

1 Algorithms for Maximum-Likelihood Bandwidth
2 Selection in Kernel Density Estimators

3 José M. Leiva-Murillo*

4 Antonio Artés-Rodríguez

5 *Dept. Signal Theory and Communication, Universidad Carlos III de Madrid,*
6 *Av. Universidad, 30, Leganés 28911 (Madrid), Spain.*
7 *Ph. +34 916248450, Fax +34 916248749*

8 **Abstract**

In machine learning and statistics, kernel density estimators are rarely used on multivariate data due to the difficulty of finding an appropriate kernel bandwidth to overcome overfitting. However, the recent advances on information-theoretic learning have revived the interest on these models. With this motivation, in this paper we revisit the classical statistical problem of data-driven bandwidth selection by cross-validation maximum likelihood for Gaussian kernels. We find a solution to the optimization problem under both the spherical and the general case where a full covariance matrix is considered for the kernel. The fixed-point algorithms proposed in this paper obtain the maximum likelihood bandwidth in few iterations, without performing an exhaustive bandwidth search, which is unfeasible in the multivariate case. The convergence of the methods proposed is proved. [A set of classification experiments are performed to prove the usefulness of the obtained models in pattern recognition.](#)

9 *Key words:* Kernel Density Estimation, Multivariate Density Modeling,

11 **1. Introduction**

12 Kernel Density Estimators (KDE) is a family of probability density func-
13 tion (PDF) models that have been intensively studied by the statistics com-
14 munity. Despite of being considered non-parametric, a *bandwidth* parameter
15 determines the scale of the kernel function and therefore the performance of
16 these models. Univariate KDEs have been object of much interest, whilst
17 multivariate models have been studied mainly in low dimensional cases. For
18 this reason, the interest of the pattern recognition and machine learning com-
19 munity in KDEs has been rather limited. For example, in the information-
20 theoretic learning literature, KDEs with a Gaussian kernel are often used
21 to estimate entropy or mutual information [18], [24], [15]. In most cases
22 the bandwidth of the kernel is treated as a hyperparameter chosen by cross
23 validation or a heuristic criterion.

24 The objective of this paper is to provide a method for the design of the
25 kernel bandwidth based on the maximum likelihood (ML) criterion, com-
26 monly applied in parametric and semi-parametric PDF modeling. Although
27 the idea of maximizing leave-one-out ML has been previously proposed in
28 the literature, an exhaustive search was needed to find the optimal band-
29 width. We avoid this inconvenience by providing fixed-point algorithms that
30 converge to the ML bandwidth even when it is characterized by a number of
31 parameters. The proof of the convergence in the general case is given by the
32 properties of the expectation-maximization (EM) algorithm [2]. Because the
33 spherical model is a particular case of the unconstrained case, its convergence

34 is automatically proved as well. However, we follow a different procedure for
 35 the convergence analysis of the spherical bandwidth, which provides us with
 36 important information about the range of values to which the bandwidth is
 37 guaranteed to belong.

38 In the next Section, we study the problem of PDF estimation by KDE
 39 models, and propose two novel and efficient algorithms for bandwidth se-
 40 lection for both the spherical and the unconstrained Gaussian kernels. In
 41 Section 3 we analyze the performance of KDEs with the obtained band-
 42 widths on both synthetic and benchmark real data. The paper finishes in
 43 Section 4 with some conclusions about the work presented.

44 2. Bandwidth Selection for Kernel Density Estimators

45 A Kernel Density Estimator (KDE) is a non-parametric PDF model that
 46 consists of a linear combination of kernel functions centered on the data (see,
 47 for example, [9])

$$\hat{p}_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N k(\mathbf{x} - \mathbf{x}_i | \boldsymbol{\theta}) \quad (1)$$

48 where $\mathbf{x} \in \mathbb{R}^D$ and $k(\mathbf{x} | \boldsymbol{\theta})$ is the kernel function with a given set of parameters
 49 $\boldsymbol{\theta}$. The kernel must be a unitary function, i.e.: $\int k(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = 1$. The success
 50 of these models in pattern recognition and machine learning is due to several
 51 reasons. First, a-priori assumptions on the distribution of the data need not
 52 be made. Secondly, the model does not need to be trained, as it only relies on
 53 the samples. Finally, it is easy to carry out transformations on the data, such
 54 as $\mathbf{z} = \mathbf{f}(\mathbf{x})$, and estimate the PDF in the \mathbf{z} -space, even if the transformation
 55 is not invertible, because $\hat{p}(\mathbf{z})$ is a KDE built from the transformed sample
 56 set $\{\mathbf{z}_i = \mathbf{f}(\mathbf{x}_i)\}$.

57 Although the KDEs are commonly considered as non-parametric models,
58 the kernel function has an adjustable bandwidth defined by $\boldsymbol{\theta}$ that determines
59 the accuracy of the model, so that it can be treated as a parameter to be
60 optimized.

61 The bandwidth selection problem has been a matter of intensive research
62 for the last 30 years in the statistics community. See [25] for an exhaustive re-
63 view of criteria for univariate data. The most extended criteria in bandwidth
64 estimation are the integrated square error (ISE), the mean ISE (MISE) and
65 the asymptotic MISE (AMISE) ¹. In the one-dimensional case, optimizing
66 these criteria with respect to the bandwidth is not problematic as it involves
67 a global search on one variable which is computationally feasible. This is
68 the main reason why multivariate bandwidth selection has been addressed
69 only in very low dimensional spaces [6]. According to the review in [25], the
70 main categories of methods for bandwidth selection are: rule-of-thumb meth-
71 ods, cross-validation methods and *plug-in* methods. When applied to criteria
72 based on asymptotic square error, these methods suffer from a serious draw-
73 back: they need to estimate not only the true density, but also its first and
74 second derivatives for a study based on Taylor expansion. As the true density
75 is unknown in general, the estimated one can be used instead by means of
76 a *pilot* bandwidth [26]. Although this can be acceptable in the univariate
77 case, in the multivariate case the inaccuracy may become dramatic. In the
78 case of rule-of-thumb methods, data are assumed Gaussian, which leads to
79 close-form expressions. Otherwise, the multivariate case has been addressed

¹Criteria based on L_1 -norm have been also proposed [4].

