

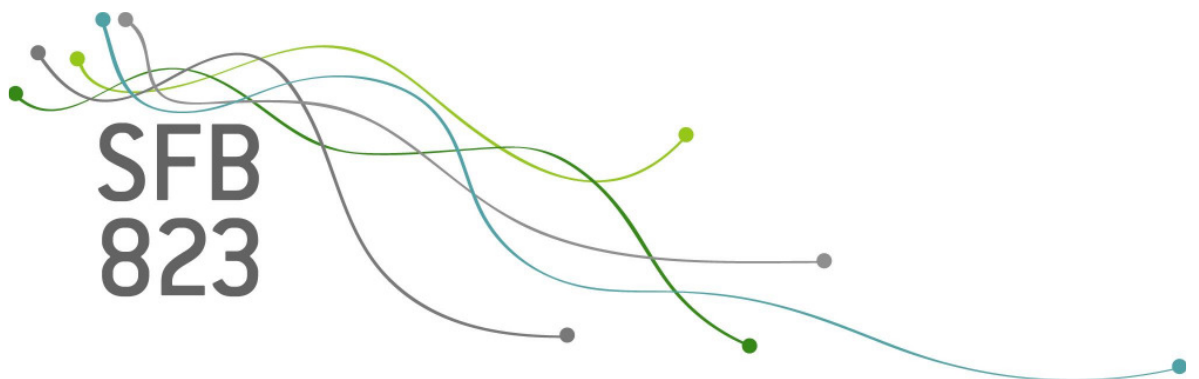
SFB
823

Combining regular and irregular histograms by penalized likelihood

Yves Rozenholc, Thoralf Mildenberger,
Ursula Gather

Nr. 31/2009

Discussion Paper



Combining Regular and Irregular Histograms by Penalized Likelihood

Yves Rozenholc
UFR de Mathématiques et d'Informatique
Université Paris Descartes

Thoralf Mildenerger*
Fakultät Statistik
Technische Universität Dortmund

Ursula Gather
Fakultät Statistik
Technische Universität Dortmund

November 23, 2009

A fully automatic procedure for the construction of histograms is proposed. It consists of constructing both a regular and an irregular histogram and then choosing between the two. For the regular histogram, only the number of bins has to be chosen. Irregular histograms can be constructed using a dynamic programming algorithm if the number of bins is known. To choose the number of bins, two different penalties motivated by recent work in model selection are proposed. A complete description of the algorithm and a proper tuning of the penalties is given. Finally, different versions of the procedure are compared to other existing proposals for a wide range of densities and sample sizes. In the simulations, the squared Hellinger risk of the procedure that chooses between regular and irregular histograms is always at most twice as large as the risk of the best of the other methods. The procedure is implemented in an R-Package.

1. Introduction

For a sample (X_1, X_2, \dots, X_n) of a real random variable X with an unknown density f w.r.t. Lebesgue measure, we denote the realizations by (x_1, x_2, \dots, x_n) and the realizations of the order statistics by $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$. The goal in nonparametric density

*Address for correspondence: Fakultät Statistik, Technische Universität Dortmund, 44221 Dortmund, Germany. E-mail: mildenbe@statistik.tu-dortmund.de. Web: http://www.statistik.tu-dortmund.de/mildenerger_en.html

estimation is to construct an estimate \hat{f} of f from the sample. In this work, we focus on estimation by histograms, which are defined as piecewise constant densities. The procedure we propose consists of constructing both a regular and an irregular histogram (both to be defined below) and then choosing between the two. Although other types of nonparametric density estimators are known to be superior to histograms according to several optimality criteria, histograms still play an important role in practice. The main reason is their simplicity and hence their interpretability (Birgé and Rozenholc, 2006). Often, the histogram is the only density estimator taught to future researchers in non-mathematical subject areas, usually introduced in an exploratory context without reference to optimality criteria.

We first introduce histograms and describe the connection to Maximum Likelihood estimation: Given (x_1, x_2, \dots, x_n) and a set of densities \mathcal{F} , the maximum likelihood estimate – if it exists – is given by an element $\hat{f} \in \mathcal{F}$ that maximizes the likelihood $\prod_{i=1}^n f(x_i)$ or equivalently its logarithm, the log-likelihood $L(f, x_1, \dots, x_n)$:

$$\hat{f} := \operatorname{argmax}_{f \in \mathcal{F}} L(f, x_1, \dots, x_n) := \operatorname{argmax}_{f \in \mathcal{F}} \sum_{i=1}^n \log(f(x_i)).$$

Without further restrictions on the class \mathcal{F} , the log-likelihood is unbounded, and hence, no maximum likelihood estimate exists. One possibility is to restrict \mathcal{F} to a set of histograms. Consider a partition $\mathcal{I} := \{I_1, \dots, I_D\}$ of a compact interval $K \subset \mathbb{R}$ into D intervals I_1, \dots, I_D , such that $I_i \cap I_j = \emptyset$ for $i \neq j$ and $\bigcup I_i = K$. Now consider the set $\mathcal{F}_{\mathcal{I}}$ of all histograms that are piecewise constant on \mathcal{I} and zero outside \mathcal{I} :

$$\mathcal{F}_{\mathcal{I}} := \left\{ f \left| f = \sum_{j=1}^D h_j \mathbf{1}_{I_j}, h_j \geq 0, j = 1, \dots, D \text{ and } \sum_{j=1}^D h_j |I_j| = 1 \right. \right\},$$

where $\mathbf{1}_A$ denotes the indicator function of a set A and $|I|$ the length of the interval I . If K contains $[x_{(1)}, x_{(n)}]$, the *Maximum Likelihood Histogram* (ML histogram) is defined as the maximizer of the log-likelihood in $\mathcal{F}_{\mathcal{I}}$ and is given by

$$\hat{f}_{\mathcal{I}} := \operatorname{argmax}_{f \in \mathcal{F}_{\mathcal{I}}} L(f, x_1, \dots, x_n) = \frac{1}{n} \sum_{j=1}^D \frac{N_j}{|I_j|} \mathbf{1}_{I_j}, \quad (1)$$

with $N_j = \sum_{i=1}^n \mathbf{1}_{I_j}(x_i)$. Its log-likelihood is

$$L(\hat{f}_{\mathcal{I}}, x_1, \dots, x_n) = \sum_{j=1}^D N_j \log \frac{N_j}{n|I_j|}. \quad (2)$$

In the following, we consider partitions $\mathcal{I} := \mathcal{I}_D := (I_1, \dots, I_D)$ of the interval $I := [x_{(1)}, x_{(n)}]$, consisting of D intervals of the form

$$I_j := \begin{cases} [t_0, t_1] & j = 1 \\ (t_{j-1}, t_j] & j = 2, \dots, D \end{cases},$$

with breakpoints $x_{(1)} =: t_0 < t_1 < \dots < t_D =: x_{(n)}$. A histogram is called *regular* if all intervals have the same length and *irregular* otherwise. The intervals are also referred to as *bins*.

We will only consider ML histograms in this work, and use the term "histogram" synonymously with "ML histogram" unless explicitly stated otherwise. We focus on finding a data-driven construction of a histogram with good risk behavior. Given a distance measure d between densities, the risk is defined as the expected distance between the true and the estimated density:

$$R_n(f, \hat{f}_{\mathcal{I}}, d) := E_f[d(f, \hat{f}_{\mathcal{I}}(X_1, \dots, X_n))].$$

We consider the risks with respect to the following loss functions:

- Squared Hellinger distance

$$d_H(f, g) = \frac{1}{2} \int (\sqrt{f(t)} - \sqrt{g(t)})^2 dt, \quad (3)$$

which has been normalized such that its maximum value is 1.

