, revised parameter estimate θ^{new} is determined by maximizing the function:

$$\theta^{\text{new}} = \arg \max_{\theta} Q(\theta, \theta^{\text{old}}).$$

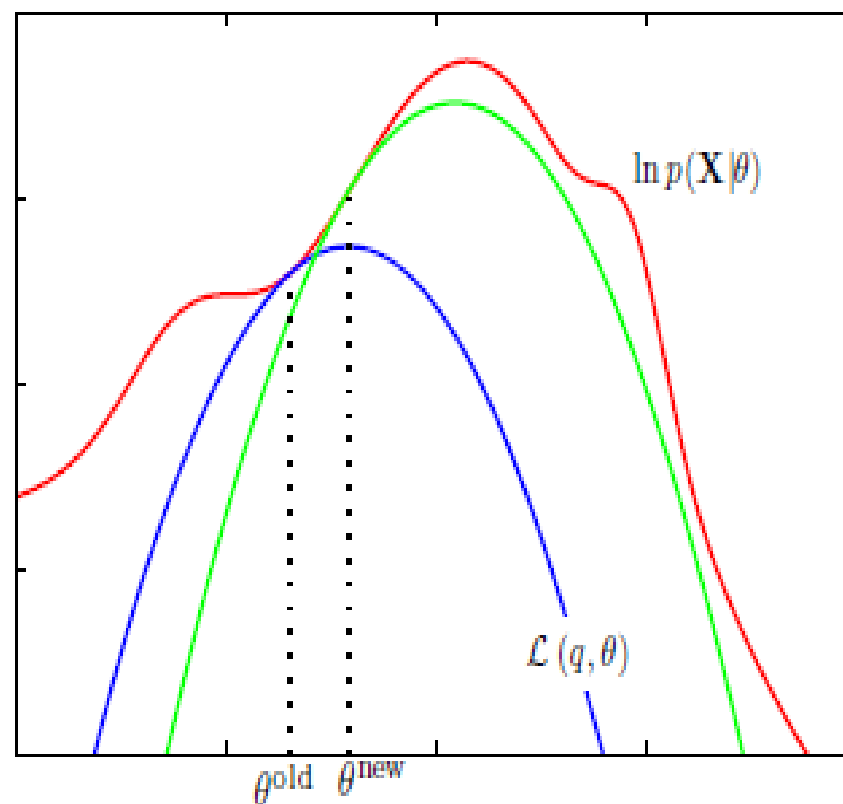
- In the definition of $Q(\theta, \theta^{\text{old}})$, the logarithm acts directly on the joint distribution $p(X, Z | \theta)$.

EM Algorithm

- So, in the M step, the quantity that is being maximized is the expectation of the complete-data log likelihood.
- Note that the variable θ over which we are optimizing appears only inside the logarithm.
- If the joint distribution $p(Z, X | \theta)$ comprises a member of the exponential family, or a product of such members, then we see that the logarithm will cancel the exponential and lead to an M step that will be typically much simpler than the maximization of the corresponding incomplete-data log likelihood function $p(X | \theta)$.

EM Algorithm

Figure 9.14 The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values. See the text for a full discussion.



EM Algorithm (explanation of fig. in last slide)

- The **red curve** depicts the (incomplete data) log likelihood function whose value we wish to maximize.
- We start with some initial parameter value θ^{old} , and in the first E step we evaluate the posterior distribution over latent variables, which gives rise to a lower bound $\mathcal{L}(\theta, \theta^{\text{old}})$ whose value equals the log likelihood at θ^{old} , as shown by the **blue curve**.
- Note that the bound makes a tangential contact with the log likelihood at θ^{old} , so that both curves have the same gradient.
- This bound is a convex function having a unique maximum (for mixture components from the exponential family).
- In the M step, the bound is maximized giving the value $\theta^{\text{(new)}}$, which gives a larger value of log likelihood than $\theta^{\text{(old)}}$.
- The subsequent E step then constructs a bound that is tangential at $\theta^{\text{(new)}}$ as shown by the **green curve**.

EM for Gaussian Mixtures

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters (comprising the means and covariances of the components and the mixing coefficients).

- 1) Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
- 2) **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \mu_j, \Sigma_j)}.$$

EM for Gaussian Mixtures

- 3) **M step.** Re-estimate the parameters using the current responsibilities

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{new}}) (\mathbf{x}_n - \mu_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}).$$

EM for Gaussian Mixtures

4) Evaluate the log likelihood:

$$\ln p(\mathbf{X}|\mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

and check for convergence of either the parameters or the log likelihood. If the convergence criterion is not satisfied return to step 2.

EM Algorithm

Algorithm 1 (Expectation-Maximization)

```
1 begin initialize  $\theta^0, T, i = 0$   
2       do  $i \leftarrow i + 1$   
3           E step : compute  $Q(\theta; \theta^i)$   
4           M step :  $\theta^{i+1} \leftarrow \arg \max_{\theta} Q(\theta; \theta^i)$   
5       until  $Q(\theta^{i+1}; \theta^i) - Q(\theta^i; \theta^{i-1}) \leq T$   
6       return  $\hat{\theta} \leftarrow \theta^{i+1}$   
7 end
```

References

- <https://www.geeksforgeeks.org/gaussian-mixture-model/>
- Pattern Classification (2nd ed) by R. O. Duda, P. E. Hart and D. G. Stork, John Wiley & Sons, 2000
- Pattern Recognition and Machine Learning by Christopher M. Bishop, Springer.

End