80 in the literature either by using a spherical kernel or by considering a scaled
 81 version of the covariance matrix of data as the kernel bandwidth [22], [20],
 82 being the scale factor the only parameter to adjust.

83 The Bayesian framework provides a natural way of obtaining the band-
 84 width by assigning it a proper prior and computing the posterior of the
 85 bandwidth given the data. However, two difficulties arise when working with
 86 KDEs. First, unlike in parametric modeling, the parameters of the band-
 87 width do not index any family of densities, since the centers of the kernels
 88 are needed. For this reason, a loss different than likelihood is considered [3],
 89 [7]. Secondly, the choice of the prior distribution determines the final choice
 90 of the bandwidth. The use of Monte Carlo Markov Chains (MCMC) has also
 91 been proposed [27].

92 In this paper, we aim at applying generative modeling to real-world ma-
 93 chine learning and pattern recognition problems, in which the dimension is
 94 higher than usually considered by the statistics community for kernel density
 95 estimation. We base our work on the maximum likelihood (ML) criterion.
 96 The risk of overfitting is inherent to ML because the complexity of the model
 97 is not considered. However, this risk can be alleviated by performing cross-
 98 validation, since the model is evaluated in points that are different from the
 99 ones used for the KDE model. The leave-one-out (LOO) is the extreme case
 100 in which $N - 1$ samples build a model evaluated on the point left

$$\hat{p}_{\boldsymbol{\theta}}(\mathbf{x}_i) = \frac{1}{N - 1} \sum_{\substack{j=1 \\ j \neq i}}^N G(\mathbf{x}_i - \mathbf{x}_j | \boldsymbol{\theta}) \quad (2)$$

101 where we make explicit the use of a Gaussian kernel. Note that a ML solution
 102 is equivalent to minimum entropy when the entropy is estimated as $\hat{h}(\mathbf{X}) =$

103 $-\frac{1}{N} \sum_i \log \hat{p}(\mathbf{x}_i)$. It has been proven in the literature that if $\hat{p}(\mathbf{x})$ is a KDE,
104 the entropy is overestimated [1]. Hence, a minimum entropy criterion for
105 bandwidth choice leads to a model that provides an estimation of the entropy
106 which is closest to the true value than any other performed with a KDE.

107 This procedure was first proposed in [5] and later studied by other au-
108 thors [22], [12]. The study of the multivariate case has been constrained
109 to problems in which the dimension is low, because in absence of a closed
110 optimization procedure, an exhaustive search of the optimal bandwidth is un-
111 feasible if the dimension of $\boldsymbol{\theta}$ is high. In this section, we present a procedure
112 for the bandwidth selection problem that overcomes these difficulties.

113 We will consider two different degrees of complexity for the covariance
114 matrix of the Gaussian kernel. In the most simple case, we assume a spher-
115 ical shape so that $\mathbf{C} = \sigma^2 \mathbf{I}_D$. In this case, $\boldsymbol{\theta}$ consists of just one parameter.
116 In the general or unconstrained case, no constraints are imposed to \mathbf{C} further
117 than its positive semi-definiteness; in this case, the number of elements in $\boldsymbol{\theta}$
118 is $(D-1)D/2$. In the following, we separately describe the bandwidth selec-
119 tion procedure for the spherical and the general case, providing two versions
120 of what we call the Maximum Likelihood Leave-One-Out (ML-LOO) algo-
121 rithm. The spherical model is a particular case of the unconstrained model,
122 so that the convergence in the former would follow from the convergence of
123 the latter. Although both cases can be studied through their connection to
124 the Expectation-Maximization algorithm, we choose to separately prove the
125 convergence in both cases. Specifically, in the spherical case we base our
126 proof in the general properties of fixed-point algorithms. This way, interest-
127 ing conclusions arise about the range of values where the optimal bandwidth

128 can be found. A preliminary version of this work was presented in [13].

129 *2.1. The spherical case*

130 The value of a kernel function centered in point \mathbf{x}_i and evaluated in \mathbf{x}_j
 131 is, for the spherical Gaussian case,

$$G_{ij}(\sigma^2) = G(\mathbf{x}_i - \mathbf{x}_j | \sigma^2) = (2\pi)^{-D/2} \sigma^{-D} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{x}_j\|^2\right).$$

132 The derivative with respect to σ is given by

$$\nabla_{\sigma} G_{ij}(\sigma^2) = \left(\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^3} - \frac{D}{\sigma}\right) G_{ij}(\sigma^2).$$

133 Let us consider now the log-likelihood of the data according to this model
 134 under the iid assumption: $\log L(\mathbf{X} | \sigma^2) = \sum_i \log \hat{p}_{\sigma}(\mathbf{x}_i)$, where $\hat{p}_{\sigma}(\mathbf{x}_i)$ is LOO
 135 estimated as in (2). The derivative of the LOO log-likelihood, is given by

$$\begin{aligned} \nabla_{\sigma} \log L(\mathbf{X} | \sigma^2) &= \sum_i \frac{1}{\hat{p}_{\sigma}(\mathbf{x}_i)} \frac{1}{N-1} \sum_{j \neq i} \frac{\partial}{\partial \sigma} G_{ij}(\sigma^2) \\ &= \frac{1}{N-1} \sum_i \frac{1}{\hat{p}_{\sigma}(\mathbf{x}_i)} \sum_{j \neq i} \left(\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^3} - \frac{D}{\sigma}\right) G_{ij}(\sigma^2). \end{aligned}$$

136

We now search for the maximum value of $\log L(\mathbf{X} | \sigma^2)$, and so the point that makes the derivative null. We have

$$\begin{aligned} \sum_i \frac{1}{\hat{p}_{\sigma}(\mathbf{x}_i)} \sum_{j \neq i} \frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{\sigma^3} G_{ij}(\sigma^2) &= \sum_i \frac{1}{\hat{p}_{\sigma}(\mathbf{x}_i)} \frac{D}{\sigma} \sum_{j \neq i} G_{ij}(\sigma^2) \\ &= \frac{N(N-1)D}{\sigma}. \end{aligned}$$

