

Research Report 95
Tutkimusraportti 95

Feature Representation and Discrimination Based on Gaussian Mixture Model Probability Densities – Practices and Algorithms

**Pekka Paalanen, Joni-Kristian Kämäräinen, Jarmo Ilonen,
Heikki Kälviäinen**

Lappeenranta University of Technology
Department of Information Technology
P.O. Box 20
FIN-53851 Lappeenranta

ISBN 952-214-021-X
ISSN 0783-8069

Lappeenranta 2005

Feature Representation and Discrimination Based on Gaussian Mixture Model Probability Densities – Practices and Algorithms

P. Paalanen, J.-K. Kamarainen, J. Ilonen, H. Kälviäinen

Department of Information Technology
Lappeenranta University of Technology
P.O.Box 20, FI-53851 Lappeenranta, Finland

Abstract

Statistical methods have certain advantages which advocate their use in pattern recognition. One central problem in statistical methods is estimation of class conditional probability density functions based on examples in a training set. In this study maximum likelihood estimation methods for Gaussian mixture models are reviewed and discussed from a practical point of view. In addition, good practices for utilizing probability densities in feature classification and selection are discussed for Bayesian, and more importantly, for non-Bayesian tasks. As a result, the use of confidence information in the classification is proposed and a method for confidence estimation is presented. The propositions are tested experimentally.

1 Introduction

In former pattern recognition (PR) studies it was common practice to divide pattern recognition into sub-categories, such as structural, statistical, and neural approaches [21]. The division was loosely based on internal representations of patterns in different methods, but later it became evident that no strict taxonomy would apply and the subject itself contains equally important sub-topics, e.g., feature extraction and classifying features. In recent literature strict categorical divisions have been evaded (e.g., [6, 26]). However, the division into structural, statistical and neural approaches describes development of PR methods; structural representations are traditionally consistent with single instances of patterns while statistical approaches generalize instances under a more general representation, and finally, neural approaches are black boxes where the representation can only be indirectly affected. Lately, black box and gray box methods have proved to be some of the most powerful methods and methods such as multi-layer perceptron neural networks [3] and support vector machines [5] are frequently applied with a great success. Furthermore, novel methods seem to embed feature selection into a classifier synthesis as for example in AdaBoost boosting algorithm [9]. These powerful methods are also state-of-the-art methods in practice and it is justifiable to ask if structural and statistical approaches are anymore relevant at all.

Drawbacks in black and gray box PR methods are often their incapability to provide confidence information for their decision or difficulty to incorporate risk and cost models into the recognition process. In many applications it is not sufficient just to assign one of predefined classes to new observations; for example, in face detection facial evidence, such as eye centers and nostrils, should be detected from a scene and provided in a ranked order (best candidates first) to next processing level in order to perform detection computationally efficiently [11]. Gray box methods may include the confidence information as an ad hoc solution, but statistical methods usually provide the information in an interpretable form along with sufficient mathematical foundations. Statistical methods thus provide some advantages over black box methods; the decision making is based on an interpretable basis where the most probable or lowest risk (expected cost) option can be chosen (e.g. Bayesian decision making [22]) and different observations can be compared based on their statistical properties.

In a typical PR problem, features from known observations, a training set, are provided and necessary statistics must be established based on them for recognition of unknown observations and estimation of confidence. A class of patterns is typically represented as a probability density function (pdf) of features. Selection of proper features is a distinct and application specific problem, but as a more general consideration finding a proper pdf estimate has a crucial impact to successful recognition. Typically form of the pdf is somehow restricted and the search reduces to a problem of fitting the restricted model to observed features. Often already simple models, such as a single Gaussian distribution (normal distributed random variable), can efficiently represent patterns but a more general model, such as a finite mixture model, must be used to approximate more complex pdf's; arbitrarily complex probability density functions can be approximated using finite mixture models. Furthermore, the finite mixture representation is natural for certain kind of observations: observations which are produced by a randomly selected source from a set of alternative sources belonging to a same main class. This kind of task is natural when object categories are identified instead of object classes. For example, features from eye centers are naturally partitioned into closed eye and open eye, or Caucasian and Asian eye sub classes. The problem arises how probability densities should be approximated with finite mixture models and how the model parameters should be estimated. Equally important is to define correct practices to use pdf's in the pattern recognition and classification tasks.

In this study Gaussian mixture model (GMM) pdf's are studied as finite mixture models. The two main considerations with GMM are estimation of number of Gaussian components and robustness of the algorithm to tolerate singularities occurred due to a small sample size. It cannot be assumed that user knows all necessary details, and thus, the estimation should be unsupervised and utilize existing approximation and statistical theories. Several estimation methods have been proposed in literature and the most prominent ones have been experimentally evaluated in this study. The methods are extended to \mathbb{C}^n since the complex domain features, such as Gabor filter responses, seem to be natural for some applications [11, 18]. Correct classification practices are analyzed and defined based on problem characteristics: i) classifying an unknown observation into one of predefined classes, ii) finding best candidates from a set of observations, iii) deciding about belonging to a single known class when other classes are unknown or their samples are insufficient, and iv) concluding what useful statistical information should be provided to the next processing level. Finally, by providing implementations ([1]) for the described methods, we aim to encourage good practices when using GMM pdf's for representation and discrimination of patterns.

2 Gaussian mixture probability density function

Finite mixture models and their typical parameter estimation methods can approximate a wide variety of pdf's and are thus attractive solutions for cases where single function forms, such as a single normal distribution, fail. However, from a practical point of view it is often sound to form the mixture using one predefined distribution type, a basic distribution. Generally the basic distribution function can be of any type, but the multivariate normal distribution, the Gaussian distribution, is undoubtedly one of the most well-known and useful distributions in statistics, playing a predominant role in many areas of applications [27]. For example, in multivariate analysis most of the existing inference procedures have been developed under the assumption of normality and in linear model problems the error vector is often assumed to be normally distributed. In addition to appearing in these areas, the multivariate normal distribution also appears in multiple comparisons, in the studies of dependence of random variables, and in many other related fields. Thus, if there exists no prior knowledge of a pdf of phenomenon, only a general model can be used and the Gaussian distribution is a good candidate due to the enormous research effort in the past. For a more detailed discussion on the theory, properties and analytical results of multivariate normal distributions we refer to [27].

2.1 Multivariate normal distribution

A non-singular multivariate normal distribution of a D dimensional random variable $\mathbf{X} \mapsto \mathbf{x}$ can be defined as

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

where $\boldsymbol{\mu}$ is the mean vector and Σ the covariance matrix of the normally distributed random variable \mathbf{X} . In Figure 1 an example of 2-dimensional Gaussian pdf is shown. Multivariate Gaussian pdf's belong to the class of elliptically contoured distributions, which is evident in Fig. 1 where the equiprobability surfaces of the Gaussian are $\boldsymbol{\mu}$ -centered hyperellipsoids [27].

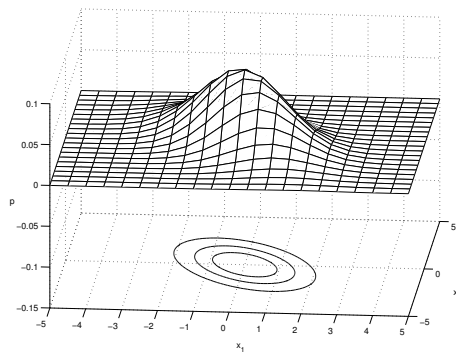


Figure 1: A two-dimensional Gaussian pdf and contour plots (equiprobability surfaces).

The Gaussian distribution in Eq. 1 can be used to describe a pdf of a real valued random vector ($\mathbf{x} \in \mathbb{R}^D$). A similar form can be derived for complex random vectors ($\mathbf{x} \in \mathbb{C}^D$) as

$$\mathcal{N}^{\mathbb{C}}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\pi^D |\boldsymbol{\Sigma}|} \exp [-(\mathbf{x} - \boldsymbol{\mu})^* \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})] \quad (2)$$

where $*$ denotes adjoint matrix (transpose and complex conjugate). [10]

2.2 Finite mixture model

Despite the fact that multivariate Gaussian pdf's have been successfully used to represent features and discriminate between different classes in many practical problems (e.g., [14, 19]) the assumption of single component leads to strict requirements for the phenomenon characteristics: a single basic class which smoothly varies around the class mean. The smooth behavior is not typically the most significant problem but the assumption of unimodality. For multimodally distributed features the unimodality assumption may cause an intolerable error to the estimated pdf and consequently into the discrimination between classes. The error is not caused only by the limited representation power but it may also lead to completely wrong interpretations (e.g. a class represented by two Gaussian distributions and another class between them). In object recognition this can be the case for such a simple thing as a human eye, which is actually an object category instead of a class since visual presence of the eye may include open eyes, closed eyes, Caucasian eyes, Asian eyes, eyes with glasses, and so on.