- Powers of the L_p -norms (for $p = 1$ and 2) defined by

$$d_p := \|f - g\|_p^p = \int |f(t) - g(t)|^p dt. \quad (4)$$

The L_2 distance is widely used mainly for its mathematical tractability. It is often possible to derive explicit expressions for the L_2 risk at least asymptotically, cf. Kogure (1987). However, as argued by Devroye and Györfi (1985, ch. 1), the L_1 distance can be considered more natural in the context of density estimation because – unlike the L_2 distance – it is defined for all densities and it has desirable invariance properties. We focus mainly on the Hellinger distance for several reasons: it is also defined for any two densities, it has important invariance properties and the results of Castellan (1999, 2000) are derived for the corresponding risk. Another widely used loss function is the Kullback-Leibler distance

$$d_{KL}(f, g) := \int \log \left(\frac{f(t)}{g(t)} \right) f(t) dt$$

which is not suitable in histogram density estimation since it is infinite whenever the estimated density is zero on an interval where the true distribution has positive mass. Hence it is excluded from consideration. For a detailed discussion on the choice of loss functions in histogram density estimation, see section 2.2. in Birgé and Rozenholc (2006) and the references given there.

Given the sample, the histogram $\hat{f}_{\mathcal{I}}$ depends only on the chosen partition $\mathcal{I} = (I_1, \dots, I_D)$. The values on the intervals of the partition are fixed, namely equal to the relative frequencies divided by the bin widths. In order to achieve good performance in terms of risk, the crucial point is thus choosing the partition. A naïve comparison of the likelihood of histograms for partitions with different numbers of bins is misleading since partitions with too many bins will result in a large likelihood without yielding a sensible estimate of f . But also without any further restrictions on the allowed partitions the likelihood can be made arbitrarily large for a fixed number of bins.

Many approaches exist for the special case of regular histograms where I is divided into D equal sized bins; the problem is then reduced to the choice of D , cf. Birgé and Rozenholc (2006) and Davies et al. (2008) and the references given there. Using irregular partitions can reduce bias and therefore can improve performance for spatially inhomogenous densities, but

the increased difficulty of choosing a good partition may lead to an increase in risk for more well-behaved densities. The idea of constructing both a regular and an irregular histogram and then choosing between the two is briefly discussed in Birgé and Rozenholc (2006). To our knowledge, this approach has not yet been put into practice.

Our recommendation is to construct a regular histogram that maximizes the penalized log-likelihood

$$L(\hat{f}_{\mathcal{I}}, x_1, \dots, x_n) - (D - 1) - \log^{2.5} D \quad (5)$$

among all regular partitions of $[x_{(1)}, x_{(n)}]$ with $D = 1, \dots, \lfloor n/\log n \rfloor$ bins (where $\lfloor x \rfloor$ denotes the largest integer not larger than x) and an irregular histogram that maximizes the penalized log-likelihood

$$L(\hat{f}_{\mathcal{I}}, x_1, \dots, x_n) - \log \binom{n-1}{D-1} - (D - 1) - \log^{2.5} D \quad (6)$$

among a set of partitions of $[x_{(1)}, x_{(n)}]$ with breakpoints equal to the sample points, where again D is the number of bins in a given partition. The final estimate is then the one with the larger penalized log-likelihood. The penalty in (5) for the regular case was proposed in Birgé and Rozenholc (2006), while the motivation for (6) is developed later in this paper, where we consider different penalty forms for the irregular case.

Several methods have been developed previously to choose a good irregular histogram. Kogure (1987) gives asymptotic results for the optimal choice of bins. His approach is based on using blocks of equisized bins, and the dependence on tuning parameters is explored via simulations in his PhD thesis (Kogure, 1986). It does not result in a fully automatic procedure. Kanazawa (1988) proposes to control the Hellinger distance between the unknown true density and the estimated histogram and introduces a dynamic programming algorithm to find the best partition with a given number of bins. Kanazawa (1992) derives the asymptotically optimal choice of the number of bins. Unfortunately this result involves the first and second derivatives of the unknown density, which leads to a construction that cannot be applied from a practical point of view. Celisse and Robin (2008) give explicit formulas for L_2 leave- p -out cross-validation for regular and irregular histograms. They only briefly comment on the case of irregular histograms and only show simulations with ad-hoc choices of the set of partitions. In our simulations, we use their explicit formula to compare risk behavior of cross-validation and our penalized likelihood approach when both are used to choose an irregular histogram from the same set of partitions. The multiresolution histogram by Engel (1997) is based on a tree of dyadic partitions to control the L_2 -error. The performance crucially depends on the finest resolution level, for which no universally usable recommendation is given. Some other tree-based procedures have been suggested for the multivariate case. They can be used for the univariate case, but they either perform a complete search over a restricted set of partitions (Blanchard et al., 2007; Klemelä, 2009) or a greedy search over a full set of partitions (Klemelä, 2007) to deal with computational problems that do not occur in the univariate case. Theoretical results on conditions for consistency of histogram estimates with data-driven and possibly irregular partitions are derived in Chen and Zhao (1987); Zhao et al. (1988) and Lugosi and Nobel (1996). Devroye and Lugosi (2004) give a construction of histograms where bin widths are allowed to vary according to a pre-specified function.

Hartigan (1996) considers regular and irregular histogram construction from a Bayesian point of view. However, we are not aware of any fully tuned automatic Bayesian procedure for irregular histogram construction. Rissanen et al. (1992) give a construction based on the Minimum Description Length (MDL) paradigm, which leads to a penalized likelihood estima-

tor. A choice of several discretization parameters is needed, and the recommendation given by the authors is to perform an exhaustive search over all possible combinations of values, which is computationally expensive. A more recent proposal by Kontkanen and Myllymäki (2007) is also based on the MDL principle; it involves a discretization which results in the estimate not being a proper density. Catoni (2002) suggests a multi-stage procedure based on coding ideas that computes a density estimate by aggregating histograms.

The taut string procedure introduced by Davies and Kovac (2004) can also be used to generate an irregular histogram as described in Davies et al. (2008). Regularization is achieved not by controlling the number of bins but by controlling the modality of the estimate. The stated aim of the authors is not to minimize some risk but to find an estimate of the density that has minimum number of modes that could have generated the data, where the latter is formalized by a criterion based on differences of Kuiper metrics between the empirical and the estimated distribution. The main idea is to construct a piecewise linear spline of minimal length (the taut string) in a tube around the empirical cdf and then take its derivative, which is piecewise constant. The histogram is then constructed using the knots of the string as the boundaries of the bins. This coincides with the derivative of the string except on intervals where the string switches from the upper to the lower boundary of the tube or vice versa. Let us emphasize that although the partition is chosen without reference to maximum likelihood, the histogram constructed in this way fulfils definition (1). The main tuning parameter is the tube width, and an automatic choice is suggested by the authors. The procedure has shown a particularly good behavior also w.r.t. classical loss functions (Davies et al., 2008), and therefore is compared with our method in our simulations.

For our construction of irregular histograms, we will focus on penalized likelihood maximization techniques. For a good data-driven histogram estimate one needs an appropriate penalization to provide an automatic choice of D as well as of the partition $\mathcal{I} = (I_1, \dots, I_D)$. Since Akaike's Information Criterion (AIC) introduced by Akaike (1973), penalized likelihood has been used with many different penalty terms. AIC aims at ensuring a good risk behavior of the resulting estimate. Another widely used criterion is the Bayesian Information Criterion (BIC) introduced by Schwarz (1978). It is constructed in such a way as to consistently estimate the smallest true model order, which in histogram density estimation would lead to very large models unless the true density is piecewise constant. In practice, criteria like AIC and BIC are routinely applied in many different statistical models, often without reference to their different conceptual backgrounds and without appropriate modifications for the model under consideration. In their original forms, both AIC and BIC do not account for multiple partitions with the same number of bins. See Chapter 7.3 of Massart (2007) for a critique of the use of AIC in histogram density estimation. Since both are widely used, we include them in our comparisons. Our penalties are motivated by recent model selection works by Barron et al. (1999), Castellán (1999, 2000) and Massart (2007). The regular histogram construction proposed in Birgé and Rozenholc (2006) with which we combine our irregular histogram is based on the same ideas.

Our paper is structured as follows: In Section 2, we review the problem of constructing an irregular histogram using penalized likelihood. Section 3 gives a description of the choice of the penalty. Section 4 gives a detailed description of the proposed procedure for irregular histograms. In Section 5, we comment on the empirical evaluation of the risks under consideration. Section 6 gives the results of a simulation study and conclusions.

