

Maximum Certainty Data Partitioning

Stephen J. Roberts, Richard Everson & Ieab Rezek

Intelligent & Interactive Systems Group
Department of Electrical & Electronic Engineering
Imperial College of Science, Technology & Medicine
Exhibition Road, London SW7 2BT, UK
s.j.roberts@ic.ac.uk

March 10, 1999

Abstract

Problems in data analysis often require the unsupervised partitioning of a data set into clusters. Many methods exist for such partitioning but most have the weakness of being model-based (most assuming hyper-ellipsoidal clusters) or computationally infeasible in anything more than a 3-dimensional data space. We re-consider the notion of cluster analysis in information-theoretic terms and show that minimisation of partition entropy can be used to estimate the number and structure of probable data generators.

Keywords: Cluster analysis, data partitioning, information theory.

1 Introduction

Many problems in data analysis, especially in signal and image processing, require the *unsupervised* partitioning of data into a set of ‘self-similar’ clusters or regions. An ideal partition unambiguously assigns each datum to a single cluster and one thinks of the data as being generated by a number of data generators, one for each cluster. Many algorithms have been proposed for such analysis and for the estimation of the optimal number of partitions. The majority of popular and computationally feasible techniques rely on the assumption that clusters are hyper-ellipsoidal in shape. In the case of Gaussian mixture modelling [1, 2, 3] this is explicit; in the case of dendrogram linkage methods (which typically rely on the L_2 norm) it is implicit [4]. For some data sets this leads to an over-partitioning. Alternative methods, based on valley seeking [2] or maxima-tracking in scale-space [5] for example, have the advantage that they are free from such assumptions. They can be, however, computationally intensive, sensitive to noise (in the case of valley seeking approaches) and unfeasible in high-dimensional spaces (indeed these methods can become prohibitive in even a 3-dimensional data space).

In this paper we re-consider the issue of data partitioning from an information-theoretic viewpoint and show that minimisation of entropy, or maximisation of partition certainty, may be used to evaluate the most probable set of data generators. The approach does not assume cluster convexity, it is shown to partition a range of data structures and to be computationally efficient.

2 Theory

The idea underlying this approach is that the observed dataset is generated by a number of data generators (classes). We first model the unconditional probability density function (pdf) of the data and then seek a number of partitions whose linear combination yields the data pdf. Densities and classifications conditioned on this partition set are then easily obtained.

2.1 Information maximisation

Consider a set of K partitions. The probability density function of a single datum \mathbf{x} , conditioned on this partition set, is given by:

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

We consider the overlap between the contribution to this density function of the k -th partition and the density $p(\mathbf{x})$. This overlap may be measured by the Kullback-Liebler measure between these two distributions. The latter is defined, for distributions $p(x)$ and $q(x)$, as:

$$KL(p(x) \parallel q(x)) = \int p(x) \ln \left(\frac{p(x)}{q(x)} \right) dx \quad (2)$$

Note that this measure reaches a minimum of zero if, and only if, $p(x) = q(x)$. For any other case it is strictly positive and increases as the overlap between the two distributions *decreases*. What we desire, therefore, is that the KL measure be maximised as this implies that the overlap between two distributions are minimised. We hence write our overlap measure as:

$$v_k = -KL(p(\mathbf{x} | k)p(k) \parallel p(\mathbf{x})) \quad (3)$$

As the this measure is strictly non-positive we may define a total overlap as the summation of all v_k :

$$\begin{aligned} V &= -\sum_k KL(p(\mathbf{x} | k)p(k) \parallel p(\mathbf{x})) \\ &= -\sum_k \int p(\mathbf{x} | k)p(k) \ln \left(\frac{p(\mathbf{x} | k)p(k)}{p(\mathbf{x})} \right) d\mathbf{x} \end{aligned} \quad (4)$$

We note, furthermore, that as $V \leq 0$, so minimisation over all data is equivalent to minimisation of V for each datum. An ‘ideal’ data partitioning separates the data such that overlap between partitions is minimal. We therefore seek the partitioning for which V is a minimum. By Bayes’ theorem we may re-write equation 4 as

$$\begin{aligned} V &= -\sum_k \int p(k | \mathbf{x})p(\mathbf{x}) \ln p(k | \mathbf{x}) d\mathbf{x} \\ &= -\int \left(\sum_k p(k | \mathbf{x}) \ln p(k | \mathbf{x}) \right) p(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (5)$$

Note that the summation term is simply the Shannon entropy, given datum \mathbf{x} , over the set of partition posteriors, i.e. $H(\mathbf{x}) = -\sum_k p(k | \mathbf{x}) \ln p(k | \mathbf{x})$. Minimising V is hence equivalent to minimising the expected (sample) entropy of the partitions over all observed data. It is this objective which we will use to form minimum-entropy, or maximum certainty, partitions. It is noted that this is achieved by having, for each datum, some partition posterior close to unity, while all the others are close to zero, which conforms to our objective for ideal partitioning.