For a multimodal random variable, whose values are generated by one of several randomly occurring independent sources instead of a single source, a finite mixture model can be used to approximate the true pdf. If the Gaussian form is sufficient for single sources, then a Gaussian mixture model (GMM) can be used in the approximation. It should be noted that this does not necessarily need be the case but GMMs can also approximate many other types of distributions (see Fig. 2).

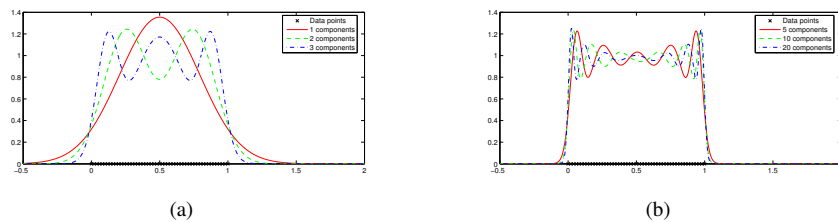


Figure 2: A uniform distribution $\mathcal{U}(0,1)$ is represented by 100 evenly spaced data points. The distribution is approximated using Gaussian mixture model and EM parameter estimation.

The GMM probability density function can be defined as a weighted sum of Gaussians

$$p(\mathbf{x}; \boldsymbol{\theta}) = \sum_{c=1}^C \alpha_c \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) \quad (3)$$

where α_c is the weight of c th component. The weight can be interpreted as *a priori* probability that a value of the random variable is generated by the c th source, and thus,

$0 \leq \alpha_c \leq 1$ and $\sum_{c=1}^C \alpha_c = 1$. Now, a Gaussian mixture model probability density function is completely defined by a parameter list [7]

$$\theta = \{\alpha_1, \mu_1, \Sigma_1, \dots, \alpha_C, \mu_C, \Sigma_C\} . \quad (4)$$

An example of Gaussian mixture model pdf is shown in Fig. 3.

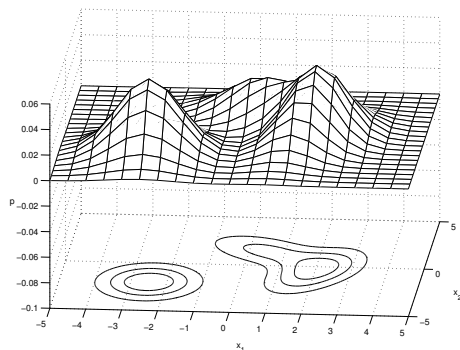


Figure 3: Surface of two-dimensional Gaussian mixture model pdf with three components and contour plots (equiprobability surfaces).

2.3 Estimating mixture model parameters

A vital question with GMM pdf's is how to estimate the model parameters θ . For a mixture of C components and a D dimensional random variable \mathbf{X} ($\mathbf{X} \mapsto \mathbf{x} \in \mathcal{R}^D$) the total number of parameters to be estimated is presented in Table 1. The number of free parameters is lower for pure complex data than for data where real and imaginary parts are concatenated to form a real vector of double length. However, if computing with real numbers is necessary, the mapping from \mathbb{C} to \mathbb{R}^2 should be chosen according to data. If magnitude and phase representation is used the phase may introduce a discontinuity into features [17]. Still, using purely complex representation may be advantageous requiring less training examples in parameter estimation. In the calculations the degree of freedom for a single real variable is 1 and for a single complex variable 2. The same identity has been used in the degrees of freedom for complex covariance matrix since it applies as an upper bound.

Table 1: Number of free parameters in a Gaussian mixture model.

Type	α_c	μ_c	Σ_c	Tot.
$\mathbf{x} \in \mathbb{R}^D$	1	D	$\frac{1}{2}D^2 + \frac{1}{2}D$	$C(\frac{1}{2}D^2 + \frac{3}{2}D) + C - 1$
$\mathbf{x} \in \mathbb{C}^D$	1	$2D$	D^2	$C(D^2 + 2D) + C - 1$
$\mathbf{x} \in \mathbb{C}^D \rightarrow \mathbb{R}^{2D}$	1	$2D$	$2D^2 + D$	$C(2D^2 + 3D) + C - 1$

In literature there exists two principal approaches for estimating the parameters: maximum-likelihood estimation and Bayesian estimation. While there are strong theoretical and methodological arguments supporting Bayesian estimation, in practice the maximum-likelihood estimation is simpler and, when used for designing classifiers,

can lead to classifiers nearly as accurate; many implementation issues support the selection of maximum-likelihood estimation. In this study the maximum-likelihood estimation is selected based on purely practical reasons.

2.3.1 Maximum-likelihood estimation

Assume that there is a set of independent samples $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ drawn from a single distribution described by a probability density function $p(\mathbf{x}; \boldsymbol{\theta})$ where $\boldsymbol{\theta}$ is the pdf parameter list. The likelihood function

$$\mathcal{L}(\mathbf{X}; \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{x}_n; \boldsymbol{\theta}) \quad (5)$$

tells the likelihood of the data \mathbf{X} given the distribution or, more specifically, given the distribution parameters $\boldsymbol{\theta}$. The goal is to find $\hat{\boldsymbol{\theta}}$ that maximizes the likelihood:

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}; \boldsymbol{\theta}) . \quad (6)$$

Usually this function is not maximized directly but the logarithm

$$L(\mathbf{X}; \boldsymbol{\theta}) = \ln \mathcal{L}(\mathbf{X}; \boldsymbol{\theta}) = \sum_{n=1}^N \ln p(\mathbf{x}_n; \boldsymbol{\theta}) \quad (7)$$

called the log-likelihood function which is analytically easier to handle. Because of the monotonicity of the logarithm function the solution to Eq. (6) is the same using $\mathcal{L}(\mathbf{X}; \boldsymbol{\theta})$ or $L(\mathbf{X}; \boldsymbol{\theta})$.

Depending on $p(\mathbf{x}; \boldsymbol{\theta})$ it might be possible to find the maximum analytically by setting the derivatives of the log-likelihood function to zero and by solving $\boldsymbol{\theta}$. For a Gaussian pdf the analytical solution leads to the well-known estimates of mean and variance, but usually the analytical approach is intractable. In practice an iterative method such as the expectation maximization (EM) algorithm is used. Maximizing the likelihood may in some cases lead to singular estimates, which is the fundamental problem of maximum likelihood methods with Gaussian mixture models [7].

If the parameters of the Gaussian mixture model pdf must be estimated for K different classes it is typical to assume independence, i.e., instances belonging to one class do not reveal anything about other classes. In the case of independent classes, the estimation problem of K class-conditional pdf's can be divided into K separate estimation problems.

2.3.2 Basic EM estimation

The expectation maximization (EM) algorithm is an iterative method for calculating maximum likelihood distribution parameter estimates from incomplete data (elements missing in feature vectors) [2]. The algorithm can also be used to handle cases where an analytical approach for maximum likelihood estimation is infeasible, such as Gaussian mixtures with unknown and unrestricted covariance matrices and means. In the following the notation and derivations correspond to the ones used by Duda et al. [6] and Figueiredo and Jain [8].

Assume that each training sample contains known features and missing or unknown features. Existing features are represented by X and all unknown features by Y . The expectation step (E-step) for the EM algorithm is to form the function

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i) \equiv E_Y[\ln \mathcal{L}(X, Y; \boldsymbol{\theta}) | X; \boldsymbol{\theta}^i] \quad (8)$$

where $\boldsymbol{\theta}^i$ is the previous estimate for the distribution parameters and $\boldsymbol{\theta}$ is the variable for a new estimate describing the (full) distribution. \mathcal{L} is the likelihood function in Eq. (5). The function calculates the likelihood of the data, including the unknown feature Y marginalized with respect to the current estimate of the distribution described by $\boldsymbol{\theta}^i$. The maximization step (M-step) is to maximize $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i)$ with respect to $\boldsymbol{\theta}$ and set

$$\boldsymbol{\theta}^{i+1} \leftarrow \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^i). \quad (9)$$

The steps are repeated until a convergence criterion is met.