2. Penalized likelihood construction of irregular histograms

Constructing an irregular histogram by penalized likelihood means maximizing

$$L(\hat{f}_{\mathcal{I}}, x_1, \dots, x_n) - \text{pen}_n(\mathcal{I}) \quad (7)$$

w.r.t. partitions $\mathcal{I} = (I_1, \dots, I_{|\mathcal{I}|})$ of $[x_{(1)}, x_{(n)}]$, where $\text{pen}_n(\mathcal{I})$ is a penalty term depending only on the partition \mathcal{I} and possibly on the sample (data-driven). We will introduce a new choice here motivated by work of Barron et al. (1999), Castellan (1999, 2000) and Massart (2007).

Optimizing w.r.t. the partition \mathcal{I} with $|\mathcal{I}|$ fixed in (7) leaves us with a continuous optimization problem. Without further restrictions, for $|\mathcal{I}| \geq 2$ the likelihood is unbounded. The partition

$$\{[x_{(1)}, x_{(1)} + \eta), [x_{(1)} + \eta, x_{(n)}]\}$$

for small η (such that $N_1 = 1$ and $N_2 = n - 1$) leads to a log-likelihood equal to

$$N_1 \log \frac{N_1}{n|I_1|} + N_2 \log \frac{N_2}{n|I_2|} = \log \frac{1}{n\eta} + (n - 1) \log \frac{n - 1}{n(x_{(n)} - x_{(1)} - \eta)}$$

which can be arbitrarily large when η goes to 0.

One possibility is to restrict to all partitions that are built with endpoints on the observations; the optimization problem (7) can then be solved using a dynamic programming algorithm first used for histogram construction by Kanazawa (1988). More details are given in Section 4.

With $D = |\mathcal{I}|$, we propose the following families of penalties parametrized by two constants c and α :

$$\text{pen}_n^A(\mathcal{I}) = c \log \binom{n-1}{D-1} + \alpha(D-1) + \varepsilon_{c,\alpha}^{(1)}(D), \quad (8)$$

$$\text{pen}_n^B(\mathcal{I}) = c \log \binom{n-1}{D-1} + \alpha(D-1) + \varepsilon^{(2)}(D) \quad (9)$$

$$(10)$$

and

$$\text{pen}_n^R(\mathcal{I}) = c \log \binom{n-1}{D-1} + \frac{\alpha}{n} \sum_{j=1}^D \frac{N_j}{|I_j|} + \varepsilon^{(2)}(D). \quad (11)$$

where

$$\varepsilon_{c,\alpha}^{(1)}(D) = ck \log D + 2\sqrt{c\alpha(D-1)(\log \binom{n-1}{D-1} + k \log D)} \quad (12)$$

and

$$\varepsilon^{(2)}(D) = \log^{2.5} D. \quad (13)$$

The precise choices for c and α obtained by simulations are described in Section 3. Note that, while the penalties given in (8) and (9) depend only on the number of bins of the partition, the penalty in formula (11) is a *random penalty* in the sense that it also depends on the data.

We now give arguments to explain the origins of these penalties. The penalty defined by (8) is derived from Theorem 3.2 in Castellan (1999), which is also stated as Theorem 7.9 in Massart (2007, p. 232) and from eq. (7.32) in Theorem 7.7 in Massart (2007, p.219). From the penalty form in Theorem 7.9 in Massart (2007) we derive $\varepsilon^{(1)}$:

$$\text{pen}_n(\mathcal{I}) = c_1(\sqrt{D-1} + \sqrt{c_2 x_{\mathcal{I}}})^2, \quad (14)$$

where the weights $x_{\mathcal{I}}$ are chosen such that

$$\sum_D \sum_{|\mathcal{I}|=D} e^{-x_{\mathcal{I}}} \leq \Sigma \quad (15)$$

for an absolute constant Σ . Because the endpoints of our partitions are fixed, there are $\binom{n-1}{D-1}$ different partitions with cardinality D , and we assign equal weights x_D to every partition \mathcal{I} with $|\mathcal{I}| = D$ such that

$$\sum_D \binom{n-1}{D-1} e^{-x_D} \leq \Sigma.$$

To achieve this, we set

$$x_D = \log \binom{n-1}{D-1} + \varepsilon(D).$$

Then (15) becomes

$$\sum_D e^{-\varepsilon(D)} \leq \Sigma.$$

Choosing $\varepsilon(D)$ of the form $k \log D$ with $k > 1$ ensures that the sum is converging and that Σ is finite. Finally for $k > 1$ we have

$$x_D = \log \binom{n-1}{D-1} + k \log D.$$

Substitution into (14) gives

$$\begin{aligned} \text{pen}_n(\mathcal{I}) &= c_1 \left(D-1 + c_2 \left(\log \binom{n-1}{D-1} + k \log D \right) \right) \\ &\quad + 2 \sqrt{c_2(D-1) \left(\log \binom{n-1}{D-1} + k \log D \right)}. \end{aligned} \quad (16)$$

Let us emphasize that Theorem 7.9 in Massart (2007, p. 232) requires $c_1 > 1/2$ and $c_2 = 2(1 + 1/c_1)$. Coming back to our notations, with $\alpha = c_1$, $c = c_1 c_2$ we obtain Equation (12).

We now use Theorem 7.7 in Massart (2007, p. 219) to justify the random penalty in (11). The orthonormal basis considered in this theorem for a given partition \mathcal{I} consists of all $\mathbf{1}_I/\sqrt{|I|}$ for all I in \mathcal{I} . The least squares contrast used in this theorem in our framework is $-n^{-2} \sum_{I \in \mathcal{I}} N_I^2/|I|$. To link the minimization of the least squares contrast and the maximization of the log-likelihood, we consider the following approximation:

$$L(\hat{f}_{\mathcal{I}}, x_1, \dots, x_n) = \sum_{j=1}^D N_j \log \left(\frac{N_j}{n|I_j|} \right) \approx \sum_{j=1}^D N_j \left(\frac{N_j}{n|I_j|} - 1 \right) = \frac{1}{n} \sum_{j=1}^D \frac{N_j^2}{|I_j|} - n.$$

From the penalty form (7.32) and the use of $M = 1$ and $\varepsilon = 0$ in Theorem 7.7 in Massart (2007, p. 219), following the same derivation as for $\varepsilon^{(1)}$, we find the penalty in (16) with $c_1 = 1$ and $c_2 = 2$.

Using the least squares approximation, we can use the random penalty (7.33) in Theorem 7.7 in Massart (2007). Let us emphasize that \widehat{V}_m defined by Massart is in our framework $\sum_{I \in \mathcal{I}} N_I/n|I|$ with $m = \mathcal{I}$. To derive $\varepsilon^{(2)}$ in (11) we start from the penalty defined in (7.33) in Massart (2007):

$$\text{pen}_n(\mathcal{I}) = (1 + \varepsilon)^5 \left(\sqrt{\widehat{V}_{\mathcal{I}}} + \sqrt{2ML_{\mathcal{I}}D} \right)^2.$$

Following the same derivations as for the penalty (14), setting $M = 1$, $\varepsilon = 0$ and $L_{\mathcal{I}} = D^{-1}(\log \binom{n-1}{D-1} + k \log D)$ we obtain:

$$\begin{aligned} \text{pen}_n(\mathcal{I}) &= \widehat{V}_{\mathcal{I}} + 2 \log \binom{n-1}{D-1} + 2k \log D \\ &\quad + 2 \sqrt{2\widehat{V}_{\mathcal{I}} \left(\log \binom{n-1}{D-1} + k \log D \right)}. \end{aligned}$$

Let us emphasize that, because of terms of the form $\varphi(D)\widehat{V}_{\mathcal{I}}$, the expression in the square root above prevents the use of dynamic programming to compute the maximum of the penalized log-likelihood defined in (7). To avoid this problem we propose, following penalty forms proposed in Birgé and Rozenholc (2006) and Comte and Rozenholc (2004), to replace the remainder expression

$$2k \log D + 2 \sqrt{2\widehat{V}_{\mathcal{I}} \left(\log \binom{n-1}{D-1} + k \log D \right)}$$

by a power of $\log D$. We have tried several values of the power and found that Formula (13) leads to a good choice. Finally, we also replaced $\varepsilon_{c,\alpha}^{(1)}$ in formula (8) by $\varepsilon^{(2)}$, leading to the penalty given in (9).