2.2 Mixture models

2.2.1 Kernel-based density estimators

We restrict ourselves in this paper to considering a set of kernels or basis functions which model the probability density function (pdf) of the data and thence of each data partition. It is worth at this point considering the major approaches to estimation of a density function using a finite set of kernel functions.

1. Parametric representation: in this case a strong assumption is made regarding a model for the data generators, namely that a single kernel represents each data generator. The Gaussian kernel is a popular choice, leading to Gaussian mixture modelling.
2. Semi-parametric representation: arbitrary density functions may be adequately represented (i.e. to any finite precision) with a finite number of kernels (if the kernel is chosen to be a function with universal approximation properties). The density function of each data generator is thus represented by a finite mixture of kernels (typically Gaussians).
3. Non-parametric representation: in this case each datum serves as the prototype (normally the location) for a single kernel function. As with semi-parametric representations, if the kernels are chosen so as to have universal approximation properties, then arbitrary density functions may be approximated.

Further details of density estimation approaches may be found in [6]. As we wish to decompose overly-complex models of data generation into simpler partitionings we require either semi- or non-parametric models of the data pdf. In all the examples presented in this paper we find little difference between partitioning results using non- and semi-parametric estimators, and we have chosen the latter (with a heuristically-chosen ‘moderate’ number of kernels) in all examples. This, clearly, gives a computational advantage although this is arguably offset by having *any* heuristic values in our analysis (which is, in principle, free from any) and by the necessity of optimising the kernel set. We perform the latter using the standard EM (expectation-maximisation) algorithm (see [7], for example). It is worth commenting that the EM algorithm has a well-known failure mode, namely that a single kernel can reduce its variance to zero and hence maximise the data likelihood to infinity. This situation arises when the mixture model is over-complex given the available data. We avoid this by careful initialisation of the EM algorithm using K-means clustering and by early stopping [7]. The latter performs a natural regularisation of the system and we observe (empirically) no over-fitting problems on all the data presented in this paper.

2.2.2 Partitions as mixture models

We consider a set of partitions of the data. We may model the density function of the data, conditioned on the i -th partition, via a semi-parametric mixture model of the form:

$$p(\mathbf{x} | i) = \sum_{j=1}^J p^*(\mathbf{x} | j) \pi^*(j) \quad (6)$$

where J is the number of kernels forming the mixture and π_j^* are a set of (unknown) mixture coefficients which sum to unity. Each mixture component may be, for example, a Gaussian kernel and hence each candidate partition of the data is represented in this approach as a different mixture of these kernels. The use of the ‘star’ notation, i.e. $p^*(j)$, denotes that this set of probabilities is evaluated over the kernel representation, rather than over the set of data partitions. Equation 6 may be written, via Bayes’ theorem, as a linear transformation to a set of partition *posteriors* of the form:

$$\mathbf{p} = \mathbf{W} \mathbf{p}^* \quad (7)$$

where \mathbf{p} is the set of partition posterior probabilities (in vector form), \mathbf{W} is some transform, or mixing, matrix (not assumed to be square) and \mathbf{p}^* is the set of kernel posteriors (in vector form). Hence the i -th partition posterior may be written as:

$$p_i = p(i | \mathbf{x}) = \sum_j W_{ij} p^*(j | \mathbf{x}) \quad (8)$$

If we are to interpret \mathbf{p} as a set of posterior probabilities we require:

$$p_i \in [0, 1] \forall i \text{ and } \sum_i p_i = 1 \quad (9)$$

As $p^*(j | \mathbf{x}) \in [0, 1]$ so the first of these conditions is seen to be met if each $W_{ij} \in [0, 1]$ (as $\sum_i a_i b_i \in [0, 1]$ if $a_i, b_i \in [0, 1]$). The second condition is met when

$$\begin{aligned} 1 &= \sum_i p(i | \mathbf{x}) = \sum_i \sum_j W_{ij} p^*(j | \mathbf{x}) \\ &= \sum_j p^*(j | \mathbf{x}) \sum_i W_{ij} = \sum_i W_{ij} \end{aligned} \quad (10)$$

i.e. each column of \mathbf{W} sums to unity.