For the convergence criterion it is suggested that (e.g. [6])

$$Q(\boldsymbol{\theta}^{i+1}; \boldsymbol{\theta}^i) - Q(\boldsymbol{\theta}^i; \boldsymbol{\theta}^{i-1}) \leq T \quad (10)$$

with a suitably selected T and that (e.g. [26])

$$\|\boldsymbol{\theta}^{i+1} - \boldsymbol{\theta}^i\| \leq \epsilon \quad (11)$$

for an appropriately chosen vector norm and ϵ . Common for the both of these criteria is that iterations are stopped when the change in the values falls below a threshold. A more sophisticated criterion can be derived from Eq. (10) by using a relative instead of an absolute rate of change.

The EM algorithm starts from an initial guess $\boldsymbol{\theta}^0$ for the distribution parameters and the log-likelihood is guaranteed to increase on each iteration until it converges. The convergence leads to a local or global maximum, but it can also lead to singular estimates, which is true particularly for Gaussian mixture distributions with arbitrary (not restricted) covariance matrices.

The initialization is one of the problems of the EM algorithm. The selection of $\boldsymbol{\theta}^0$ (partly) determines where the algorithm converges or hits the boundary of the parameter space producing singular, meaningless results. Some solutions use multiple random starts or a clustering algorithm for initialization [8].

In the case of Gaussian mixture models the known data X is interpreted as incomplete data. The missing part Y is the knowledge of which component produced each sample \boldsymbol{x}_n . For each \boldsymbol{x}_n there is a binary vector $\boldsymbol{y}_n = \{y_{n,1}, \dots, y_{n,C}\}$, where $y_{n,c} = 1$, if the sample was produced by the component c , or zero otherwise. The complete data log-likelihood is

$$\ln \mathcal{L}(X, Y; \boldsymbol{\theta}) = \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \ln(\alpha_c p(\boldsymbol{x}_n | c; \boldsymbol{\theta})) \quad (12)$$

The E-step contains computation of the conditional expectation of the complete data log-likelihood, the Q -function, given X and the current estimate $\boldsymbol{\theta}^i$ of the parameters. Since the complete data log-likelihood $\ln \mathcal{L}(X, Y; \boldsymbol{\theta})$ is linear with respect to the missing Y , the conditional expectation $W \equiv E[Y | X, \boldsymbol{\theta}^i]$ has simply to be computed and put into $\ln \mathcal{L}(X, Y; \boldsymbol{\theta})$. Therefore

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^i) \equiv E[\ln \mathcal{L}(X, Y; \boldsymbol{\theta}) | X, \boldsymbol{\theta}^i] = \ln \mathcal{L}(X, W; \boldsymbol{\theta}) \quad (13)$$

where the elements of W are defined as

$$w_{n,c} \equiv \mathbb{E} [y_{n,c} | \mathbf{X}, \boldsymbol{\theta}^i] = \Pr [y_{n,c} = 1 | \mathbf{x}_n, \boldsymbol{\theta}^i]. \quad (14)$$

The probability can be calculated with the Bayes law ([8])

$$w_{n,c} = \frac{\alpha_c^i p(\mathbf{x}_n | c; \boldsymbol{\theta}^i)}{\sum_{j=1}^C \alpha_j^i p(\mathbf{x}_n | j; \boldsymbol{\theta}^i)} \quad (15)$$

where α_c^i is the a priori probability (of estimate $\boldsymbol{\theta}^i$) and $w_{n,c}$ is the a posteriori probability that $y_{n,c} = 1$ after observing \mathbf{x}_n . In other words, $w_{n,c}$ is the probability that \mathbf{x}_n was produced by component c .

Applying the M-step to the problem of estimating the distribution parameters for C -component Gaussian mixture with arbitrary covariance matrices, the resulting iteration formulas are as follows:

$$\alpha_c^{i+1} = \frac{1}{N} \sum_{n=1}^N w_{n,c} \quad (16)$$

$$\boldsymbol{\mu}_c^{i+1} = \frac{\sum_{n=1}^N \mathbf{x}_n w_{n,c}}{\sum_{n=1}^N w_{n,c}} \quad (17)$$

$$\Sigma_c^{i+1} = \frac{\sum_{n=1}^N w_{n,c} (\mathbf{x}_n - \boldsymbol{\mu}_c^{i+1})(\mathbf{x}_n - \boldsymbol{\mu}_c^{i+1})^T}{\sum_{n=1}^N w_{n,c}}. \quad (18)$$

The new estimates are gathered to $\boldsymbol{\theta}^{i+1}$ (Eq. 4). If the convergence criterion (Eqs. 10 or 11) is not satisfied, $i \leftarrow i + 1$ and Eqs. (15)–(18) are evaluated again with new estimates.

The interpretation of the Eqs. (16)–(18) is actually quite intuitive. The weight α_c of a component is the portion of samples belonging to that component. It is computed by approximating the component-conditional pdf with the previous parameter estimates and taking the posterior probability of each sample point belonging to the component c (Eq. 15). The component mean $\boldsymbol{\mu}_c$ and covariance matrix Σ_c are estimated in the same way. The samples are weighted with their probabilities of belonging to the component, and then the sample mean and sample covariance matrix are computed.

It is worthwhile to note that hitherto the number of components C was assumed to be known. Clustering techniques try to find the true clusters and components from a training set, but our task of training a classifier only needs a good enough approximation of the distribution of each class. Therefore, C does not need to be guessed accurately, it is just a parameter defining the complexity of the approximating distribution. Too small C prevents the classifier from learning the sample distributions well enough and too large C may lead to an overfitted classifier. More importantly, too large C will definitely lead to singularities when the amount of training data becomes insufficient.

2.3.3 Figueiredo-Jain Algorithm

The Figueiredo-Jain (FJ) algorithm tries to overcome three major weaknesses of the basic EM algorithm [8]. The EM algorithm presented in Section 2.3.2 requires the user to set the number of components and the number remains fixed during the estimation

process. The FJ algorithm adjusts the number of components during estimation by annihilating components that are not supported by the data. This leads to the other EM failure point, the boundary of the parameter space. FJ avoids the boundary when it annihilates components that are becoming singular. FJ also allows starting with an arbitrarily large number of components, which tackles the initialization issue with the EM algorithm. The initial guesses for component means can be distributed into the whole space occupied by training samples, even setting one component for every single training sample.

The classical way to select the number of mixture components is to adopt the "model-class/model" hierarchy, where some candidate models (mixture pdf's) are computed for each model-class (number of components), and then select the "best" model. The idea behind the FJ algorithm is to abandon such hierarchy and to find the "best" overall model directly. Using the minimum message length criterion and applying it to mixture models leads to the objective function [8]

$$\Lambda(\boldsymbol{\theta}, \mathbf{X}) = \frac{V}{2} \sum_{c: \alpha_c > 0} \ln\left(\frac{N\alpha_c}{12}\right) + \frac{C_{\text{nz}}}{2} \ln \frac{N}{12} + \frac{C_{\text{nz}}(V+1)}{2} - \ln \mathcal{L}(\mathbf{X}, \boldsymbol{\theta}) \quad (19)$$

where N is the number of training points, V is the number of free parameters specifying a component, and C_{nz} is the number of components with nonzero weight in the mixture ($\alpha_c > 0$). $\boldsymbol{\theta}$ in the case of Gaussian mixture is the same as in Eq. (4). The last term $\ln \mathcal{L}(\mathbf{X}, \boldsymbol{\theta})$ is the log-likelihood of the training data given the distribution parameters $\boldsymbol{\theta}$ (Eq. 7).

The EM algorithm can be used to minimize Eq. (19) with a fixed C_{nz} [8]. It leads to the M-step with component weight updating formula

$$\alpha_c^{i+1} = \frac{\max\left\{0, \left(\sum_{n=1}^N w_{n,c}\right) - \frac{V}{2}\right\}}{\sum_{j=1}^C \max\left\{0, \left(\sum_{n=1}^N w_{n,j}\right) - \frac{V}{2}\right\}}. \quad (20)$$

This formula contains an explicit rule of annihilating components by setting their weights to zero. Other distribution parameters are updated as in Eqs. (17) and (18).

The above M-step is not suitable for the basic EM algorithm though. When initial C is high, it can happen that all weights become zero because none of the components have enough support from the data. Therefore a component-wise EM algorithm (CEM) is adopted. CEM updates the components one by one, computing the E-step (updating W) after each component update, where the basic EM updates all components "simultaneously". When a component is annihilated its probability mass is immediately redistributed strengthening the remaining components [8].