137 The second equality has been obtained by the fact that, by definition,
 138 $\sum_{j \neq i} G_{ij} = (N-1)\hat{p}_\sigma(\mathbf{x}_i)$. By isolating the σ^2 we obtain the following fixed-
 139 point rule

$$\sigma_{n+1}^2 = \frac{1}{N(N-1)D} \sum_i \frac{1}{\hat{p}_{\sigma_n}(\mathbf{x}_i)} \sum_{j \neq i} \|\mathbf{x}_i - \mathbf{x}_j\|^2 G_{ij}(\sigma_n^2). \quad (3)$$

140 We prove the convergence of the algorithm in (3) by the following Theo-
 141 rem:

Theorem 1. *There is a fixed point in the interval $\left(\frac{\overline{d_{NN}^2}}{D}, \frac{\overline{d^2}}{D}\right)$, being $\overline{d_{NN}^2}$ the mean quadratic distance to the nearest neighbor and $\overline{d^2}$ the mean distance among data points, so that $\left(\frac{\overline{d_{NN}^2}}{D} \leq \frac{\overline{d^2}}{D}\right)$. Besides, the fixed point is unique and the algorithm converges to it in the mentioned interval if the following condition holds*

$$\frac{1}{4\sigma^4 N(N-1)^2 D} \sum_i \frac{1}{\left(\sum_{l \neq i} \exp\left(-\frac{d_{il}^2}{2\sigma^2}\right)\right)^2} \times \sum_{j \neq i} \sum_{k \neq i, j} (d_{ij}^2 - d_{ik}^2)^2 \exp\left(-\frac{d_{ij}^2 + d_{ik}^2}{2\sigma^2}\right) < 1 \quad (4)$$

142 where $d_{ij}^2 = \|\mathbf{x}_i - \mathbf{x}_j\|^2$.

143 *Proof.* Let $\sigma^2 = g(\sigma^2)$ be the function in (3), whose fixed point is to be
 144 obtained. The proof of the fixed point existence is based on the search of
 145 an interval (a, b) such that $a < g(\sigma^2) < b$ if $\sigma^2 \in (a, b)$, which is called a
 146 *contractive map* [8].

147 In order to demonstrate that the interval $\left(\frac{\overline{d_{NN}^2}}{D}, \frac{\overline{d^2}}{D}\right)$ holds the property
 148 stated by the Theorem, we need to prove these three conditions:

149 1. $\lim_{\sigma^2 \rightarrow 0} g(\sigma^2) = \frac{\overline{d_{NN}^2}}{D}$.

- 150 2. $\lim_{\sigma^2 \rightarrow \infty} g(\sigma^2) = \frac{\overline{d^2}}{D}$.
- 151 3. $g(\sigma^2)$ is monotonic in the interval.

152 These conditions are graphically summarized in Fig. 1. This way we are
 153 guaranteed that the interval is a contractive map, and there is at least one
 154 crossing point between the function $g(\sigma^2)$ and the line $g(\sigma^2) = \sigma^2$.

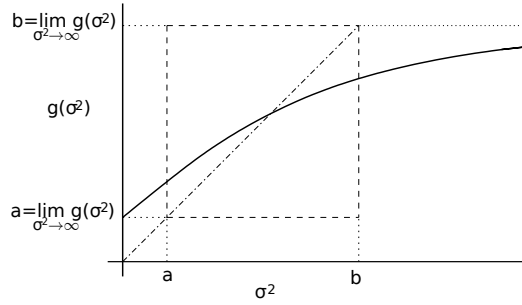


Figure 1: Contractive Map $\sigma^2 = g(\sigma^2)$

155 To prove the first point, we rewrite (3) as

$$g(\sigma^2) = \frac{1}{ND} \sum_i \sum_{j \neq i} d_{ij}^2 \frac{1}{1 + \sum_{k \neq i, j} \exp(\frac{d_{ij}^2 - d_{ik}^2}{2\sigma^2})}. \quad (5)$$

The limit at 0 is given by

$$\lim_{\sigma^2 \rightarrow 0} f(\sigma^2) = \frac{1}{ND} \sum_i \min_{j \neq i} d_{ij}^2 = \frac{\overline{d_{NN}^2}}{D}$$

156 because the elements in the denominator of (5) are null (the exponentials
 157 tend to infinite) in exception of the cases in which $d_{ij}^2 < d_{ik}^2, \forall k \neq j$, i.e. \mathbf{x}_j
 158 is the nearest neighbor of \mathbf{x}_i . The first condition is then proven.

To prove the second condition, we take the limit

$$\begin{aligned}\lim_{\sigma^2 \rightarrow \infty} g(\sigma^2) &= \lim_{\sigma^2 \rightarrow \infty} \frac{1}{ND} \sum_i \sum_{j \neq i} d_{ij}^2 \frac{\exp(-\frac{d_{ij}^2}{2\sigma^2})}{\sum_{k \neq i} \exp(-\frac{d_{ik}^2}{2\sigma^2})} \\ &= \frac{1}{ND} \sum_i \sum_{j \neq i} d_{ij}^2 \frac{1}{N-1} = \frac{\bar{d}^2}{D}.\end{aligned}$$

159 Then, the second condition is proven. The average distance can be effi-
160 ciently estimated from the covariance matrix of the data as $\bar{d}^2 = 2 \text{tr}\{\Sigma_x\}$,
161 by simple properties of linear algebra.

162 To demonstrate the last condition, we compute the derivative of $g'(\sigma^2)$
163 and check out that it is positive

$$\begin{aligned}\frac{dg(\sigma^2)}{d\sigma^2} &= \frac{1}{2\sigma^4 ND} \sum_i \sum_{j \neq i} d_{ij}^2 \frac{\sum_k (d_{ij}^2 - d_{ik}^2) \exp(-\frac{d_{ij}^2 + d_{ik}^2}{2\sigma^2})}{(\sum_{l \neq i} \exp(-\frac{d_{il}^2}{2\sigma^2}))^2} \\ &= \frac{1}{2\sigma^4 ND} \sum_i \frac{1}{(\sum_{l \neq i} \exp(-\frac{d_{il}^2}{2\sigma^2}))^2} \times \sum_{j \neq i} \sum_{k \neq i} d_{ij}^2 (d_{ij}^2 - d_{ik}^2) \exp(-\frac{d_{ij}^2 + d_{ik}^2}{2\sigma^2}) \\ &= \frac{1}{4\sigma^4 N(N-1)^2 D} \sum_i \frac{1}{(\sum_{l \neq i} \exp(-\frac{d_{il}^2}{2\sigma^2}))^2} \\ &\quad \times \sum_{j \neq i} \sum_{k \neq i, j} (d_{ij}^2 - d_{ik}^2)^2 \exp(-\frac{d_{ij}^2 + d_{ik}^2}{2\sigma^2}) \geq 0.\end{aligned}\tag{6}$$