Given a set of partition posteriors, the partition priors may be re-evaluated as their maximum-likelihood estimates, i.e.

$$p(i) = \langle p(i | \mathbf{x}) \rangle = \frac{1}{N} \sum_n p(i | \mathbf{x}_n) \quad (11)$$

hence we may also form partition (class) conditional likelihoods via Bayes' theorem, i.e.

$$p(\mathbf{x} | i) = \frac{p(i | \mathbf{x}) p(\mathbf{x})}{p(i)} \quad (12)$$

We note once more that each partition-conditional density in the transformed space is represented by a mixture of kernels from the original space. If 'hard' partitioning is required, each datum is assigned to the partition with the largest posterior in the transformed space.

Centroid locations of each partition are simply obtained if required (and believed to be meaningful; the centroid of a ring structure lies in a region of no data density, for example). The i -th such point, \mathbf{m}_i , which we may consider as the centroid of the i -th data partition, is estimated by its maximum likelihood value namely:

$$\mathbf{m}_i = \frac{\sum_n \mathbf{x}_n p(i | \mathbf{x}_n)}{\sum_n p(i | \mathbf{x}_n)} \quad (13)$$

2.3 Entropy minimisation

Given that we represent each partition via a fixed set of kernels, we wish to adjust the elements of the matrix \mathbf{W} such that the entropy over the partition posteriors is minimised. We must also, however, take into account the constraints on the elements of \mathbf{W} (that they are bounded in $[0, 1]$ and the sum of each column of \mathbf{W} is unity). We may achieve this by introducing a set of dummy variables, which will be optimised, such that \mathbf{W} is represented by a generalised logistic function (the so-called 'softmax' function) of the form:

$$W_{ij} = \frac{\exp(\theta_{ij})}{\sum_{i'} \exp(\theta_{i'j})} \quad (14)$$

The gradient of the entropy with respect to each dummy variable, θ_{ij} , is given via the chain rule as

$$\frac{\partial H}{\partial \theta_{ij}} = \sum_{i'} \frac{\partial H}{\partial W_{i'j}} \cdot \frac{\partial W_{i'j}}{\partial \theta_{ij}} \quad (15)$$

The summation term is easily evaluated noting that

$$\frac{\partial W_{i'j}}{\partial \theta_{ij}} = W_{i'j} \delta_{i'i} - W_{i'j} W_{ij} \quad (16)$$

where $\delta_{i'i} = 1$ if $i = i'$ and zero otherwise. The term $\partial H / \partial W_{i'j}$ is evaluated by writing the expectation of the entropy (of Equation 5) as a sample mean over all N data, i.e.

$$\frac{\partial H}{\partial W_{i'j}} = \frac{1}{N} \sum_n \frac{\partial H(\mathbf{x}_n)}{\partial p(i' | \mathbf{x}_n)} \cdot \frac{\partial p(i' | \mathbf{x}_n)}{\partial W_{i'j}} \quad (17)$$

As $p(i' | \mathbf{x}_n) = \sum_j W_{i'j} p^*(j | \mathbf{x}_n)$ the above is easily evaluated. In all the experiments reported in this paper we optimise \mathbf{W} using the above formalism via the BFGS quasi-Newton method [8].

2.4 Model-order estimation

Since the number of partitions is not known *a priori* it is useful to be able to discover the most probable number of partitions. To this end we evaluate the entropy change, per partition, as a result of observing the data set, X . This quantity is given as,

$$\Delta H(M_K | X) = H(M_K) - H(M_K | X) \quad (18)$$

where M_K is the K -partitions model. The first term on the right-hand side of the above Equation is simply the entropy of the model priors *before* data is observed and is the Shannon entropy taking the partition probabilities to be uniform and equal to $1/K$. The second term is the entropy associated with the posterior partition probabilities having observed X . It is noted that the prior entropy is constant for any M_K and hence our objective of minimising the posterior entropy will result in a maximum of $\Delta H(M_K | X)$ at the most-probable partition number. Noting that $H(X) - H(X | M_K) = H(M_K) - H(M_K | X)$ and that $H(X | M_K)$ is the expectation of the negative log-likelihood of X given M_K so the likelihood (evidence) of X given M_K may be written as:

$$p(X | M_K) \propto \exp \{ \Delta H(M_K | X) \} \quad (19)$$

in which the data entropy term, $H(X)$, is ignored as it is constant for all models. Choosing the model with the largest value of this likelihood is equivalent, via Bayes' theorem, to choosing the model with the highest probability, $p(M_K | X)$ if we assume flat prior beliefs, $p(M_K)$ for each model. We hence obtain a posterior belief measure for each candidate partitioning:

$$p(M_K | X) = \frac{\exp \{ \Delta H(M_K | X) \}}{\sum_{K'} \exp \{ \Delta H(M_{K'} | X) \}} \quad (20)$$

and it is this measure which we use to assess the model order, choosing the order K for which it is maximal.

3 Results

3.1 Simple data set

We first present results in detail from a data set in which clusters are simple and distinct; the data are generated from four Gaussian distributed sources with 30 data drawn from each. Each component has the same (spherical) covariance. These data are shown in Figure 1. As an illustration of a simple kernel set, ten Gaussian components are fitted to the data using the EM algorithm. The resultant set of posterior probabilities, $p^*(j | \mathbf{x}_n)$, is shown in Figure 2. Figure 3(a) shows $\ln p(M_K | X)$ and plot (b) $p(M_K | X)$. Note that a set of four partitions are clearly favoured. Choosing the $K = 4$ model we obtain, for this example, \mathbf{W} as a 4×10 matrix. The set of four partition probabilities, $p(i | \mathbf{x}_n)$, is shown in Figure 4. The resultant partitioning of the data set gives the results of Figure 5. There are no errors in the partitioning for this simple data set.

3.2 Ring data

The next (synthetic) data set we investigate is drawn from two generator distributions; an isotropic Gaussian and a uniform 'ring' distribution. A total of 100 data points were drawn from each distribution (hence $N = 200$). A 20-kernel Gaussian mixture model was fitted to the data (using again the EM algorithm). Figure 6(a) shows that $p(M_K | X)$ gives greatest support for the two partition model. Plot (b) of the same figure depicts the resultant data partitioning. For this example there are no errors. Note that, due to the pathological structure of this example, a Gaussian mixture model *per se* fails to estimate the 'correct' number of partitions and provide a reasonable data clustering.

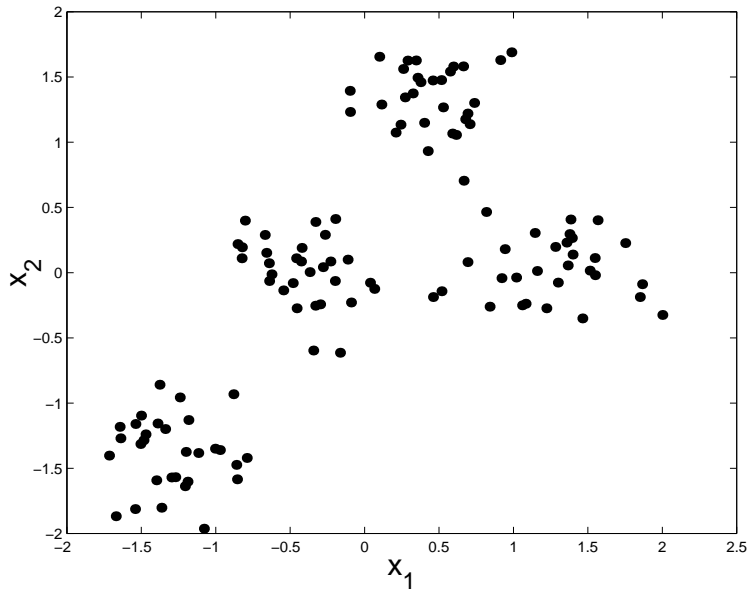


Figure 1: **Simple data set:** 30 points are drawn from each of four well-separated Gaussian sources.