When CEM converges, it is not guaranteed that the minimum of $\Lambda(\boldsymbol{\theta}, \mathbf{X})$ is found, because the annihilation rule (Eq. 20) does not take into account the decrease caused by decreasing C_{nz} . After convergence the component with the smallest weight is removed and the CEM is run again, repeating until $C_{\text{nz}} = 1$. Then the estimate with the smallest $\Lambda(\boldsymbol{\theta}, \mathbf{X})$ is chosen. The implementation of the FJ algorithm uses a modified cost function instead of $\Lambda(\boldsymbol{\theta}, \mathbf{X})$ [8].

Hitherto the only assumptions for the mixture distribution are that the EM algorithm can be written for it and all components are parameterized the same way (number of

parameters V for a component). With Gaussian mixture model the number of parameters per component is $V = D + \frac{1}{2}D^2 + \frac{1}{2}D$ in the case of real valued data and arbitrary covariance matrices. With complex valued data the number of real free parameters $V = 2D + D^2$ where D is the dimensionality of the (complex) data. For complex data the number of free parameters should be replaced by the new value for V instead the one by Figueiredo and Jain who derived the rule for real valued data [8]. As can be seen in Eq. (20) the new value amplifies the annihilation which makes sense because there are more degrees of freedom in a component. On the other hand there are more training data with the same number of training points, because the data is complex (two values in one variable as opposed to one value).

2.3.4 Greedy EM algorithm

The greedy algorithm starts with a single component and then adds components into the mixture one by one [28]. The optimal starting component for a Gaussian mixture is trivially computed, optimal meaning the highest training data likelihood. The algorithm repeats two steps: insert a component into the mixture, and run EM until convergence. Inserting a component that increases the likelihood the most is thought to be an easier problem than initializing a whole near-optimal distribution. Component insertion involves searching for the parameters for only one component at a time. Recall that EM finds a local optimum for the distribution parameters, not necessarily the global optimum which makes it initialization dependent method.

Let p_C denote a C -component mixture with parameters θ_C . The general greedy algorithm for Gaussian mixture is as follows [28]:

Algorithm 1 Greedy EM

- 1: Compute the optimal (in the ML sense) one-component mixture p_1 and set $C \leftarrow 1$.
- 2: Find a new component $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}', \Sigma')$ and corresponding mixing weight α' that increase the likelihood the most:

$$\{\boldsymbol{\mu}', \Sigma', \alpha'\} = \arg \max_{\{\boldsymbol{\mu}, \Sigma, \alpha\}} \sum_{n=1}^N \ln[(1 - \alpha)p_C(\mathbf{x}_n) + \alpha \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}, \Sigma)] \quad (21)$$

while keeping p_C fixed.

- 3: Set $p_{C+1}(\mathbf{x}) \leftarrow (1 - \alpha')p_C(\mathbf{x}) + \alpha' \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}', \Sigma')$ and then $C \leftarrow C + 1$.
- 4: Update p_C using EM (or some other method) until convergence. [optional]
- 5: Evaluate some stopping criterion; go to Step 2 or quit.

The stopping criterion in Step 5 can be for example any kind of model selection criterion, wanted number of components, or the minimum message length criterion.

The crucial point is of course Step 2. Finding the optimal new component requires a global search, which is performed by creating $C N_{\text{cand}}$ candidate components. The number of candidates will increase linearly with the number of components C , having N_{cand} candidates per each existing component. The candidate resulting in the highest likelihood when inserted into the (previous) mixture is selected. The parameters and weight of the best candidate are then used in Step 3 instead of the truly optimal values.

The candidates for executing Step 2 are initialized as follows: the training data set \mathbf{X} is partitioned into C disjoint data sets $\{A_c\}$, $c = 1 \dots C$, according to the posterior

probabilities of individual components; the data set is Bayesian classified by the mixture components. From each A_c number of N_{cand} candidates are initialized by picking uniformly randomly two data points \mathbf{x}_l and \mathbf{x}_r in A_c . The set A_c is then partitioned into two using the smallest distance selection with respect to \mathbf{x}_l and \mathbf{x}_r . The mean and covariance of these two new subsets are the parameters for two new candidates. The candidate weights are set to half of the weight of the component that produced the set A_c . Then new \mathbf{x}_l and \mathbf{x}_r are drawn until N_{cand} candidates are initialized with A_c . The partial EM algorithm is then used on each of the candidates. The partial EM differs from the EM and CEM algorithms by optimizing (updating) only one component of a mixture, it does not change any other components. In order to reduce the time complexity of the algorithm a lower bound on the log-likelihood is used instead of the true log-likelihood. The lower-bound log-likelihood is calculated with only the points in the respective set A_c . The partial EM update equations are as follows:

$$w_{n,C+1} = \frac{\alpha \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}, \Sigma)}{(1 - \alpha)p_C(\mathbf{x}) + \alpha \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma)} \quad (22)$$

$$\alpha = \frac{1}{\aleph(A_c)} \sum_{n \in A_c} w_{n,C+1} \quad (23)$$

$$\boldsymbol{\mu} = \frac{\sum_{n \in A_c} w_{n,C+1} \mathbf{x}_n}{\sum_{n \in A_c} w_{n,C+1}} \quad (24)$$

$$\Sigma = \frac{\sum_{n \in A_c} w_{n,C+1} (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T}{\sum_{n \in A_c} w_{n,C+1}} \quad (25)$$

where $\aleph(A_c)$ is the number of training samples in the set A_c . These equations are much like the basic EM update equations in Eqs. (16)–(18). The partial EM iterations are stopped when the relative change in log-likelihood of the resulting $C + 1$ -component mixture drops below threshold or maximum number of iterations is reached. When the partial EM has converged the candidate is ready to be evaluated.

2.3.5 Avoiding covariance matrix singularities

All above estimation methods may fail due to singularities appearing during computation of new estimates. Several heuristic explicit covariance matrix fixing procedures can be applied in order to prevent singularities. Also, Nagy et al. [20] present an algorithm using factorized covariance matrices that avoids singularities altogether. The method used with the EM and FJ algorithms is described next.

Computation numerics may introduce inaccuracies so that the covariance matrix is no longer strictly Hermitian (complex conjugate symmetric, $\Sigma = \Sigma^*$). The matrix is forced to Hermitian by

$$\Sigma \leftarrow \Sigma - \frac{\Sigma - \Sigma^*}{2} \quad (26)$$

and removing imaginary part from the diagonal.

The matrix is tested for positive definiteness with Cholesky factorization (Chol). While the test is negative, the diagonal of the matrix is modified. If the diagonal contains elements that are less than 10ϵ the diagonal values are grown by an amount relative to the largest absolute value on the diagonal, also taking account that negative diagonal elements become positive enough. Otherwise the diagonal is grown by 1 percent.

2.4 Experiments

All the above methods are publicly available in the Matlab software package provided by the authors [1]. In the experiments discussed next the methods were applied in classification tasks with several publicly available data sets. A comparison between three methods, EM, FJ and GEM was performed to inspect their accuracy and robustness. The experiments were conducted by varying both the parameter related to number of components and the amount of data used for training. For FJ and GEM the number of components denoted the maximum number of components. Mean and maximum accuracies and algorithm crash rate were inspected. A crash is a situation where the algorithm does not produce a final mixture estimate and cannot continue, e.g., in EM a covariance matrix becomes undefined or in FJ all mixture components are annihilated.

In the first experiment forest spectral data first introduced by Jaaskelainen et al. [15] was used. Since the data was of high dimensionality it was projected to fewer dimensions as proposed in [16]. For training, varying amounts of the data (20–70%) were randomly selected and 30% were used for classification by Bayes rule. In Fig. 4(a) it can be seen how all methods performed equally well and minimum and maximum accuracies were both near the mean accuracy. It should be noted that in Fig. 4(a) FJ and GEM always converged to the same number of components, that is, a single Gaussian. GEM was the most stable in terms of crash rate while EM frequently crashed for a non-optimal number of components and FJ algorithm did not succeed until enough training examples were included (Fig. 4(b)). The crash rate is represented as a proportion to the total amount of re-executions of the algorithm and re-executions with reselected data. Furthermore, there were no significant differences in maximal accuracies of different methods when they succeeded in estimation as shown in Figs. 4(c) and 4(d). However, for a good result EM must be executed several times.

Next, the same experiment was conducted for the forest data projected into 10 dimensions and results are shown in Fig. 5. Fewer dimensions produced smoother data, which can be seen from the results as all methods performed more reliably and even the FJ did not crash. However, more discriminating information was lost and consequently the maximal accuracy decreased. Otherwise the methods behaved equally as in the first experiment.