164

165 In the last step, the terms with $j = k$ are removed, since they are null, and
166 the rest have been regrouped.

167 The existence of the fixed point is then proved. To demonstrate the con-
168 vergence of the algorithm in such interval, we need to check out the condition
169 $|g'(\sigma^2)| < 1$ [8]. In that case, we are guaranteed that only a crossing point
170 between $g(\sigma^2)$ and the line $g(\sigma^2) = \sigma^2$ exists. The convergence condition (4)
171 means that the value of the derivative obtained in (6) is lesser than 1. \square

172 Unfortunately, the computational complexity of obtaining the bound (4)
 173 is $O(N^3)$. In practice we have found the first condition as the most critical.
 174 The reason is that for quantized data it is possible that $\overline{d_{NN}^2}/D \approx 0$. It is easy
 175 to see from Fig. 1 that, in that case, the fixed-point location dangerously
 176 approaches to zero. In that case, the data should be treated as discrete
 177 instead, and a distribution model based on the histogram can be a better
 178 choice than a KDE.

179 The advantage of using the fixed-point iteration instead of an alternative
 180 optimization procedure is double. First, the practitioner does not have to
 181 worry about the optimization procedure for the log-likelihood, which is al-
 182 ready implicit in the fixed-point algorithm. Secondly, Theorem 1 provides
 183 us with a range of values where the optimal bandwidth is guaranteed to be
 184 found, which is of both theoretical and practical interest.

185 *2.2. The unconstrained case*

The general expression for a Gaussian kernel is

$$G_{ij}(\mathbf{C}) = |2\pi\mathbf{C}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{C}^{-1}(\mathbf{x}_i - \mathbf{x}_j)\right)$$

and its derivative w.r.t. \mathbf{C}

$$\nabla_{\mathbf{C}} G_{ij}(\mathbf{C}) = \frac{1}{2} (\mathbf{C}^{-1}(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T - \mathbf{I}) \mathbf{C}^{-1} G_{ij}(\mathbf{C}).$$

186 As in the previous cases, we take the derivative of the log-likelihood and
 187 make it equal to zero

$$\begin{aligned} \sum_i \frac{1}{\hat{p}_{\mathbf{c}}(\mathbf{x}_i)} \frac{1}{N-1} \sum_{j \neq i} \frac{1}{2} \mathbf{C}^{-1}(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T \mathbf{C}^{-1} G_{ij} \\ = \sum_i \frac{1}{\hat{p}_{\mathbf{c}}(\mathbf{x}_i)} \frac{1}{N-1} \sum_{j \neq i} \frac{1}{2} \mathbf{C}^{-1} G_{ij}. \end{aligned}$$

By multiplying both members by \mathbf{C} , we obtain

$$\sum_i \frac{1}{\hat{p}_{\mathbf{c}}(\mathbf{x}_i)} \sum_{j \neq i} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T G_{ij} = \mathbf{C} \sum_i \frac{1}{\hat{p}_{\mathbf{c}}(\mathbf{x}_i)} \sum_{j \neq i} G_{ij}.$$

188 After some simplifications as in the spherical case, we have

$$\sum_i \frac{1}{\hat{p}_{\mathbf{c}}(\mathbf{x}_i)} \sum_{j \neq i} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T G_{ij} = \mathbf{C}N(N-1)$$

189 leading to the following fixed-point rule

$$\mathbf{C}_{n+1} = \frac{1}{N(N-1)} \sum_i \frac{1}{\hat{p}_{\mathbf{c}_n}(\mathbf{x}_i)} \sum_{j \neq i} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T G_{ij}(\mathbf{C}_n). \quad (7)$$

190 The expression in (7) suggests a relationship with the expectation max-
 191 imization (EM) result for Gaussian mixture models (GMM). A GMM is a
 192 PDF estimator given by the expression $p(\mathbf{x}) = \sum_{k=1}^K \alpha_k G(\mathbf{x}|\boldsymbol{\mu}_k, \mathbf{C}_k)$. The
 193 weights of the K components of the mixture are given by the α_k , and each
 194 Gaussian is characterized by its mean vector $\boldsymbol{\mu}_k$ and its covariance matrix
 195 \mathbf{C}_k . The solution provided by the EM algorithm consists of an iterative pro-
 196 cedure where the values at step t are estimated from their corresponding
 197 values at step $t-1$. To do so, a matrix of auxiliary variables r_{ki} is used, each
 198 of which stands for the likelihood that the i -th sample has been generated
 199 by the k -th component of the mixture, with $\sum_k r_{ki} = 1$. The EM solution
 200 establishes the following updating rule for the covariance matrix at step t

$$\mathbf{C}_k^t = \sum_k \sum_i r_{ki}^t \frac{(\mathbf{x}_i - \boldsymbol{\mu}_k^t)(\mathbf{x}_i - \boldsymbol{\mu}_k^t)^T}{N} \quad (8)$$

201 where the r_{ki}^t and $\boldsymbol{\mu}_k^t$ are also iteratively updated. Note that our KDE model
 202 can be considered as a special case of GMM where there are as many mix-
 203 tures as samples ($K = N$) with the same weights ($\alpha_k = 1/N$) and fixed

204 mean vectors: $\boldsymbol{\mu}_k = \mathbf{x}_k$. The covariance matrix is the same for each of the
 205 components. Besides, the leave-one-out scheme imposes the values $r_{ki} = 0$ if
 206 $k = i$ and $r_{ki} = \frac{1}{N-1}$ if $k \neq i$. Then, the updating rule in (8) is equal to
 207 the one given by the iteration (7).