3.3 Iris data

Anderson’s ‘iris’ data set is well-known [9]. The data we analysed consisted of 50 samples for each of the three classes present in the data, *Iris Versicolor*, *Iris Virginica* and *Iris Setosa*. Each datum is four-dimensional and consists of measures of the plants morphology. Once more a 20-kernel model was fitted to the data set. Figure 7(a) shows the model-order measure, shown in this case on a linear y-scale. Although support is greatest for (correct) the $K = 3$ partitioning it is clear that a two-partition model has support. We regard this as sensible given the nature of the data set, i.e. it naturally splits into two partitions. As in previous figures plot(b) depicts the data partitioning. This plot shows the projection onto the first two principal components of the data set. The partitioning has three errors in 150 samples giving an accuracy of 98%. This is slightly better than that quoted in [3] and the same as that presented for Bayesian Gaussian mixture models in [1].

3.4 Wine recognition data

As a final example we present results from a wine recognition problem. The data set consists of 178 13-dimensional exemplars which are a set of chemical analyses of three types of wine. Once more we fit a 20-kernel model and perform a minimum-entropy clustering. Figure 8(a) shows $\ln P(M_K | X)$. There is a clear maximum at the ‘correct’ partitioning ($K = 3$). Plot (b) shows this partitioning projected onto the first two components of the data set. For this example there are 4 errors, corresponding to an equivalent classification performance of 97.75%. This data set has not (to the authors’ knowledge) been analysed using an *unsupervised* classifier, but *supervised* analysis has been reported. Our result is surprisingly good considering that *supervised* first-nearest neighbour classification achieves only 96.1%, and multivariate linear-discriminant analysis 98.9% [10]. It should be commented that the same partitioning is obtained via analysis of the first two data principal components *alone*, rather than the full 13-D data set.

4 Conclusions

We have presented a computationally simple technique for data partitioning based on a linear mixing of a set of fixed kernels. The technique is shown to give excellent results on a range of problems. For computational parsimony we have used an initial semi-parametric approach to kernel fitting although, as mentioned, the results from a non-parametric analysis are near identical in all cases.

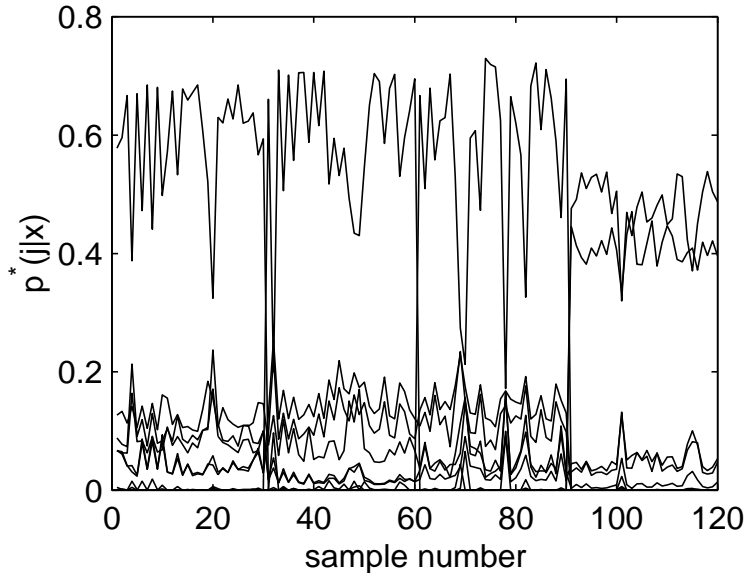


Figure 2: **Simple data set:** the plot shows the posteriors $p^*(j | \mathbf{x})$ from EM fitting of the set of ten kernels.