In the third experiment well-known waveform data was used. Waveform data consisted of 5000 40-element vectors, where the first attributes are meaningful attributes with added noise and the last 19 elements are pure noise. The optimal Bayes classification accuracy is 86% [4]. FJ reached the optimum with one Gaussian component and the other methods came very close (Fig. 6). In Fig. 6(a) it seems that there is enough data for slight overfitting. The Greedy EM (GEM) shows that it adds components quite greedily. Otherwise all methods performed equally.

Based on the experiments it can be said that the standard EM algorithm outperforms FJ and GEM if a good prior knowledge exists about the number of components. Both FJ and GEM automatically estimate the number of components and it seems that the FJ produces more accurate results, but also requires more training samples than GEM. As an implementation the GEM is the most stable while the standard EM is the most unstable. The unsupervised nature of FJ and GEM still motivates their use in many practical applications.

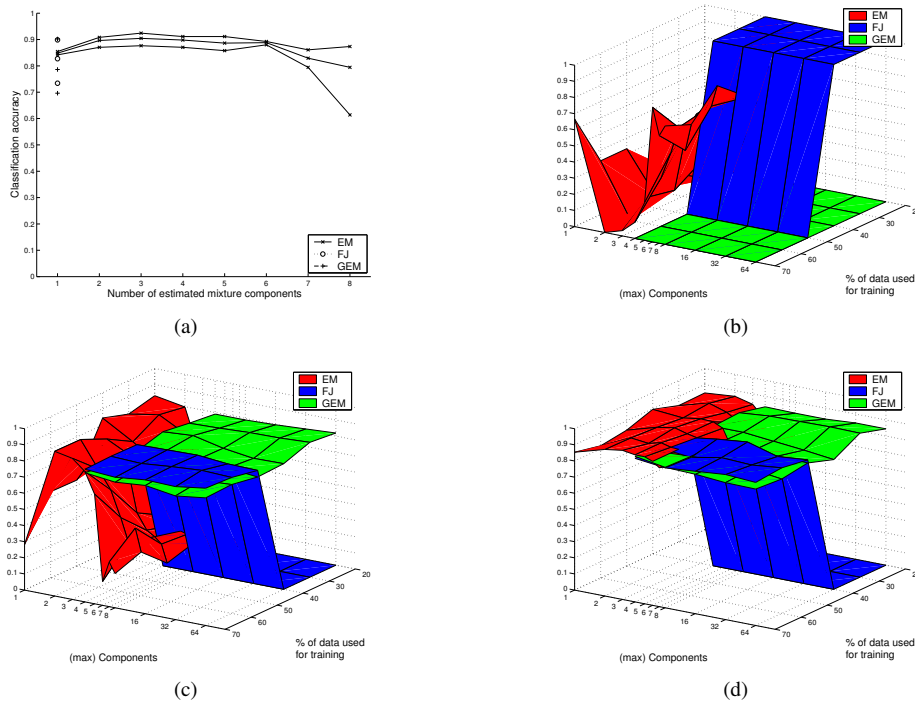


Figure 4: Results for the 23-d forest data as functions of amount of training data: (a) classification accuracy (70% of data used for training and min., max., and mean curves plotted); (b) algorithm crash rate; (c) mean accuracy; (d) maximum accuracy.

3 Feature discrimination and classification

Bayesian decision making is the most standard way to classify observations if class conditional pdf's are known, but the classification can be applied only to one of the four items mentioned in the introduction, i.e., i) classification of an unknown observation into one of predefined classes. It should be noted that the plain use of the Bayes formula does not guarantee Bayesian optimal decision making and there are numerous pattern recognition problems which cannot be formulated as a Bayesian task [22].

In this study Gaussian mixture models are considered for representing class conditional pdf's and the Bayesian rule can be used for selecting which class to assign for an unclassified observation. The weakness of the approach is that a posteriori probabilities can be used only to make a decision between known classes for a single observation. It is important to note that posteriors cannot be used to compare different observations, for example, to decide which one of them is more reliably from a certain class. The comparison between instances must be based on some other statistical reasoning.

The Bayesian decision making was useful only in the first problem defined in the introduction, but the other two problems: ii) finding best candidates from a set of observations and iii) deciding about belonging to a single known class when other classes are unknown or their samples are insufficient must be formulated in another way. In this study use of confidence information is proposed and confidence additionally covers the

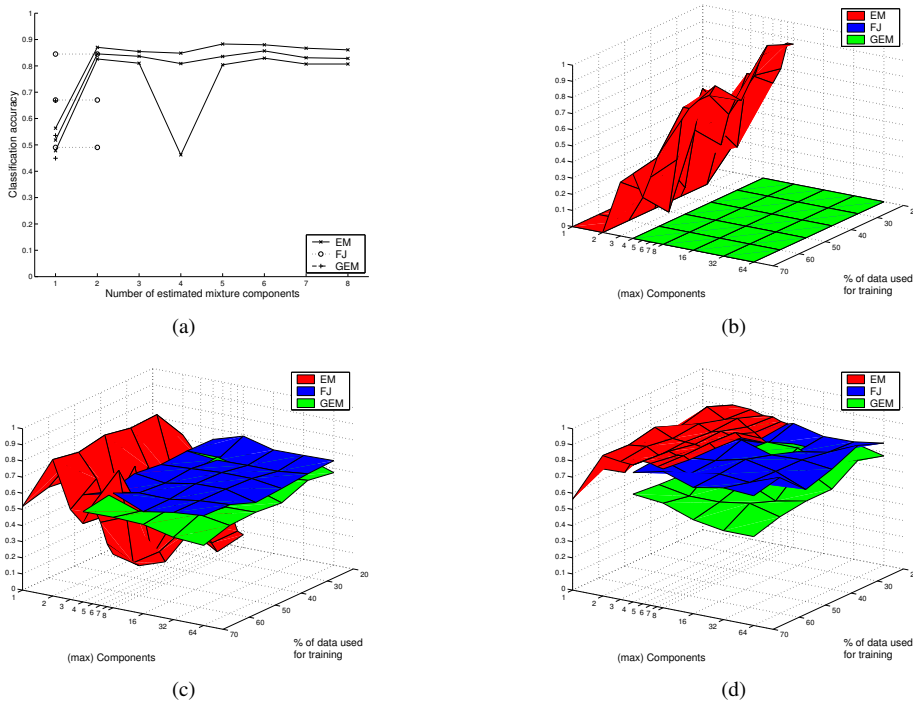


Figure 5: Results for the 10-d forest data as functions of amount of training data: (a) classification accuracy (70% of data used for training and min., max., and mean curves plotted); (b) algorithm crash rate; (c) mean accuracy; (d) maximum accuracy.

last problem, iv) concluding what useful statistical information should be provided to the next processing level.

3.1 Classification using the Bayesian decision rule

Bayesian classification and decision making are based on probability theory and the principle of choosing the most probable or the lowest risk (expected cost) option [22, 26]. Assume that there is a classification task to classify feature vectors (observations) to K different classes. A feature vector is denoted as $\mathbf{x} = [x_1, x_2, \dots, x_D]^T$ where D is the dimension of a vector. Probability that a feature vector \mathbf{x} belongs to a class ω_k is $P(\omega_k|\mathbf{x})$, and it is often referred to as a posteriori probability. The classification of the vector is done according to posterior probabilities or decision risks calculated from the probabilities [22].