3. Choice of the Penalty

Using histograms with the endpoints of the partition placed on the observations as described later in Section 4, we ran empirical risk estimation in order to fix our penalty using the losses defined by (3) and (4) for $p = 1$ and 2. We used the same densities for calibration as in the simulations described in Section 6 but different samples and a smaller number of replications. We focused on the Hellinger risk to obtain good choices of the penalties, but the behavior w.r.t. L_1 loss is very similar. For minimizing L_2 risk, other choices may be preferable. Since no single penalty is best in all cases, the calibration of a penalty always leads to some compromise. We describe in the following what we consider to be a good proposal.

In formula (8) we tried:

- $c = 2(\alpha + 1)$ and $\alpha \in \{0.5, 1\}$ following Theorem 7.9 in Massart (2007).
- $c = 2$ and $\alpha = 1$ following Theorem 7.7 eq. (7.32) in Massart (2007) with $M = 1$ and $\varepsilon = 0$.

- $c = 1$ and $\alpha \in \{0.5, 1\}$.

We always set $k = 2$. From these experiments, the most satisfactory choice is $c = 1$ and $\alpha = 0.5$. We also ran experiments replacing $\varepsilon_{c,\alpha}^{(1)}$ by $\varepsilon^{(2)}$, leading to the penalty given in (9). In this case, we have found that the most satisfactory choice is $c = 1$ and $\alpha = 1$, and this choice is even better than $\varepsilon_{2,1}^{(1)}$. Note that the resulting penalty, given in (6), exactly corresponds to the penalty in (5) proposed in Birgé and Rozenholc (2006) for the regular case, except for the additional term $\log \binom{n-1}{D-1}$ that is needed to account for multiple partitions with the same number of bins. This term is zero for a histogram with just one bin, so the penalized likelihoods for regular and irregular histograms can directly be compared in this case. Because (8) and (9) are very similar, we only use this version in our simulations in Section 6.

For the random penalty in formula (11) we ran risk evaluation experiments using all combinations of $c \in \{0.5, 1, 2\}$ and $\alpha \in \{0.5, 1\}$. Let us emphasize that $c = 2$ and $\alpha = 1$ correspond to formula (7.33) in Massart (2007) up to our choice of $\varepsilon^{(2)}$ defined in (13). From our point of view, the most satisfactory choice is $c = 1$ and $\alpha = 0.5$. When comparing the log-likelihood penalized in this way to the one in (5), care has to be taken to ensure that both give the same value for a histogram with just one bin.

To conclude this section, we remark that the results are very close. Only for the trimodal uniform density have we found differences in favor of the penalty (9). For all other densities, the absolute values of the relative differences $\left| \frac{\widehat{R}_n^R - \widehat{R}_n^D}{\widehat{R}_n^D} \right|$ of the risks are less than 0.163.

4. Algorithm for constructing irregular histograms

We maximize (7) w.r.t. partitions \mathcal{I} built with endpoints on the observations:

$$\mathcal{I} = ([x_{(1)}, x_{(k_1)}], [x_{(k_1)}, x_{(k_2)}], [x_{(k_2)}, x_{(k_3)}], \dots, [x_{(k_{D-2})}, x_{(k_{D-1})}], [x_{(k_{D-1})}, x_{(n)}]),$$

where $1 < k_1 < \dots < k_{D-1} < n$. We start from a "finest" partition \mathcal{I}_{\max} defined by $D_{\max} < n$ and the choice $1 < k_1 < \dots < k_{D_{\max}-1} < n$. Let us write this partition as

$$\mathcal{I}_{\max} = (I_1^0, \dots, I_{D_{\max}}^0),$$

where $I_d^0 = (t_{d-1}, t_d]$ for $d = 1$ to D_{\max} and where $t_0 = x_{(1)} - \text{eps}$, $t_{D_{\max}} = x_{(n)}$ and $t_d = x_{(k_d)}$ for $0 < d < D_{\max}$. Here eps represents the machine precision and is used only to allow for the use of left-open, right-closed intervals. Our aim is to build a sub-partition \mathcal{I} of \mathcal{I}_{\max} which maximizes (7). This problem is solved in polynomial time by a dynamic programming (DP) algorithm as used in Kanazawa (1988) and Comte and Rozenholc (2004). We briefly describe the algorithm in our context of penalized histograms. Let us assume that, given the sample, (7) can be rewritten as $\Phi^0(\mathcal{I}) + \Psi(D, n)$, where Φ^0 is an additive function with respect to the partition in the sense that

$$\Phi^0(\mathcal{I}) = \Phi(I_1) + \dots + \Phi(I_D) \text{ if } \mathcal{I} = (I_1, \dots, I_D).$$

In our case, $\Phi(I)$ depends only on the number N_I of observations in interval I and on its length $|I|$. More precisely, for a penalty of the form (8) or (9), we have

$$\Phi(I) = N_I \log \frac{N_I}{n|I|}, \tag{17}$$

and $\Psi(D, n) = \text{pen}_n^A(\mathcal{I})$ or $\Psi(D, n) = \text{pen}_n^B(\mathcal{I})$, respectively. For a penalty of the form (11) we have

$$\Phi(I) = N_I \log \frac{N_I}{n|I|} - \frac{\alpha N_I}{n|I|},$$

and $\Psi(D, n) = c \log \binom{n-1}{D-1} + \varepsilon^{(2)}(D)$.

We denote $p_1(i, j) = \Phi((t_i, t_j])$ and $p_1(j) := p_1(0, j)$. Finally, let us define $i_1(j) = 0$. Assume that we have already computed all $p_1(i, j)$ for $0 \leq i < j \leq D_{\max}$ (which needs $O(D_{\max}^2)$ operations). The dynamic programming algorithm works as follows. First, the maxima of Φ_0 for partitions with $D = 1, \dots, D_{\max}$ bins are calculated:

- For $D = 2 \dots D_{\max}$
 - For $j = D \dots D_{\max}$,
 - $i_D(j) = \arg_i \max_{D-1 \leq i < j} [p_{D-1}(i) + p_1(i, j)]$;
 - $p_D(j) = p_{D-1}(i_D(j)) + p_1(i_D(j), j)$

For each D , $p_D(D_{\max})$ is the maximum of $\Phi^0(\mathcal{I})$ for all sub-partitions \mathcal{I} of our finest partition \mathcal{I}_{\max} with D bins. The partition which achieves the maximum of $\Phi^0(\mathcal{I}) + \Psi(D, n)$ can be constructed in the following way:

- Compute $\hat{D} = \arg_D \max_{1 \leq D \leq D_{\max}} p_D(D_{\max}) + \Psi(D, n)$.
- Initialize $L = D_{\max}$
- For $j = \hat{D}, \dots, 1$: $L = c(i_j(L[1]), L)$

The vector L defines the indices of the t_j 's which are the endpoints of the best partition in the sense of (7). The notation $L[1]$ denotes the first coordinate of the vector L and $c(u, L)$ denotes concatenation of the vector u with L .

The computation of $i_D(j) = \arg_i \max_{D-1 \leq i < j} [p_{D-1}(i) + p_1(i, j)]$ requires $O(j - D + 1)$ operations and the total complexity of this algorithm is of order D_{\max}^3 . Hence the total number of operations may be of order n if we start from a finest partition with D_{\max} of order $n^{1/3}$. We propose to use a greedy algorithm in order to build a finest partition with $\lfloor \max\{n^{1/3}, 100\} \rfloor$ bins if this number is smaller than n . Let us call $\mathcal{E}(\mathcal{I})$ the set of endpoints of partition \mathcal{I} . Starting with the partition $\mathcal{I}_0 = ([x_{(1)}, x_{(n)}])$, we produce a sequence of partitions \mathcal{I}_D satisfying :

$$\mathcal{I}_{D+1} = \arg \max \Phi^0(\mathcal{I}),$$

where the maximum is taken over all partitions \mathcal{I} with $\mathcal{E}(\mathcal{I}) = \mathcal{E}(\mathcal{I}_D) \cup \{t\}$ with t in $\{x_1, \dots, x_n\} \setminus \mathcal{E}(\mathcal{I}_D)$. For all three penalty forms, we use a greedy maximization of the likelihood to obtain this partition, i.e. we always use Φ as in (17).