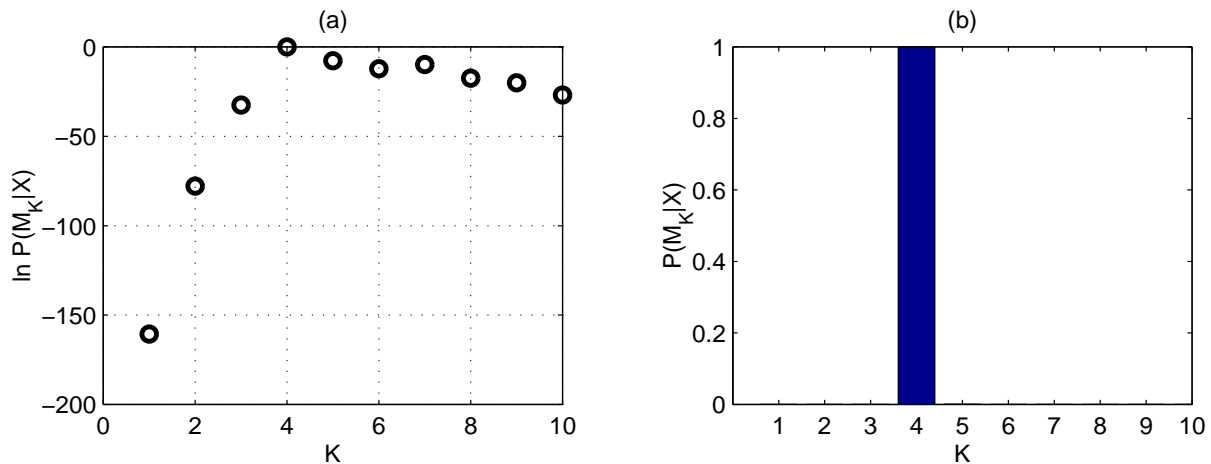


Figure 3: **Simple data set:** (a) $\ln p(M_K | X)$ and (b) $p(M_K | X)$. Note the clear maxima at $K = 4$.

The methodology is general and non-Gaussian kernels may be employed in which case the estimated partition-conditional densities will be mixture models of the chosen kernel functions. The method, furthermore, scales favourably with the dimensionality of the data space and the entropy-minimisation algorithm is efficient even with large numbers of samples.

5 Acknowledgements

IR and RE are funded respectively via grants from the commission of the European Community (project SIESTA, grant BMH4-CT97-2040) and British Aerospace plc. whose support we gratefully acknowledge. The *iris* and *wine* data sets are available from the UCI machine-learning repository. The authors would also like to thank the anonymous reviewers of this paper for insightful comments and suggestions.

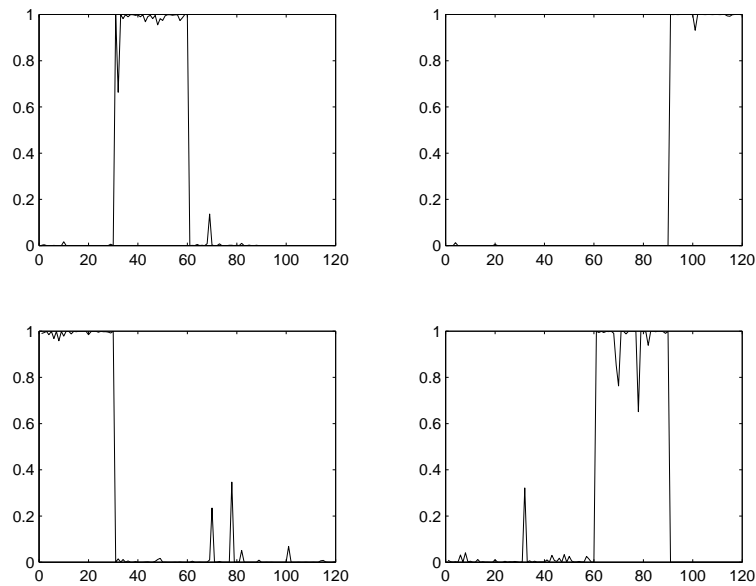


Figure 4: **Simple data set:** The four posterior probabilities, $p(i | \mathbf{x})$, for the most likely partitioning.

References

- [1] S.J. Roberts, D. Husmeier, I. Rezek, and W. Penny. Bayesian approaches to mixture modelling. *IEEE Transaction on Pattern Analysis and Machine Intelligence*, 20(11):1133–1142, 1998.
- [2] K. Fukunaga. *An Introduction to Statistical Pattern Recognition*. Academic Press, 1990.
- [3] I. Gath and B. Geva. Unsupervised Optimal Fuzzy clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(7):773–781, 1989.
- [4] A.K. Jain and R.C. Dubes. *Algorithms for Clustering Data*. Prentice Hall, 1988.
- [5] S.J. Roberts. Parametric and non-parametric unsupervised cluster analysis. *Pattern Recognition*, 30(2):261–272, 1997.
- [6] B.W. Silverman. *Density Estimation for Statistics and Data Analysis*. Number 26 in Monographs on statistics and applied probability. Chapman and Hall, London, 1986.
- [7] C.M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, Oxford, 1995.
- [8] W.H. Press, B.P. Flannery, S.A. Teukolsky, and W.T. Vetterling. *Numerical Recipes in C*. Cambridge University Press, 1991.
- [9] E. Anderson. The Irises of the Gaspé peninsula. *Bull. Amer. Iris Soc.*, 59:2–5, 1935.
- [10] S. Aeberhard, D. Coomans, and O. de Vel. Comparative-Analysis of Statistical Pattern-Recognition Methods in High-Dimensional Settings. *Pattern Recognition*, 27(8):1065–1077, 1994.