208 We conclude that the optimization for the unconstrained case can be set
 209 up as a particular case of the EM algorithm. The EM algorithm is known
 210 to monotonically increase the likelihood, so that its convergence to a local
 211 minimum has been proven in the literature [14]. According to this, the
 212 algorithm given in (7) is immediately shown to converge.

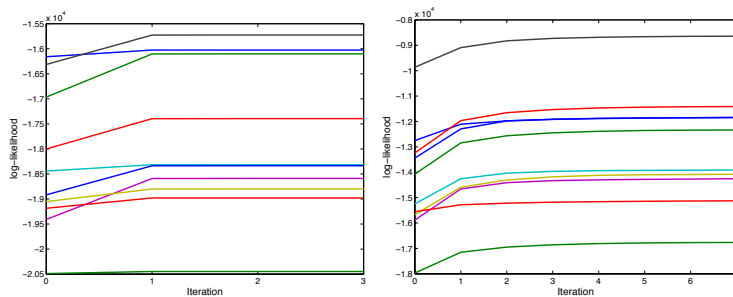


Figure 2: Log-likelihood values across iteration for the 10 different classes in Optdigits dataset. Left: spherical model; right: full (unconstrained) model.

213 In Fig. 2 we illustrate the convergence speed of both versions of ML-LOO
 214 on real data from the Optdigits dataset from the UCI repository. When the
 215 bandwidth is initialized with Scott’s rule [20], it only takes one iteration to
 216 get the optimal bandwidth when the spherical model is considered. In the
 217 unconstrained case, the likelihood values are higher than in the spherical one
 218 as expected, and the convergence is slower, but we note that no relevant
 219 increase in the log-likelihood is obtained beyond the 5th or 6th iteration.

220 A graphical example of the performance of both algorithms in (3) and (7)
 221 is displayed in Fig. 3 for a small set of synthetic data. Note that, in spite
 222 of using few samples to build the model, the proposed algorithm reaches
 223 solutions that are smooth yet descriptive. Besides, the unconstrained kernel
 224 is more flexible than the spherical one, in the sense that it is more accurate
 225 for distributions in which the covariance matrix of data is far from spherical.

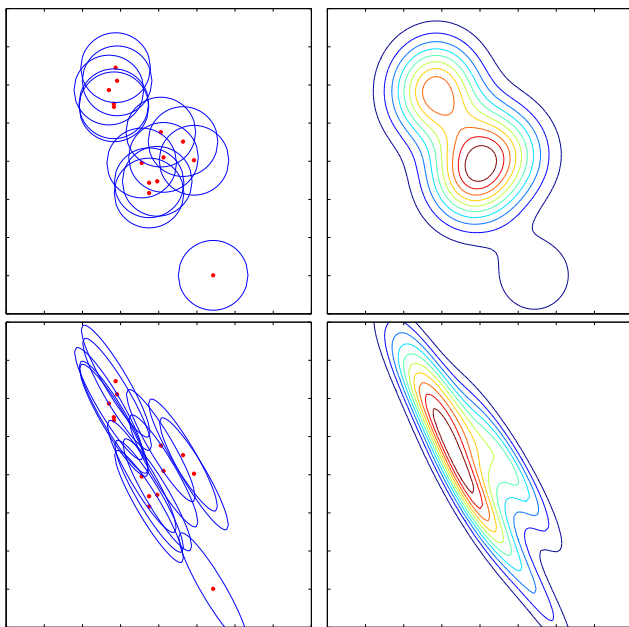


Figure 3: Examples of models obtained for the spherical (top) and the general (bottom) cases. Left: individual kernels, plotted at their -3dB level w.r.t. the mode; right: KDE models as the average of kernels on the left.

226 3. Experimental Results

227 Probability density estimation belongs to the family of unsupervised learn-
 228 ing methods. This means that it is not possible to compare different methods

229 unless the experiment has been designed on synthetic data whose distribution
230 is known. The estimated density is then compared to the true one according
231 to some criterion. Obviously, the criterion chosen for comparison benefits the
232 methods that have made use that same criterion for bandwidth selection.

233 As an alternative, we propose to apply the bandwidth selection criteria to
234 models that are used in supervised learning tasks. In particular, we consider
235 a set of classification or pattern recognition problems. In classification, the
236 general problem is to estimate the relationship between an input \mathbf{x} and an
237 output y from a dataset $\{\mathbf{x}_i, y_i\}$, $i = 1, \dots, L$, $\mathbf{x}_i \in \mathbb{R}^D$, $y \in \{1, 2, \dots, N_c\}$.
238 When the density of the data is modeled, classification is said to be *genera-*
239 *tive*. In particular, when KDEs are used for density modeling, the procedure
240 is referred to as Parzen classification.

241 We first describe the concept of Parzen classification, and then show the
242 performance of the classifier on both synthetic and real data.

243 3.1. Parzen Classification

244 A Parzen classifier assigns a sample \mathbf{x} a label according to a criterion
245 based either on the likelihood $p(\mathbf{x}|y)$ or the posterior $p(y|\mathbf{x})$, using the density
246 models estimated for each class. In the former case, the decision is taken
247 according to

$$\hat{y} = \arg \max_l \hat{p}(\mathbf{x}|c_l) \quad (9)$$

248 where $\hat{p}(\mathbf{x}|c_l)$ is a Parzen or KDE model built from the training data be-
249 longing to class l [9]. In the literature, generative models have rarely been
250 used in real world problems due to their high dimensionality, unless some

251 factorization is imposed on the variables to decompose the problem in low
252 dimensional tasks [16].

253 When choosing between a spherical or a full Gaussian kernel to perform
254 the classification rule (9), a tradeoff between advantages and weaknesses
255 of each option must be considered. In the spherical case, the model can
256 be inaccurate if the training data are distributed according to a strongly
257 non-spherical pattern. If the components of \mathbf{x} are very different in their
258 range or variance, taking the same bandwidth for each dimension may be
259 not reasonable. On the other hand, if we use a full (or unconstrained) model
260 we must be aware that the number of parameters involved can be very high
261 and lead to overfitting.