Let us remark that the theoretical results by Castellán (1999, 2000) and Massart (2007, ch. 7), are derived for the case of a finest regular grid with bin sizes not smaller than a constant times $\log^2(n)/n$. In particular, the set of partitions is fixed beforehand and may depend on n but not on the sample. This also means that that no bins are possible that are shorter than a constant times $\log^2(n)/n$. However, we found that, in practice, we can improve performance drastically for some densities by using a data-dependent finest grid imposing no restrictions on the smallest bins without losing much at other densities. More comments on this are given in Section 6.

5. Risk evaluation

The risks of the procedures are evaluated empirically by means of simulations. For each density f and each sample size n , N samples $x^{(j)} := (x_1^{(j)}, \dots, x_n^{(j)})$, $j = 1, \dots, N$ are generated and the loss functions $d = d_H, d_1, d_2$ are evaluated for every histogram procedure \hat{f} . We estimate the risks $R_n(f, \hat{f}, d)$ by

$$\hat{R}_n(f, \hat{f}, d) := \frac{1}{N} \sum_{j=1}^N d(f, \hat{f}(x_1^{(j)}, \dots, x_n^{(j)})).$$

We now describe how we computed our loss functions (3) and (4) to obtain empirical risk evaluation. To estimate the risks, we evaluate the losses $d(f, \hat{f}(x_1^{(j)}, \dots, x_n^{(j)}))$ for every simulation run j by numerical integration. First note that the integrals appearing in (3) and (4) are all of the form

$$\int \delta(t) dt := \int \tilde{\delta}(f(t), g(t)) dt$$

for continuous functions $\tilde{\delta}$. Care has to be taken of discontinuities in both the true densities f and the histogram estimates \hat{f} and furthermore the bilogarithmic peak density has infinite peaks. For given f and \hat{f} , let $\tau_1 < \dots < \tau_{L-1}$ denote the points where f or \hat{f} is discontinuous or infinite. Defining the intervals $J_0 := (-\infty, \tau_1)$, $J_l := (\tau_l, \tau_{l+1})$, $l = 1, \dots, L-1$, $J_L := (\tau_L, \infty)$, we split up the integrals into sums of integrals over open intervals where δ is continuous:

$$\int_{\mathbb{R}} \delta(t) dt := \sum_{l=0}^L \int_{J_l} \delta(t) dt.$$

Note that we use open intervals to allow both f and \hat{f} to take any (possibly infinite) value in the point τ_1, \dots, τ_L . To evaluate the integrals on $J = J_1, \dots, J_{L-1}$ we use the trapeze rule

$$\int_{J_l} \delta(t) dt \approx (\kappa_K^l - \kappa_1^l) \left(\frac{1}{2} \delta(\kappa_1^l) + \delta(\kappa_2^l) + \dots + \delta(\kappa_{K-1}^l) + \frac{1}{2} \delta(\kappa_K^l) \right)$$

for equispaced grid points $\kappa_1^l = \tau_l + \varepsilon$, $\kappa_K^l = \tau_{l+1} - \varepsilon$ and $\kappa_\nu^l = \kappa_1^l + (\nu-1)h$ for $\nu = 2, \dots, K-1$ with $h = \frac{\tau_{l+1} - \tau_l - 2\varepsilon}{K-1}$. We set $\varepsilon = 10^{-11}$ to integrate over open intervals.

Note that on the unbounded intervals J_0 and J_L for $d = d_H$ and $d = d_1$ we have $\int_{J_l} \delta(t) dt = \int_{J_l} f(t) dt$ since \hat{f} is zero. For d_2 we replace $\pm\infty$ in the definition of J_0 and J_L by the upper and lower 10^{-10} -quantiles of f and integrate numerically as on the other intervals, in case the support of f is unbounded. Otherwise, the integrals over J_0 and J_L are zero.

6. Simulation Study and Conclusions

In order to choose the constants in the penalties given in Section 3 and to assess the performance of the penalized likelihood histogram defined as the maximizer of (7) with penalty defined by (8)-(11), we conduct a simulation study involving empirical risk estimation with respect to the losses (3) and (4) (for $p=1,2$). The choices we arrive at are given in Section 3. Then we compare our choices for the penalized maximum likelihood to other available methods in a separate simulation study using the same densities.

The performance of the methods is compared for 12 of the 28 test-bed densities introduced by Berlinet and Devroye (1994). We used densities 1 (uniform), 4 (double exponential), 11 (normal), 12 (lognormal), 21-24 (mixtures of normals) and 25-28 (various other multimodal densities). We denote these by f_1, \dots, f_{12} . We also added 4 histogram densities:

- 5 bin regular histogram:

$$f_{13}(x) := 0.15u_{[0,0.2]}(x) + 0.35u_{(0.2,0.4]}(x) + 0.2u_{(0.4,0.6]}(x) \\ + 0.1u_{(0.6,0.8]}(x) + 0.2u_{(0.8,1.0]}(x)$$

- 5 bin irregular histogram:

$$f_{14}(x) := 0.15u_{[0,0.13]}(x) + 0.35u_{(0.13,0.34]}(x) + 0.2u_{(0.34,0.61]}(x) \\ + 0.1u_{(0.61,0.65]}(x) + 0.2u_{(0.65,1.0]}(x)$$

- 10 bin regular histogram:

$$f_{15}(x) := 0.01u_{[0,0.1]}(x) + 0.18u_{(0.1,0.2]}(x) + 0.16u_{(0.2,0.3]}(x) \\ + 0.07u_{(0.3,0.4]}(x) + 0.06u_{(0.4,0.5]}(x) + 0.01u_{(0.5,0.6]}(x) \\ + 0.06u_{(0.6,0.7]}(x) + 0.37u_{(0.7,0.8]}(x) + 0.06u_{(0.8,0.9]}(x) \\ + 0.02u_{(0.9,1.0]}(x)$$

- 10 bin irregular histogram:

$$f_{16}(x) := 0.01u_{[0,0.02]}(x) + 0.18u_{(0.02,0.07]}(x) + 0.16u_{(0.07,0.14]}(x) \\ + 0.07u_{(0.14,0.44]}(x) + 0.06u_{(0.44,0.53]}(x) + 0.01u_{(0.53,0.56]}(x) \\ + 0.06u_{(0.56,0.67]}(x) + 0.37u_{(0.67,0.77]}(x) + 0.06u_{(0.77,0.91]}(x) \\ + 0.02u_{(0.91,1.0]}(x)$$

where $u_I := |I|^{-1}\mathbf{1}_I$ denotes the uniform density on an interval I . All densities are implemented in the R-package `benchden` (Mildenberger et al., 2009b) and are depicted in Figure 1. Note that Castellan’s main theorem 3.2 in Castellan (1999) does not apply to all densities considered here, since she assumes for instance that the density is bounded away from zero. We include a wide range of densities in order to explore the behavior of the procedure in cases not covered by theory. The sample sizes are 50,100,500,1000,5000 and 10000. We use 500 replications for each scenario and estimate the resulting risks as described in Section 5 using $\kappa = 5000$.

The methods compared in the simulations are:

- **B** Penalized maximum likelihood using penalty (9) with $c = 1$ and $\alpha = 1$.
- **R** Penalized maximum likelihood using random penalty (11) with $c = 1$ and $\alpha = 0.5$.
- **CV** Leave-one-out cross-validation using formula (11) given in Celisse and Robin (2008). We also tried formula (12) of Celisse and Robin (2008) for different values of p without finding a big difference.

Maximization is performed over a data-driven finest grid as described in Section 3 without restrictions on the minimum bin width.