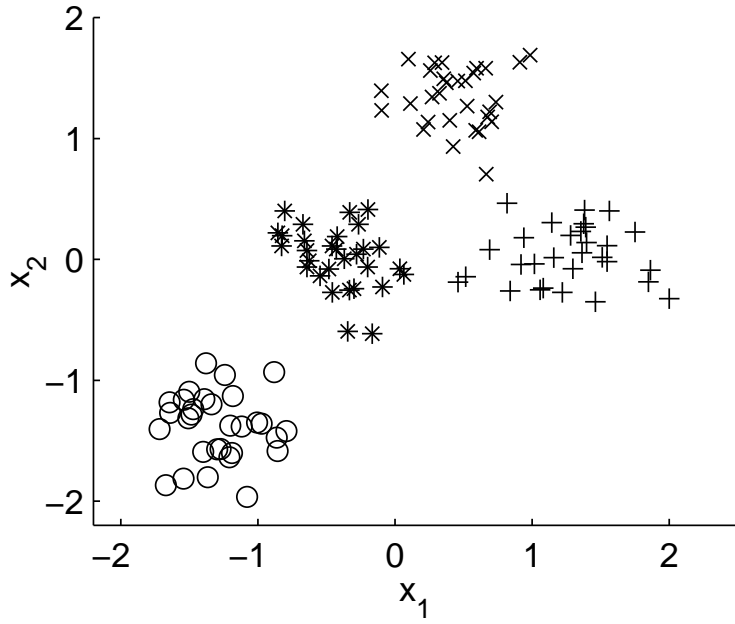


Figure 5: **Simple data set:** Data partitioning in the transformed space. For this simple example there are no errors.

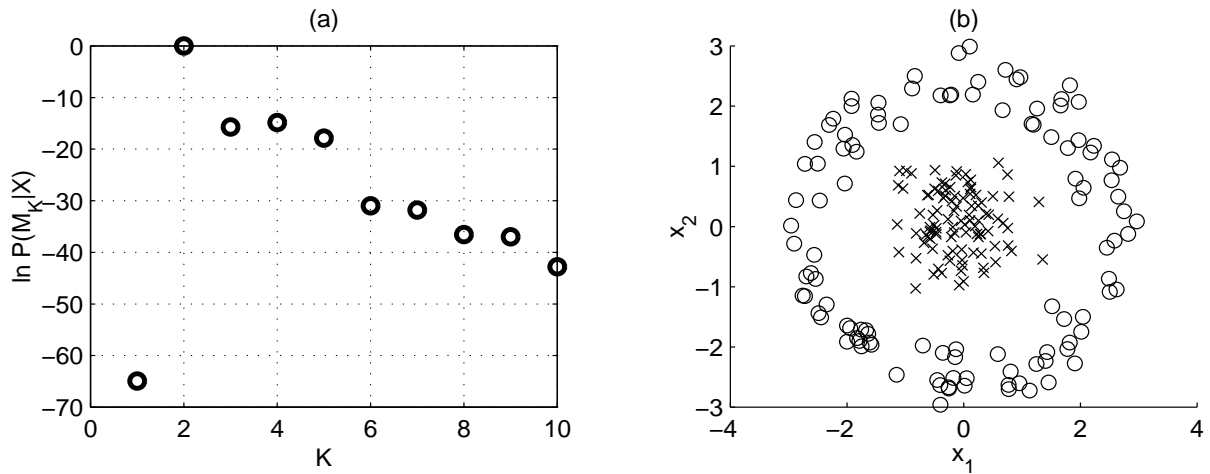


Figure 6: **Ring data set:** (a) $\ln P(M_K | X)$ & (b) resultant partitioning. For this example there are no errors.