262 In order to overcome the disadvantages of both the spherical and the full
263 kernels, we propose a hybrid scheme that combines the advantages from both
264 approaches. The idea behind this procedure is to first whitening the data to
265 make the spherical approach appropriate per-class. Then, we evaluate the
266 value of the probability density back in the original feature space by means
267 of the linear transformation property of PDFs: $p(\mathbf{A}^T \mathbf{x}) = |\mathbf{A}|^{-1} p(\mathbf{x})$. Let \mathbf{B}_l
268 be the whitening matrix for each class, so that the components of $\mathbf{B}_l^T \mathbf{x}$ are
269 uncorrelated and normalized. The classification algorithm according to this
270 procedure is described by the algorithm in Fig. 4.

271 If the empirical covariance matrix of any of the classes is singular, we
272 can consider a non square matrix \mathbf{B}_l , with as many rows as non negative
273 eigenvalues of the covariance matrix. In order to compute the determinant
274 of \mathbf{B}_l as required by the algorithm, we propose to consider the heuristic
275 estimation $|\mathbf{B}| \approx \sqrt{|\mathbf{B}^T \mathbf{B}|}$.

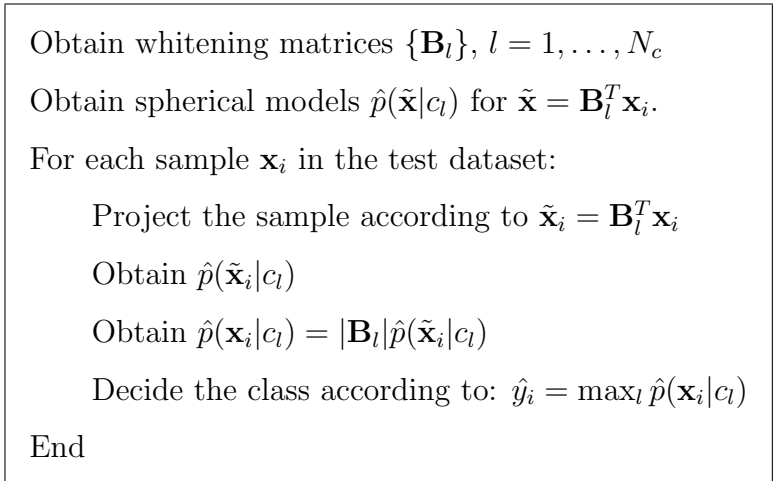


Figure 4: Hybrid method for Parzen classification.

276 Given that the three methods described have a different complexity, an
277 objective criterion is needed to choose a priori among the models. In addi-
278 tion to the classification performance, we have evaluated the leave-one-out
279 likelihood, i.e. the criterion according to which the models have been build.
280 Also, we have evaluated the Bayesian Information Criterion (BIC) [and the](#)
281 [Akaike Information Criterion \(AIC\)](#) to penalize the likelihood according to
282 the complexity of the model. Although originally developed for models of the
283 exponential family [21], BIC is usually applied to mixture models too [23].
284 The BIC suggests choosing the model for which the following expression is
285 maximum

$$BIC_j = \log L(\mathbf{X}) - \frac{K_j}{2} \log N$$

286 where j indexes the model, K_j is the number of parameters of the model and
287 N is the number of observations. In our case, the number of parameters is
288 $K_{sph} = K_{hyb} = 1$ in the spherical and hybrid models, and $K_{full} = D(D-1)/2$

289 for the unconstrained case, which is the number of elements in the bandwidth
290 matrix.

291 AIC, on the other hand, chooses the model based on a different penaliza-
292 tion of the likelihood. In particular, the *constrained* AIC or AIC^c accounts
293 for finite sample sizes

$$AIC_j^c = \log L(\mathbf{X}) - K_j - \frac{2K_j(K_j + 1)}{N - K_j - 1}$$

294

295 In addition to the classification results on data under different dimension
296 reduction degrees and the use of a spherical kernel, in this section we ex-
297 plore the performance of KDE-based classifiers with different kernel degrees.
298 First, we use a simple yet illustrative synthetic experiment to measure the
299 dependency of the classification performance with respect to the bandwidth.
300 Secondly, we provide a set of results on real data, using the spherical and full
301 models provided by ML-LOO, as well as the hybrid one, as described above.

302 3.2. Experiment on Synthetic Data

303 We illustrate the importance of an appropriate bandwidth choice by a
304 simple synthetic binary classification experiment. Data have been generated
305 following a 5 dimensional normal distribution. The labels have been gener-
306 ated according to the classification function $y = \text{sign}(x_4x_5)$. By means of
307 1000 training data, the spherical bandwidth has been obtained from these
308 data according to the rule (3). The $\hat{\sigma}$ obtained is shared by both classes. In
309 Fig. 5, the classification performance is shown on a test dataset with 500
310 samples in an interval of bandwidth values ranging from $\hat{\sigma}^2/10$ y $10\hat{\sigma}^2$. Note

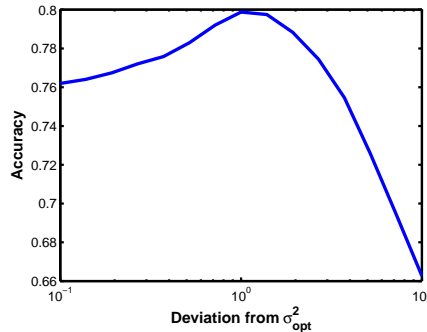


Figure 5: Classification performance averaged over 100 trials as a function of the bandwidth used.

311 that the best performance is reached when the bandwidth used is equal to
 312 $\hat{\sigma}$. The performance is shown to be very sensitive to the bandwidth choice.

313 *3.3. Experiments on Real Data*

Dataset	N_c	D	Size
Landsat	6	36	6435
Optdigits	10	64 (40)	5620
Letter	26	16	20000
Isolet	26	617 (40)	7797
Waveform	3	21	800
Segmentation	7	19 (8)	2310

Table 1: Characteristics of the public Datasets. The numbers between brackets express the dimension after dimension reduction by principal component analysis.

314 We have tested several Parzen classifiers, with spherical, unconstrained
 315 and hybrid kernels, on the same public datasets described in Table 1. We

316 have made sure that the convergence conditions hold for all the datasets,
 317 including the more complex and computationally costly condition (6). We
 318 show the results for both the whitened and the raw data in the cases in which
 319 this is possible: those in which the covariance matrix of each per-class dataset
 320 is non-singular. Otherwise, only the result on whitened data is reported.