The posterior probabilities can be computed with the Bayes formula

$$P(\omega_k|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_k)P(\omega_k)}{p(\mathbf{x})} \quad (27)$$

where $p(\mathbf{x}|\omega_k)$ is the probability density function of class ω_k in the feature space and $P(\omega_k)$ is the a priori probability which tells the probability of the class before measuring any features. If a priori probabilities are not actually known, they must be explicitly

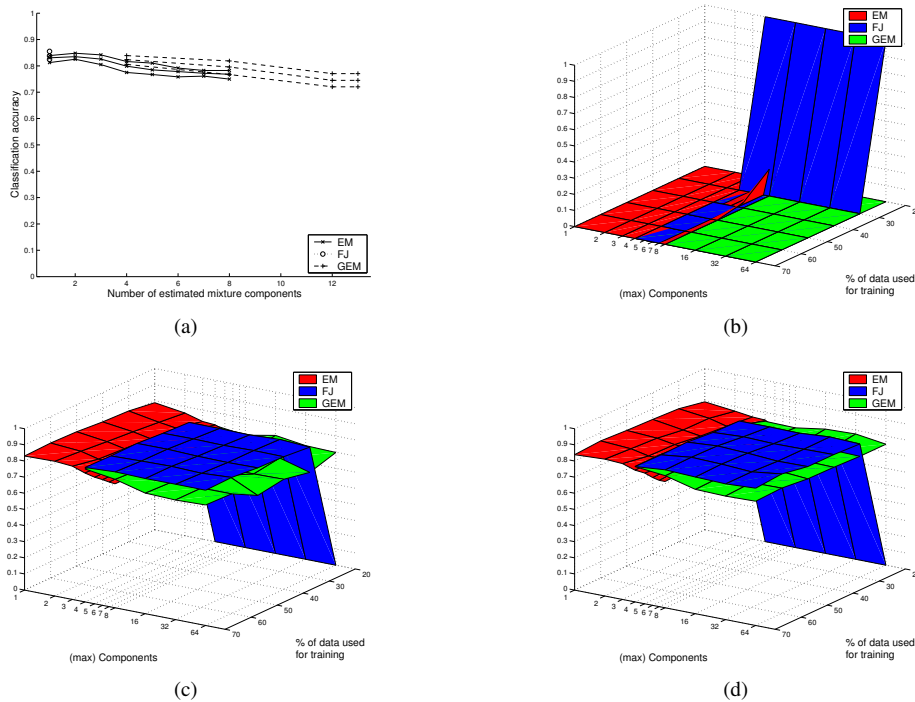


Figure 6: Results for the waveform data as functions of amount of training data: (a) classification accuracy (70% of data used for training and min., max., and mean curves plotted); (b) algorithm crash rate; (c) mean accuracy; (d) maximum accuracy.

defined or estimated from the training set. A fatal error exists if a priories do not exist [22]. The divisor

$$p(\mathbf{x}) = \sum_{i=1}^K p(\mathbf{x}|\omega_i)P(\omega_i) \quad (28)$$

is merely a scaling factor to assure that posterior probabilities are really probabilities, i.e., their sum is one.

It can be shown that choosing the class of the highest posterior probability produces the minimum error probability [6, 22, 26]. However, if the cost of making different kinds of errors is not uniform, the decision can be made with a risk function that computes the expected cost using the posterior probabilities, and choose the class with minimum risk [22].

A central component in the Bayesian classification is the class-conditional probability density function $p(\mathbf{x}|\omega_k)$. The function tells the distribution of feature vectors in the feature space inside a particular class, i.e., it describes the class model. In practice it is always unknown except in some artificial classification tasks. The distribution can be estimated from a training set with the described methods assuming Gaussian mixture models.

3.2 Classification using confidence

The term confidence may have various different meanings and formulations, but in our case confidence is used to measure reliability of a classification result where a class label is assigned to an observation. If confidence is low it is more probable that a wrong decision have been made. Intuitively value of class conditional pdf at an observation corresponds to confidence for a specific class: the higher the pdf value is, the more instances of a corresponding class appear similar as the observation.

A posteriori is a between-class measure for a single observation, but pdf values can be used as an inter-class measure to select the best representative of a class. For example in object detection it is computationally efficient to process best evidence candidates representing parts of an object first [12]. Unlikely evidence candidates can be pruned as outliers. In such tasks the use of confidence is beneficial.

3.2.1 Interpretation of confidence

A posteriori values can be used to select the most probable class for a single observation, but confidence values can be used to select the most reliable class representative among many observations. In certain tasks, confidence can be used to discard observations which cannot be sufficiently reliably assigned to any of known classes, that is, their pdf values are too low for a reliable decision making. The same approach can be used in two class problems where training examples are available only for one class, and now, observations which do not pass a predefined confidence level are assigned to the unknown class (e.g. in detection of motor failure conditions [19]). In that sense the confidence does not refer to a single value limit but to a certain allowed region in feature space.

Definition 1 Confidence value $\kappa \in [0, 1]$ and confidence region $\mathcal{R} \subseteq \Omega$ for a probability distribution function $0 \leq p(\mathbf{x}) < \infty, \forall \mathbf{x} \in \Omega$. κ is a confidence value related to a non-unique confidence region \mathcal{R} such that

$$\int_{\Omega \setminus \mathcal{R}} p(\mathbf{x}) d\mathbf{x} = \kappa \quad (29)$$

The definition of confidence in Definition 1 is easily interpretable via the confidence region \mathcal{R} . It is a region which covers a proportion $1 - \kappa$ of the probability mass of $p(\mathbf{x})$ because for probability distributions $\int_{\Omega} p(\mathbf{x}) d\mathbf{x} = 1$. It is clear that $\kappa = 1$ for \mathcal{R} contains only a finite number of individual points and $\kappa = 0$ for $\mathcal{R} = \Omega$. It should be noted that the confidence value has no use until the region \mathcal{R} is defined as a minimal volume region which satisfies Definition 1. The minimal volume region is called the highest density region (HDR) [13].

For each $k = 1, \dots, K$ different classes a class specific confidence value κ_k can be defined, but intuitively the same value is good for all classes. A confidence value corresponds to a proportion of probability mass that retains in an area \mathcal{R}_k . In a classification task where certain confidence for decision making is required the confidence value itself is not used, but actually the confidence region \mathcal{R}_k is important since a sample vector \mathbf{x} is allowed to enter the class ω_k only if $\mathbf{x} \in \mathcal{R}_k$. If a sample is not within the confidence region of any of the classes it must be classified to a garbage class. The

garbage class is a special class and samples assigned to the class need a special attention; for example, more information is needed for observations in the garbage class or in a two-class problem where data is available only from one class the garbage class may represent the other class with an unknown distribution.

There are various uses for the confidence value in classification. For example, in the search for the best evidence sets from a large set of samples the use of confidence value may decrease the work done in further processing steps. The confidence may also be directly used to guarantee the confidence of the classification. In condition monitoring mechanical failures typically evolve during the process and too early warning may lead to too pessimistic repairing procedures. Furthermore, as related to the condition monitoring, it is often easy to collect samples from normal system conditions, but all failures can never be comprehensively recorded, and thus, two class classification to normal or failure condition can be performed utilizing the confidence value: samples within the confidence region are assigned to the normal class and outside that to the failure condition whose probability distribution is unknown. Correct use and advantages of using confidence in classification will be demonstrated and further discussed in the experimental part of this study.

3.2.2 Analytical solution for confidence

As shown in Eq. 3 pdf estimated by a finite Gaussian mixture model $p(\mathbf{x})$ can be represented as

$$p(\mathbf{x}) = \sum_{c=1}^C \alpha_c p_c(\mathbf{x}) \quad (30)$$

where (given in \mathbb{R}^n)

$$p_c(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}_c|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \right] . \quad (31)$$

The problem is to

$$\mathcal{R} \leftarrow \arg \min \text{Vol}(\mathcal{R}), \text{ such that } \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x} = 1 - \kappa , \quad (32)$$

that is, it is desired to find the minimal volume region \mathcal{R} that satisfies the given equality, i.e., covers the desired probability mass.

As already mentioned the multivariate Gaussian (normal) distribution belongs to the elliptically contoured family of distributions [27]. For a single normal pdf it is straightforward to show that \mathcal{R} which satisfies Eq. (32) is a hyperellipsoid [27]

$$\mathcal{R} = \{ \mathbf{x} \in \Omega \mid (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq r^2 \} . \quad (33)$$

Distribution of the squared Mahalanobis distance $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$ in Eq. (33) is the χ^2 distribution where the number of degrees of freedom parameter is the dimensionality of \mathbf{x} . Therefore r^2 can be found by computing inverse of the cumulative χ^2 pdf.

Correspondingly a *proposition* can be made that the solution for a Gaussian mixture model pdf is

$$\mathcal{R} = \bigcup_{c=1}^C \mathcal{R}_c, \text{ where } \mathcal{R}_c = \{\mathbf{x} \in \Omega \mid (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) \leq r_c^2\}, \quad (34)$$

that is, the minimum volume region that covers the desired probability mass in a Gaussian mixture model pdf consists of hyperellipsoid regions concentrated near the mean of each mixture model component. The proposition can be only approximately correct since it claims that ellipsoidal areas span the same axis relations as the Gaussian covariance matrices, but it is not true since the probability masses of separate components affect to each other. If the proposition would however be accepted the minimal confidence region \mathcal{R} could be numerically searched by a gradient descent method. The formulation would be simple for a single multi-dimensional Gaussian, but becomes more difficult when there are several non-overlapping Gaussians and very complex when the overlapping is allowed: integration must be done over all components in every ellipsoid. The required theory for computations exist [23, 24, 25], but the optimization may still become computationally infeasible, and thus, feasible and sufficiently accurate approximating methods are needed. One such method will be described next.