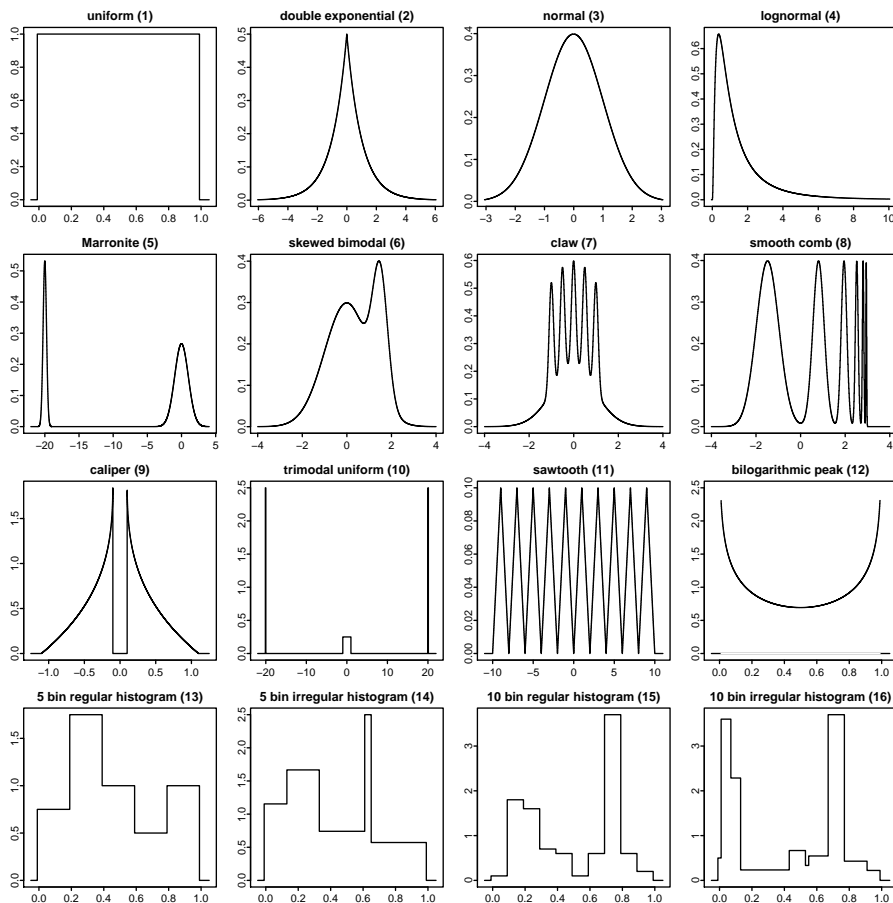


Figure 1: The densities used in the simulation study.

- Methods **B**, **R**, **CV** using the same data-driven grid but with the additional constraint that the minimum bin length allowed is $(x_{(n)} - x_{(1)}) \log^{1.5}(n)/n$. These are denoted by **B'**, **R'** and **CV'**.
- Methods **B**, **R**, **CV** but using a full optimization over a finest regular partition with bin width $(x_{(n)} - x_{(1)}) [n / \log^{1.5}(n)]^{-1}$. This is the grid considered in Castellan (1999), except that we slightly relax $\log^2(n)$ to $\log^{1.5}(n)$. These are denoted **B''**, **R''** and **CV''**.
- Penalized maximum likelihood using Information Criteria:
 - AIC** Akaike's Information Criterion (Akaike, 1973). The penalty is $\text{pen}_n^{\text{AIC}}(D) = (D - 1)$.
 - BIC** Bayesian Information Criterion (Schwarz, 1978). The penalty is $\text{pen}_n^{\text{BIC}}(D) = 0.5 \log(n)(D - 1)$.
- The taut-string method **TS** introduced by Davies and Kovac (2004). We use the function `pmden()` implemented in the R-package `fnonpar` (Davies and Kovac, 2008) with the default values except that we set `localsq=FALSE` as local squeezing of the tube does not give a ML histogram. The histogram is then constructed using the knots of the string

as the boundaries of the bins. This coincides with the derivative of the string except on intervals where the string switches from the upper to the lower boundary of the tube or vice versa.

- Regular histogram construction **BR** due to Birgé and Rozenholc (2006). The penalty is $\text{pen}_n^{\text{BR}}(D) = D + \log(D)^{2.5}$, where the log-likelihood is maximized over all regular partitions with $1, \dots, \lfloor n/\log n \rfloor$ bins.
- Combined approach: A regular histogram and an irregular histogram are constructed. After adding a constant such that all penalties are 0 for a histogram with 1 bin, the penalized log-likelihoods are compared. The histogram that has the larger value is chosen as the final estimate:
 - B*** Chooses between a regular histogram using **BR** and an irregular histogram using **B**
 - R*** Chooses between a regular histogram using **BR** and an irregular histogram using **R**.

Except for **TS**, all methods have been implemented in the R-package `histogram` (Mildenberger et al., 2009a). The program code to reproduce all simulation results will be available from the corresponding author’s website. For the discussion of the results, we focus on squared Hellinger risk. Table 1 shows the binary logarithms of relative squared Hellinger risks w.r.t. the best method for any given n and density: $\log_2(\widehat{R}_n^{\text{method}}/\widehat{R}_n^{\text{best}})$ for all 15 methods in the simulation study. Thus, a value of 0 means that the method was best in this particular setting and a value of 1 means that the risk of the method is twice as large as the risk of the best method. A table giving the absolute risk values as well as tables for the other risks can be found in the appendix. Generally, the L_1 and Hellinger risks behave similarly, while the results for L_2 are much less conclusive and no clear recommendation for using one of the methods under consideration can be given when minimizing the L_2 risk is the main aim. Table 8 in the appendix shows the modes of the number of bins chosen for all methods as well as the corresponding frequencies with which this number was chosen.

As was to be expected, the table shows no clear overall best method for all scenarios. The relative risk w.r.t. the best method is always smaller than 2 for our proposal **B***. In this sense, it can be seen as the best method. Let us remark that other methods could be better in particular situations. In many cases, **TS** or one of our proposals **B** and **R** is either the best or the binary logarithm of relative risk w.r.t. the best method is close to zero. These three methods are – except for **B*** – also the only ones in the simulation study for which this quantity is always strictly smaller than $\log_2 3 \approx 1.58$, meaning that the empirical risk is never greater than three times the risk achieved by the best method. The random penalty **R** seems to be slightly better than **B** in many cases, the most notable exception being the trimodal uniform density f_{10} for $n = 50$.

Cross-validation using the same set of partitions (i.e. a data-driven finest grid without further restrictions on minimum bin width) performs rather poorly, especially when the underlying density is a histogram ($f_1, f_{10}, f_{13}-f_{16}$) or when it has infinite peaks (f_{12}). Note that it is particularly bad for the uniform, which can be a major problem in many applications like grey level estimation of image differences. Relative performance of **CV** w.r.t. the best method becomes generally worse when sample size increases. If we compare the random and deterministic penalties and cross-validation for the case of full dynamic programming

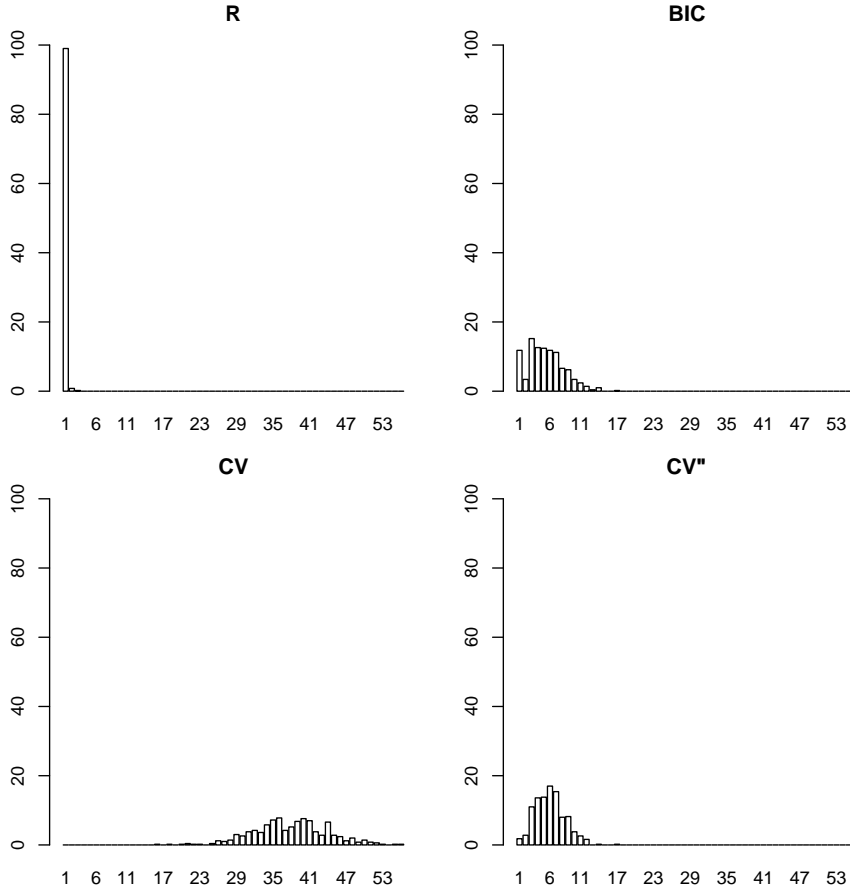


Figure 2: Barplots of number of bins chosen for the uniform density with $n = 500$ (percent of simulation runs).