321 We have taken ten independent samples for each experiment, where 75%
 322 of the data are used for training, and the rest are used for test.

Dataset	Generative				Discriminative	Reference
	Spherical KDE	Full KDE	Hybrid KDE	Scott	KNN	
Landsat (raw)	90.11± 0.74	86.34± 0.90 $\overset{ML}{PL}$	84.87± 0.73	84.34± 0.75	90.61± 0.56	97.58± 0.84
Landsat (whitened)	67.82± 0.44	85.85± 0.42 $\overset{ML}{PL}$	84.41± 0.43	84.00± 0.48	66.44± 0.84	94.77± 1.05
Optdigits (whitened)	97.95± 0.33	98.26 ± 0.26 $\overset{ML}{PL}$	98.80± 0.23	97.83± 0.23	97.39± 0.39	99.98± 0.05
Letter (raw)	95.58± 0.28	93.50 ± 0.60 $\overset{ML}{PL}$	95.34± 0.39	95.35± 0.29	95.65± 0.24	99.92 ± 0.11
Letter (whitened)	94.71± 0.25	94.30± 0.38 $\overset{ML}{PL}$	95.29± 0.34	95.25± 0.28	94.80 ± 0.25	99.93± 0.11
Isolet (whitened)	86.52 ± 0.50	90.86± 0.72 $\overset{ML}{PL}$	91.90± 0.67	91.69± 0.69	85.65± 0.55	99.74± 0.22
Waveform (raw)	78.50± 2.95 PL	74.45± 3.45 $\overset{ML}{PL}$	78.50± 3.20	77.95 ± 3.72	75.55± 2.77	86.77± 1.33
Waveform (whitened)	62.45± 3.07 PL	74.20± 3.54 $\overset{ML}{PL}$	78.90± 1.95	78.40± 1.97	57.25± 4.03	89.65± 3.89
Segmentation (whitened)	87.69± 1.19	93.29± 1.1613 $\overset{ML}{PL}$	92.46 ± 0.93	90.85± 1.05	95.67 ± 0.54	98.71 ± 0.58

Table 2: Classification performance on public datasets. The best models according to log-likelihood are marked with the "ML" label; those chosen by BIC or AIC are marked with "PL".

323 We have compared these results to the ones obtained by the generalization
 324 of Scott’s rule proposed in [11], which establishes the following covariance
 325 matrix for the kernel: $\mathbf{C} = N^{\frac{-2}{D+4}} \Sigma_x$, where Σ_x is the empirical covariance
 326 matrix of \mathbf{x} . We have also compared the generative KDE-based methods
 327 with the discriminative K-nearest neighbors (KNN, with $K = 1$). We also
 328 provide a reference value that can be considered as the highest achievable
 329 performance on each dataset. These reference values have been obtained by
 330 a state-of-the-art support vector machine (SVM) with radial basis function

331 as a kernel [19], where the hyperparameters C and σ have been chosen by a
332 grid-search and a 5-fold cross-validation on the training set.

333 The classification results are shown in Table 2. Boldfaced numbers mark
334 the highest performance on each dataset. Also, the model providing the
335 highest likelihood value is marked with "ML" label, and the one chosen by
336 the BIC or AIC^c criteria (which have agreed in all cases) is marked with a
337 "PL" (from *penalized likelihood*) label. According to the results, when using
338 generative modeling there is a version of ML-LOO that performs better than
339 Scott's rule for all cases, although the difference is not statistically significant
340 in some cases. Furthermore, better classification performance is obtained
341 by generative methods when compared to a discriminative method as KNN
342 in most cases. Note that, when performance is similar, generative models
343 have the advantage that they provide values of the joint distribution, which
344 allows both sampling from the distribution and obtaining likelihood values,
345 as opposed to discriminative methods like KNN or SVM.

346 Regarding the comparison between the bandwidth of the spherical KDE-
347 based classifier and the RBF width of the SVM, the latter is larger than the
348 former in all the cases evaluated (details are not given for lack of space). We
349 can find the reason in the inverse proportionality between number of sam-
350 ples and optimal bandwidth observed in both KDEs and SVMs. Therefore,
351 because the SVM provides a sparse solution (only a subset of the data is
352 involved), its optimal bandwidth becomes larger.

353 The results are not concluding regarding whether the spherical or the
354 full model provides better results. However, the fact that the hybrid ap-
355 proach provides the best result in the majority of the experiments reveals

356 that the method incorporates advantages from both the spherical and the
357 unconstrained methods, as described in Section 3.1. Unfortunately, neither
358 BIC nor AIC seem to give a conclusive clue about the kind of method that
359 provides the best performance - they agree only in two of the nine cases. From
360 this, we conclude that the ML-LOO procedure, in its different versions, pro-
361 vides competitive classification performance although neither likelihood, BIC
362 nor AIC provide an useful criterion to choose a priori among the different
363 options.

364 Surprisingly, in some cases the performance of spherical Parzen is worse
365 on the whitened data than in the original datasets, while the performance
366 with full, hybrid and Scott bandwidth is almost unchanged. This suggests
367 that i) whitening the data only helps in cases in which is unavoidable, i.e.
368 when the covariance matrix of the data is singular and a dimension reduction
369 is needed; and ii) whitening the whole dataset does not imply that the data
370 belonging to each class are whitened. The bad results of the spherical KDE
371 on the whitened versions of Landsat and Waveform are two examples of that.

372 4. Conclusions

373 We have provided a methodology for bandwidth selection in kernel density
374 estimators. Although the spherical Gaussian is the most widespread kernel
375 in these models, we have also described the procedure for a full Gaussian
376 kernel. The application of a maximum likelihood leave-one-out criterion has
377 led to a set of fixed-point algorithms that prevent us from carrying out an
378 exhaustive search of the optimal bandwidth parameter. The conditions for
379 the convergence of the algorithms proposed have also been established. In

380 particular, for the spherical approach the range of bandwidth values to which
381 the optimum belongs has been provided. The convergence in that range is
382 guaranteed.