3.2.3 Rank-order statistics based confidence estimation

To find the confidence region a reverse approach can be used to find a pdf value τ which is at the border of the confidence region. It is assumed that the pdf is continuous and the gradient of the pdf is never zero in a whole neighborhood of any point where the pdf value is nonzero. τ must be equal everywhere in the border, otherwise the region cannot be the minimal volume region [13]. τ can be computed by rank-order statistics using the density quantile ([13]) $F(\tau)$ and by generating data according to the pdf function. Density quantile is the opposite of confidence $\kappa = 1 - F(\tau)$.

$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ distributed data can be generated by generating D -dimensional vectors \mathbf{y}_i , each element i.i.d. $\mathcal{N}(0, 1)$, and computing $\mathbf{x}_i = \mathbf{y}_i \text{ Chol } \boldsymbol{\Sigma} + \boldsymbol{\mu}$ [27]. Complex valued data can be generated by generating D -dimensional vectors \mathbf{z}_i with each element i.i.d. $\mathcal{U}(0, \pi)$ and computing $\mathbf{x}_i = (\mathbf{y}_i \cos \mathbf{z}_i + \mathbf{y}_i \mathbf{i} \sin \mathbf{z}_i) \text{ Chol } \boldsymbol{\Sigma} + \boldsymbol{\mu}$. Gaussian mixture data can be generated by first randomly selecting one of the components using component weights as probabilities and then generating a point with the parameters of the selected component.

Definition 2 *The density quantile $F(\tau)$ corresponding to a pdf value τ is*

$$F(\tau) = \int_{p(\mathbf{x}) \geq \tau} p(\mathbf{x}) d\mathbf{x} . \quad (35)$$

The density quantile corresponds to HDR probability mass, an integral over a region where pdf values remain above the given minimal level τ . Clearly $F(\tau) \in [0, 1]$ since it is similar to a cumulative pdf. The pdf value τ has no interpretation which could be used, but it is easy to define a density quantile, e.g., quantile of 0.9 would accept the most typical 90% of instances generated by the pdf. The quantile is a non-increasing function, and thus, a reverse mapping $h(F) = \tau$ can be defined, where $F \in [0, 1]$ is the desired quantile.

The density quantile has a well-known connection to confidence bounds of a normal distribution and its pdf. When extended to mixtures of normal distributions the confidence bounds may become arbitrarily complex as the number of mixture components increases. An example of one-dimensional case is illustrated in Figure 7. Regions turn to be even more difficult as moving to multidimensional complex spaces \mathbb{C}^D .

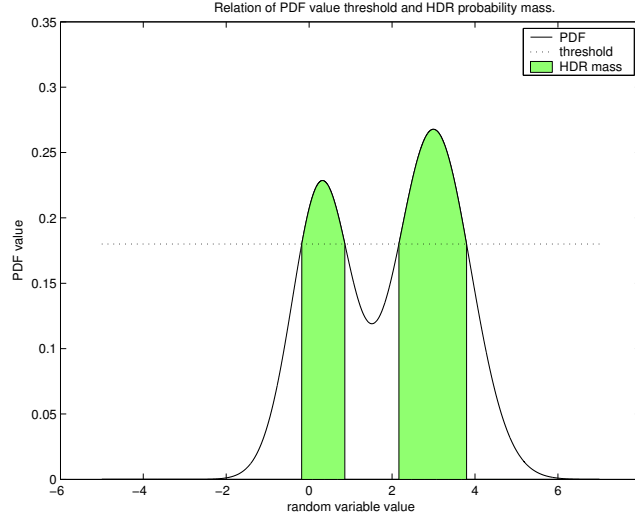


Figure 7: A highest density region (HDR) of a two-component GMM pdf and the corresponding threshold in one dimension. The confidence region is no longer a simple connected set.

Analytical solution for $F(\tau)$ or $h(F)$ can be very difficult, but a Monte Carlo method can be used to approximate the functions. Computation utilizes random sampling by generating N random points $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ from distribution described by $p(\mathbf{x})$. For every point pdf value is computed $p_i = p(\mathbf{x}_i)$ and all p_i 's are ordered into a sequence $Y = (y_1, \dots, y_N) = \text{sort} \{p_i\}$ in an ascending order. Y represents now a non-decreasing finite series which can be used to estimate $F(\tau)$ and $h(F)$. We have now reduced the original space \mathbb{C}^D into (a discrete approximation of) \mathbb{R}^+ . A similar method was proposed by Hyndman [13].

Now, value of $F(\tau)$ can be approximated by

$$i = \arg \max_i \{y_i \mid y_i \leq \tau\} \quad (36)$$

and, with linear interpolation,

$$F(\tau) \approx \begin{cases} 1 & \text{if } \tau < y_1 \\ 0 & \text{if } \tau \geq y_N \\ 1 - \frac{i-1+l(i,\tau)}{N-1} & \text{otherwise} \end{cases} \quad (37)$$

where the linear interpolation term

$$l(i, \tau) = \begin{cases} 0.5 & \text{if } y_{i+1} - y_i = 0 \\ \frac{\tau - y_i}{y_{i+1} - y_i} & \text{otherwise} \end{cases} \quad (38)$$

The reverse mapping $h(F)$ can be approximated by

$$\tau = h(F) \approx \begin{cases} y_N & \text{if } i = N \\ y_i + ((N - 1)(1 - F) + 1 - i)(y_{i+1} - y_i) & \text{otherwise} \end{cases} \quad (39)$$

where $i = \lfloor (N - 1)(1 - F) + 1 \rfloor$.

Now, the estimated pdf value τ can be used as a limit pdf value where observations falling below are directed to the “garbage class”. Hitherto, only Gaussian mixture models have been applied, but the proposed estimation approach applies to any continuous pdf’s that fulfill the gradient condition given in this section.

3.3 Experiments

The Bayesian decision making was already demonstrated in the experiments where different estimation algorithms were compared and surely in literature there is an enormous number of experimental studies utilizing Bayesian decision making. Here two examples are represented in order to demonstrate the use of confidence.

In the first experiment the use of confidence information was studied in an application where its usability is apparent. In face detection methods where detection of a whole face is based on smaller facial parts, evidence, features are extracted from every spatial location (e.g., [11, 12]) and ranked as possible evidence from different classes. The detection of the whole face is based on inspection of spatial connections of different evidence. In Fig. 8(a) an example facial image is shown and the 10 ground truth evidence are marked [11]. In the training phase Gabor features were extracted from training set images and class conditional pdf’s were estimated for each evidence class by the FJ method. In Fig. 8(b) a pdf surface is shown for the evidence class number 7 corresponding to the left (left on the image) nostril. Figs. 8(c) and 8(d) display the confidence regions of density quantiles 0.5 and 0.05 corresponding to 0.5 and 0.95 confidence values respectively. It is clear that the correct evidence location, the left nostril, was already included in the 0.95 confidence region, and thus, evidence search space for the next processing level was reduced dramatically. Using the confidence information it is possible to rank evidence from each class and provide at most a requested number of best candidates to a next processing level [11], but also allowing the possibility that evidence is not found on the image.

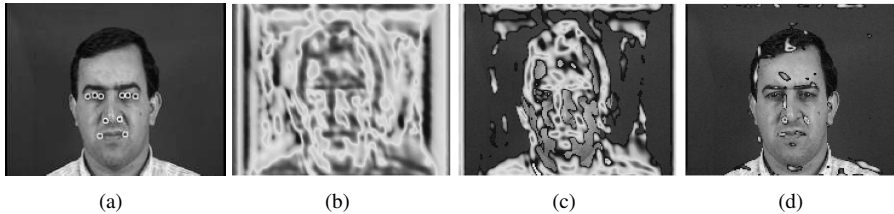


Figure 8: Example of using density quantile for defining confidence regions : (a) face image and 10 marked evidence classes; (b) pdf value surface for the left (in image) nostril class; (c) confidence threshold 0.5 ($F(\tau) = 0.5$); (d) confidence threshold 0.95 ($F(\tau) = 0.05$).

Lindh et al. [19] have proposed a statistical solution for detecting motor bearing failures from stator current signals of electric motors, but they criticized the problem that not only normal condition but also failure condition measurements are needed in their solution. In practice, providing failure condition measurements is often economically impossible and in general it is difficult to cover all failure situations while normal condition data is produced continuously. This experiment used the same features as proposed in the original study [19], but only normal condition measurements were used to form a pdf and confidence was used to decide between normal and failure conditions. In Fig. 9 a ROC curve is shown for the conducted experiment. From the curve it can be seen how by decreasing the confidence more normal condition measurements were correctly identified (true positives), but also an increasing number of failure conditions were considered as normal (false positives). The optimal trade-off depends on application. On the other hand, there was only a minor difference comparing to the result where also failure condition pdf was used in Bayesian classification (a priories were estimated from the training set, which does not correspond to the situation in practice).