the density has very sharp peaks (especially the trimodal uniform f_{10}). The intermediate case, i.e. using both penalties and cross-validation (\mathbf{B}' , \mathbf{R}' , \mathbf{CV}') for a data-driven finest grid but adding the constraint that bin widths have to be at least $(x_{(n)} - x_{(1)}) \log^{1.5}(n)/n$ could be suspected to give a compromise between the finest regular grid suggested by theory and the greedy algorithm for a data-driven grid. However, \mathbf{B}' and \mathbf{R}' share the catastrophic behavior of \mathbf{B}'' and \mathbf{R}'' at the trimodal uniform density f_{10} without offering a real improvement over \mathbf{B} and \mathbf{R} at the more well-behaved densities. On the other hand, \mathbf{CV}' is a good compromise between \mathbf{CV} and \mathbf{CV}'' , as it is in many cases either better than both or not far from the better of the two. It still shows bad behavior for the uniform and trimodal uniform densities. Table 8 shows that cross-validation has a pronounced tendency to choose histograms with a much larger number of bins than the penalized likelihood methods for all three sets of partitions. This is also illustrated by Fig. 2: For the uniform distribution with $n = 500$, our proposal \mathbf{R} often chooses only one bin, which is the best possible for the uniform. \mathbf{CV} chooses by far too many bins, resulting in a bad risk behavior. This is less extreme for \mathbf{CV}'' , but the number of bins chosen is still too large and the risk is more than 3 times larger than the best achieved for this setting.

The taut string method \mathbf{TS} shows a particularly good behavior in terms of Hellinger risk.

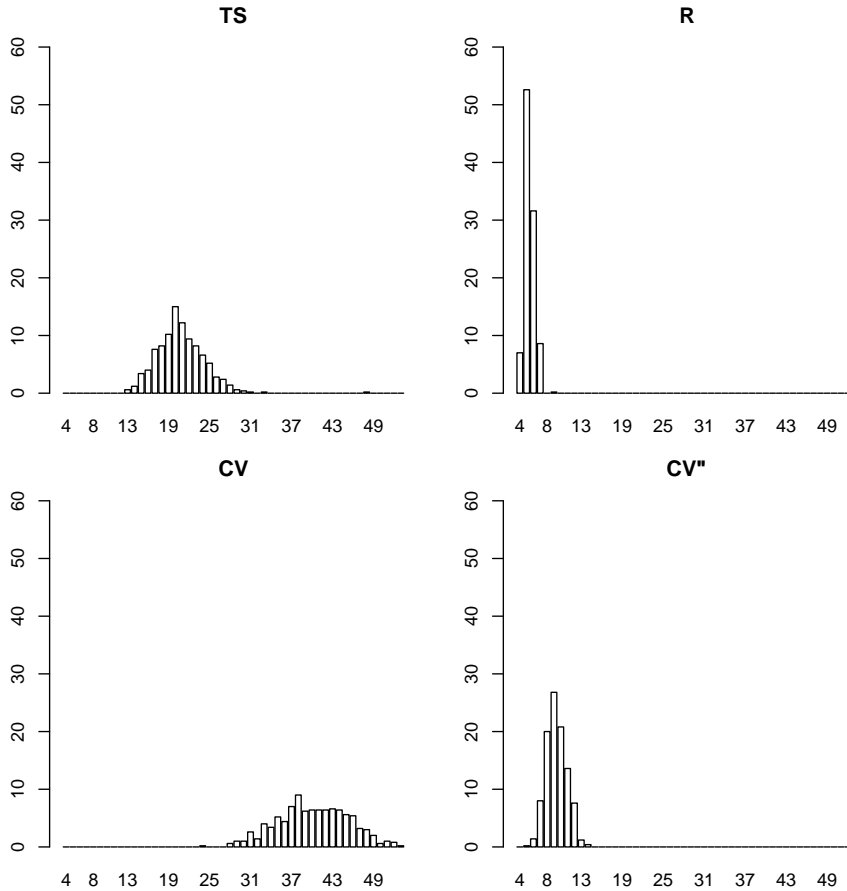


Figure 3: Barplots of number of bins chosen for the lognormal density $n = 500$ (percent of simulation runs).

One should note here that it does not control the number of bins but the modality of the estimate, thereby avoiding overfitting while still being able to choose a large number of bins to give sufficient detail. An example is given in Fig. 3, where the number of bins chosen for the lognormal distribution f_4 (with $n = 500$) is shown. Of the four methods shown, **CV** performs worst, choosing again a large number of bins. **TS** is best in this scenario and uses a larger number of bins than both **R** and **CV''**. One should note that the taut string method was originally derived for different aims than achieving a good behavior w.r.t. a given loss function (Davies and Kovac, 2004), and many questions regarding behavior in a more classical framework remain open.

Using **AIC** leads to very bad results. It has already been shown theoretically in Castellan (1999) and Massart (2007) and from a more practical point of view in Birgé and Rozenholc (2006) that **AIC** underpenalizes even for regular histograms. The tendency to overfit becomes even worse when using **AIC** for irregular histograms, since now there are many partitions with the same number of bins, leading to problems similar to those arising in multiple testing. As argued by Castellan (1999) and Massart (2007), in this case, an additional penalty is needed. Table 8 shows that the number of bins chosen by **AIC** on average is very often the largest among all methods considered. In many cases, the ratio of the Hellinger risk to the best risk

achieved by any method in the simulation study is at least 4, often even much larger. **BIC** is a criterion which does not aim for a good control of risk but at asymptotically identifying the "smallest true model", if it exists. It shares the problem of **AIC** of not taking into account the number of models of the same size. It shows some good behavior in particular for small sample sizes that deteriorates when samples become larger. Particularly noteworthy is the bad performance for "simple" models like the uniform (which is also shown in Fig. 2) and the 5 bin regular histogram density f_{13} .

The regular histogram method **BR**, which improves on Akaike's penalization, is the best method for f_3 , f_{11} and f_{13} , at least when the sample size is not very small. This shows that the greater flexibility of an irregular histogram over a regular one may be outweighed by the greater difficulty in choosing a good partition, as was already remarked by Birgé and Rozenholc (2006). Regular histograms are inferior for spatially inhomogeneous densities like f_4 and f_{10} .