383 Generative classification methods are usually in disadvantage compared
384 to discriminative methods for which the classification performance is the ob-
385 jective to maximize. However, we have explored the performance of a Parzen
386 classifier and obtained a performance that is less sensitive to the *curse of*
387 *the dimensionality* than traditionally attributed to KDE-based classifiers.
388 Moreover, higher performance than discriminative KNN is obtained in most
389 of the cases explored. Although both spherical and full versions of ML-
390 LOO converge in few iterations, we have also observed that the alternative
391 and computationally cheap bandwidth criterion based on Scott's rule pro-
392 vides reasonable results which can justify its use in large datasets. Finally,
393 although some clues have been given about the conditions in which the dif-
394 ferent versions ML-LOO perform better in classification, [a criterion based](#)
395 [on likelihood penalized by the number of parameters does not give success-](#)
396 [ful clues about the best model for classification.](#) This is explained by the
397 [fact that penalized likelihood evaluates a *full picture*, while the classification](#)
398 [accuracy is mainly determined by the accuracy of the models only close to](#)
399 [the classification boundaries.](#) The hybrid method seems to incorporate pos-
400 [itive elements from both the spherical and the unconstrained approaches,](#)
401 [and therefore it is the methods that have obtained the best result in a larger](#)
402 [number of experiments.](#)

403 **References**

- 404 [1] Ahmad, I., Lin, P., 1976. A nonparametric estimation of the entropy for
405 absolutely continuous distributions. *IEEE Transactions on Information*
406 *Theory* 22 (3), 372–375.
- 407 [2] Bilmes, J. A., 1997. A gentle tutorial of the EM algorithm and its appli-
408 cation to parameter estimation for gaussian mixture and hidden markov
409 models. Tech. Rep. tr-97-021, International Computer Science Institute,
410 Berkeley, California, USA.
- 411 [3] de Lima, M. S., Atuncar, G. S., 2011. A bayesian method estimate opti-
412 mal bandwidth multivariate kernel estimator. *Journal of Nonparametric*
413 *Statistics* 23 (1), 137–148.
- 414 [4] Devroye, L., Lugosi, G., Dec. 1996. A universally acceptable smoothing
415 factor for kernel density estimates (24), 2499–2512.
- 416 [5] Duin, R., 1976. On the choice of smoothing parameters for Parzen es-
417 timators of probability density functions. *IEEE Transactions on Com-*
418 *puters* 25 (11).
- 419 [6] Duong, T., Hazelton, M. L., 2005. Cross-validation bandwidth matri-
420 ces for multivariate kernel density estimation. *Scandinavian Journal of*
421 *Statistics* 32, 485–506.
- 422 [7] Filippone, M., Sanguinetti, G., 2011. Approximate inference of the band-
423 width in multivariate kernel density estimation. *Computational Statis-*
424 *tics & Data Analysis* 55 (12), 3104–3122.

- 425 [8] Fletcher, R., 1995. Practical Methods of Optimization (2nd Edition).
426 John Wiley & Sons, New York.
- 427 [9] Fukunaga, K., 1990. Introduction to Statistical Pattern Recognition.
428 Academic Press, New York.
- 429 [10] Goldberger, J., Roweis, S., Hinton, G., Salakhutdinov, R., 2005. Neigh-
430 bourhood components analysis. In: Advances in Neural Information
431 Processing Systems, NIPS 17. Vancouver, Canada, pp. 513–520.
- 432 [11] Haerdle, W., Mueller, M., Sperlich, S., Werwatz, A., 2004. Nonparamet-
433 ric and Semiparametric Models. Springer, New York.
- 434 [12] Hall, P., 1982. Cross-validation in density estimation. *Biometrika* 69 (2),
435 383–390.
- 436 [13] Leiva-Murillo, J., Artés-Rodríguez, A., 2008. Algorithms for gaussian
437 bandwidth selection in kernel density estimators. In: Advances in Neural
438 Information Processing System, NIPS, Optimization Workshop. Wistler,
439 Canada.
- 440 [14] McLachlan, G., 1997. The EM algorithm and extensions. John Wiley &
441 Sons, New York.
- 442 [15] Peltonen, J., Kaski, S., 2005. Discriminative components of data. *IEEE*
443 *Transactions on Neural Networks* 16 (1), 68–83.
- 444 [16] Pérez, A., Larrañaga, P., Inza, I., 2009. Bayesian classifiers based on
445 kernel density estimation: Flexible classifiers. *International Journal of*
446 *Approximate Reasoning* (50), 341–362.

- 447 [17] Platt, J. C., Platt, J. C., 1999. Probabilistic outputs for support vector
448 machines and comparisons to regularized likelihood methods. In: Ad-
449 vances in Large Margin Classifiers. MIT Press, pp. 61–74.
- 450 [18] Principe, J., Xu, D., Fischer, J., 2000. Information-Theoretic Learning.
451 Vol. 1 of Unsupervised Adaptive Filtering. John Wiley & Sons, New
452 York.
- 453 [19] Scholkopf, B., Smola, A., 2002. Learning with Kernels. The MIT Press,
454 Cambridge, MA.
- 455 [20] Scott, D., 1992. Multivariate Density Estimation. John Wiley & Sons,
456 New York.
- 457 [21] Schwarz, G., 1978. Estimating the dimension of a model. *Ann. Stat.* 14,
458 461–464.
- 459 [22] Silverman, B. W., 1986. Density Estimation for Statistics and Data
460 Analysis. Chapman & Hall, Londres.
- 461 [23] Steele, R. J., Raftery, A. E. Performance of Bayesian Model Selection
462 Criteria for Gaussian Mixture Models. Technical Report no. 559. De-
463 partment of Statistics, University of Washington.
- 464 [24] Torkkola, K., 2003. Feature extraction by non-parametric mutual infor-
465 mation maximization. *Journal on Machine Learning Research* 3, 1415–
466 1438.
- 467 [25] Turlach, B. A., 1993. Bandwidth selection in kernel density estimation:
468 A review. In: CORE and Institut de Statistique. pp. 1–33.

- 469 [26] Wand, M. P., Jones, M. C., 1995. Kernel Smoothing. Chapman and Hall,
470 London.
- 471 [27] Zhang, X., King, M. L., Hyndman, R. J., 2006. A bayesian approach to
472 bandwidth selection for multivariate kernel density estimation. *Compu-
473 tational Statistics & Data Analysis* 50 (11), 3009–3031.