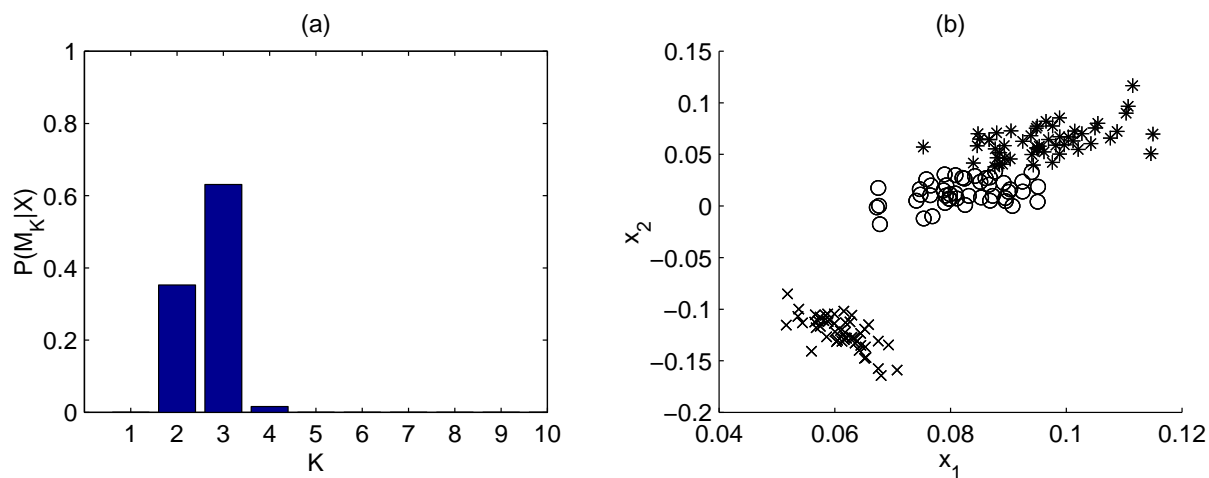


Figure 7: **Iris data set:** (a) $P(M_K | X)$ & (b) resultant partitioning. For this example there are three errors, corresponding to an accuracy of 98%.

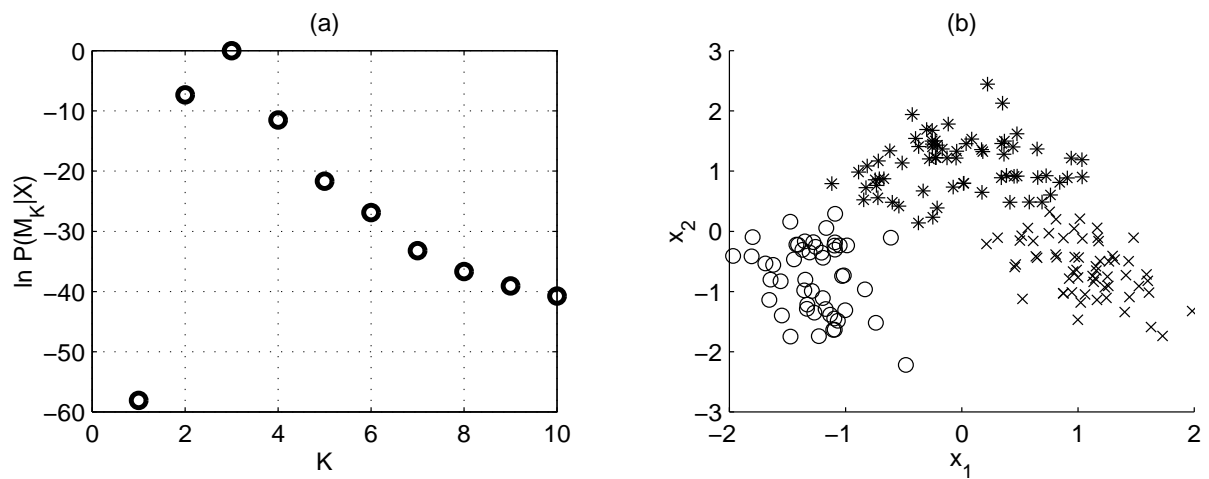


Figure 8: **Wine recognition data set:** (a) $\ln P(M_K | X)$ & (b) resultant partitioning. For this example there are four errors, corresponding to an accuracy of 97.75%.