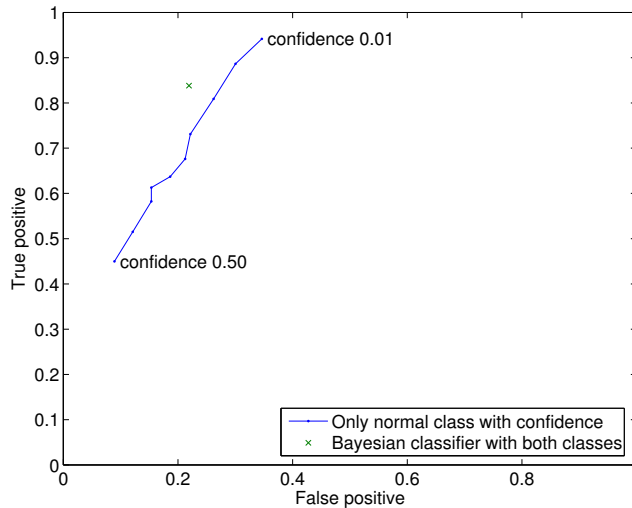


Figure 9: Receiver operating characteristic (ROC) curve for using confidence value in two class classification of electric motors. Positive test result corresponds to normal operating condition of a motor.

4 Conclusions

Our motivation for this study was the reason that statistical methods, especially methods based on class conditional probability density functions of features, are suitable in certain classification tasks. The study was aimed to be useful in feature classification and evidence based object detection, which was not considered in this study but only briefly demonstrated in the experiments.

The estimation of unknown pdf's is a general problem and in this study Gaussian mixture models (GMM) were analyzed as a model and expectation maximization (EM), Figueiredo-Jain (FJ) and greedy EM (GEM) algorithms were studied for estimating GMM parameters. The study covered also estimation of complex data in multiple dimensions in which case EM and FJ have been used. It seems that EM is the most suitable if the number of components is known or can be reliably estimated. The optimal number of components can also be estimated using a separate evaluation set and by pruning. However, in this sense the two unsupervised methods FJ and greedy EM also seemed to be robust and reliable algorithms for GMM parameter estimation. In practice singularities however occur, and thus, algorithms must tolerate that for example by enforcing covariance matrices to non-singular.

Estimated pdf's can be used in classification and the second part of this study concentrated on three different classification problems and in addition to what kind of information is needed in the next levels of processing. The first problem was solved in a traditional way by utilizing Bayes formula, but the next two problems required non-Bayesian approach. The use of confidence, as defined in this study, was proposed and a method to establish confidence information was introduced. Finally the use of confidence information was demonstrated by experiments.

This work represents a more broader view of using feature pdf's in a classification and ranking of extracted features, observations. The work consisted of a study how pdf's can be estimated and a study how the classification should be performed. In future, the results will be applied to an invariant detection of object evidence by utilizing simple Gabor features and statistical ranking as proposed.

References

- [1] GMMBAYES Matlab Toolbox. <http://www.it.lut.fi/project/gmmbayes>.
- [2] J. Bilmes. A gentle tutorial on the EM algorithm and its application to parameter estimation for Gaussian mixture and hidden Markov models, 1997.
- [3] C.M. Bishop. *Neural Networks for Pattern Recognition*. Oxford Univ. Press, 1995.
- [4] L. Breiman, J.H. Friedman, R.A. Olshen, and C.J. Stone. *Classification and Regression Trees*, pages 43–49. Wadsworth International Group, Belmont, California, 1984.
- [5] N. Cristianini and J. Shawe-Taylor. *An introduction to Support Vector Machines and other Kernel-based learning methods*. Cambridge Univ. Press, 2000.
- [6] R.O. Duda, P.E. Hart, and D.G. Stork. *Pattern Classification* John Wiley & Sons, Inc., 2nd edition, 2001.
- [7] B.S. Everitt and D.J. Hand. *Finite Mixture Distributions*. Monographs on Applied Probability and Statistics. Chapman and Hall, 1981.
- [8] M.A.T. Figueiredo and A.K. Jain. Unsupervised learning of finite mixture models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(3):381–396, Mar 2002.

- [9] Y. Freund and R. E. Schapire. A decision theoretic generalization of on-line learning and an application to boosting. *J. of Computer and System Sciences*, 55(1):119–139, 1997.
- [10] N.R. Goodman. Statistical analysis based on a certain multivariate complex Gaussian distribution (an introduction). *The Annals of Mathematical Statistics*, 34(1):152–177, March 1963.
- [11] M. Hamouz, J. Kittler, J.-K. Kamarainen, P. Paalanen, and H. Kälviäinen. Affine-invariant face detection and localization using GMM-based feature detector and enhanced appearance model. In *Proc. of the 6th Int. Conf. on Automatic Face and Gesture Recognition*, pages 67–72, Seoul, Korea, 2004.
- [12] M. Hamouz, J. Kittler, J.-K. Kamarainen, P. Paalanen, H. Kälviäinen, and J. Matas. Feature-based affine-invariant detection and localization of faces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*. to be published.
- [13] R. J. Hyndman. Computing and graphing highest density regions. *The American Statistician*, 50(2):120–126, May 1996.
- [14] J. Ilonen, J.-K. Kamarainen, H. Kälviäinen, and O. Anttalainen. Automatic detection and recognition of hazardous chemical agents. In *Proc. of the 14th Int. Conf. on Digital Signal Processing*, volume 2, pages 1345–1348, 2002.
- [15] T. Jaaskelainen, R. Silvennoinen, J. Hiltunen, and J. P. S. Parkkinen. Classification of the reflectance spectra of pine, spruce, and birch. *Applied Optics*, 33(12):2356–2362, April 1994.
- [16] J.-K. Kamarainen, V. Kyrki, J. Ilonen, and H. Kälviäinen. Improving similarity measures of histograms using smoothing projections. *Pattern Recogn. Lett.*, 24(12):2009–2019, 2003.
- [17] Joni-Kristian Kämäräinen. *Feature Extraction Using Gabor Filters*. PhD thesis, Lappeenranta University of Technology, 2003.
- [18] V. Kyrki, J.-K. Kamarainen, and H. Kälviäinen. Simple Gabor feature space for invariant object recognition. *Pattern Recognition Letters*, 25(3):311–318, 2003.
- [19] T. Lindh, J. Ahola, J.-K. Kamarainen, V. Kyrki, and J. Partanen. Bearing damage detection based on statistical discrimination of stator current. In *Proc. of the 4th IEEE Int. Symp. on Diagnostics for Electric Machines, Power Electronics and Drives*, pages 177–181, Atlanta, Georgia, USA, 2003.
- [20] I. Nagy, P. Nedoma, and M. Karny. Factorized EM algorithm for mixture estimation. In *Proc. of the Int. Conf. on Artificial Neural Networks and Genetic Algorithms*, Wien, 2001.
- [21] Robert J. Schalkoff. *Pattern Recognition: Statistical, Structural and Neural Approaches*. John Wiley & Sons, Inc., 1991.
- [22] M.I. Schlesinger and V. Hlavac. *Ten Lectures on Statistical and Structural Pattern Recognition*. Computational Imaging and Vision. Kluwer Academic Publishers, 2002.

- [23] P. N. Somerville. Numerical computation of multivariate normal and multivariate-t probabilities over convex regions. *J. of Computational and Graphical Statistics*, 7(4):529–544, 1998.
- [24] Paul N. Somerville. Numerical computation of multivariate normal and multivariate t probabilities over ellipsoidal regions. *Journal of Statistical Software*, 2001.
- [25] Garry J. Tee. Surface area and capacity of ellipsoids in n dimensions. Technical report, Department of Mathematics, University of Auckland, 2004.
- [26] S. Theodoridis and K. Koutroumbas. *Pattern Recognition*. Academic Press, 1999.
- [27] Y.L. Tong. *The Multivariate Normal Distribution*. Springer Series in Statistics. Springer-Verlag, 1990.
- [28] J. J. Verbeek, N. Vlassis, and B. Kröse. Efficient greedy learning of Gaussian mixture models. *Neural Computation*, 5(2):469–485, Feb 2003.