The simulations show that one can successfully combine the advantages of regular and irregular histograms (**B*** and **R***). The Hellinger risk for **B*** is always within twice the risk of the best method for a given situation. While **R*** is even better in most cases, it shares the bad behaviour of **BR** for f_{10} and $n = 50$. Table 2 shows that for larger sample sizes both **B*** and **R*** almost always choose an irregular histogram for densities where this is advantageous (the spatially inhomogeneous densities f_4 , f_5 , f_{10} , f_{14} and f_{16}) and a regular partition in most of the other cases.

n	method	f_1	f_2	f_3	f_4	f_5	f_6	f_7	f_8	f_9	f_{10}	f_{11}	f_{12}	f_{13}	f_{14}	f_{15}	f_{16}
50	B*	0.08	0.08	0.05	0.75	1.00	0.05	0.14	0.07	0.11	1.00	0.15	0.39	0.06	0.11	0.22	0.51
	R*	0.02	0.13	0.10	0.68	1.00	0.11	0.27	0.06	0.12	0.00	0.02	0.16	0.03	0.07	0.22	0.43
100	B*	0.04	0.07	0.04	0.74	1.00	0.11	0.23	0.02	0.19	1.00	0.12	0.43	0.01	0.12	0.24	0.65
	R*	0.00	0.10	0.10	0.64	1.00	0.24	0.40	0.05	0.36	1.00	0.05	0.28	0.01	0.11	0.33	0.65
500	B*	0.03	0.00	0.00	0.73	1.00	0.01	0.03	0.00	0.42	1.00	0.00	0.47	0.00	0.69	0.36	0.99
	R*	0.00	0.00	0.00	0.47	1.00	0.04	0.07	0.00	0.53	1.00	0.00	0.45	0.00	0.75	0.48	0.99
1000	B*	0.02	0.00	0.00	0.89	1.00	0.00	0.00	0.00	0.58	1.00	0.00	0.41	0.00	0.96	0.49	1.00
	R*	0.01	0.00	0.00	0.63	1.00	0.01	0.00	0.00	0.64	1.00	0.00	0.41	0.00	0.98	0.57	1.00
5000	B*	0.01	0.00	0.00	1.00	1.00	0.00	0.00	0.00	0.65	1.00	0.00	0.26	0.00	1.00	0.43	1.00
	R*	0.00	0.00	0.00	1.00	1.00	0.00	0.00	0.00	0.70	1.00	0.00	0.28	0.00	1.00	0.47	1.00
10000	B*	0.00	0.00	0.00	1.00	1.00	0.00	0.00	0.00	0.66	1.00	0.00	0.18	0.00	1.00	0.35	1.00
	R*	0.00	0.00	0.00	1.00	0.98	0.00	0.00	0.00	0.71	1.00	0.00	0.20	0.00	1.00	0.39	1.00

Table 2: Frequency of choosing an irregular partition.

To summarize, we propose a practical method of irregular histogram construction inspired by theoretical works by Barron et al. (1999), Castellan (1999, 2000) and Massart (2007). It can be easily implemented using a dynamic programming algorithm and it performs well for a wide range of different densities and sample sizes, even for some cases not covered by the underlying theory. Performance is shown to be improved when combined with the regular histogram approach proposed in Birgé and Rozenholc (2006). All procedures proposed here are available in the R-package **histogram** (Mildenberger et al., 2009a).

Acknowledgments

This work has been supported in part by the Collaborative Research Centers "Reduction of Complexity in Multivariate Data Structures" (SFB 475) and "Statistical Modelling of Nonlinear Dynamic Processes" (SFB 823) of the German Research Foundation (DFG). The authors also wish to thank Henrike Weinert for discussions and programming in earlier stages of the work as well as two anonymous referees for giving comments that greatly improved the paper.

References

- Akaike, H., 1973. A new look at the statistical model identification. *IEEE Transactions on Automatic Control* 19, 716-723.
- Barron, A., Birgé, L., Massart, P., 1999. Risk bounds for model selection via penalization. *Probability Theory and Related Fields* 113, 301-413.
- Berlinet, A., Devroye, L., 1994. A comparison of kernel density estimates. *Publications de l'Institut de Statistique de l'Université de Paris* 38, 3-59.
- Birgé, L., Rozenholc, Y., 2006. How many bins should be put in a regular histogram? *ESAIM: Probability and Statistics* 10, 24-45.
- Blanchard, G., Schäfer, C., Rozenholc, Y., Müller, K.-R., 2007. Optimal dyadic decision trees. *Machine Learning* 66, 209-241.
- Castellan, G., 1999. Modified Akaike's criterion for histogram density estimation. Technical Report 99.61, Université de Paris-Sud.
- Castellan, G., 2000. Sélection d'histogrammes à l'aide d'un critère de type Akaike. *Comptes rendus de l'Académie des sciences Paris* 330, Série I, 729-732.
- Catoni, O., 2002. Data compression and adaptive histograms, in: Cucker, F., Rojas J.M. (Eds.), *Foundations of Computational Mathematics, Proceedings of the Smalefest 2000*, World Scientific, Singapore, pp. 35-60.
- Celisse, A., Robin, S., 2008. Nonparametric density estimation by exact leave-p-out cross-validation. *Computational Statistics and Data Analysis* 52, 2350-2368.
- Chen, X.R., Zhao, L.C., 1987. Almost sure L_1 -norm convergence for data-based histogram density estimators. *Journal of Multivariate Analysis* 21, 179-188.
- Comte, F., Rozenholc, Y., 2004. A new algorithm for fixed design regression and denoising. *Annals of the Institute of Statistical Mathematics* 56, 449-473.
- Davies, P. L., Gather, U., Nordman, D. J., Weinert, H., 2008. A comparison of automatic histogram constructions. *ESAIM: Probability and Statistics* 13, 181-196.
- Davies, P. L., Kovac, A., 2004. Densities, spectral densities and modality. *The Annals of Statistics* 32, 1093-1136.
- Davies, P.L., Kovac, A., 2008. `ftnonpar`: Features and strings for nonparametric regression. R package version 0.1-83.
- Devroye, L., Györfi, L., 1985. *Nonparametric density estimation: the L_1 view*. Wiley, New York.
- Devroye, L., Lugosi, G., 2004. Bin width selection in multivariate histograms by the combinatorial method, *Test* 13, 129-145.
- Engel, J., 1997. The multiresolution histogram. *Metrika* 46, 41-57.

- Hartigan, J.A., 1996. Bayesian histograms, in: Bernardo, J.M., Berger, J.O., Dawid, A.P., Smith, A.F.M. (Eds.), *Bayesian Statistics 5*, Oxford University Press, Oxford, pp. 211-222.
- Kanazawa, Y., 1988. An optimal variable cell histogram. *Communications in Statistics - Theory and Methods* 17, 1401-1422.
- Kanazawa, Y., 1992. An optimal variable cell histogram based on the sample spacings. *The Annals of Statistics* 20,219-304.
- Klemelä, J., 2007. Density estimation with stagewise optimization of the empirical risk. *Machine Learning* 67, 169-195.
- Klemelä, J., 2009. Multivariate histograms with data-dependent partitions. *Statistica Sinica* 19, 159-176.
- Kogure, A., 1986. Optimal cells for a histogram. PhD thesis, Yale University.
- Kogure, A., 1987. Asymptotically optimal cells for a histogram. *The Annals of Statistics* 15, 1023-1030.
- Kontkanen, P., Myllymäki, P., 2007. MDL histogram density estimation. In: Meila M., Shen S. (Eds.), *Proc. 11th International Conference on Artificial Intelligence and Statistics (AISTATS 2007)*, Puerto Rico, March 2007. <http://www.stat.umn.edu/~aistat/proceedings/start.htm>
- Lugosi, G., Nobel, A., 1996. Consistency of data-driven histogram methods for density estimation and classification. *The Annals of Statistics* 24, 687-706.
- Massart, P., 2007. Concentration inequalities and model selection. *Lecture Notes in Mathematics* Vol. 1896, Springer, New York.
- Mildenberger, T., Rozenholc, Y., Zasada, D., 2009a. histogram: Construction of regular and irregular histograms with different options for automatic choice of bins. R package version 0.0-20.
- Mildenberger, T., Weinert, H., Tiemeyer, S., 2009b. benchden: 28 benchmark densities from Berlinet/Devroye (1994). R package version 1.0.2.
- Rissanen, J., Speed, T. P., Yu, B., 1992. Density estimation by stochastic complexity. *IEEE Transactions on Information Theory* 38, 315-323.
- Schwarz, G., 1978. Estimating the dimension of a model. *The Annals of Statistics* 6, 461-464.
- Zhao, L.C, Krishnaiah, P.R., Chen, X.R., 1988. Almost sure L_r -norm convergence for data-based histogram estimates. *Theory of Probability and its Applications* 35, 396-403.

	n	B	R	CV	B'	R'	CV'	B''	R''	CV''	AIC	BIC	TS	BR	B*	R*
f ₁	50	3.88	0.22	4.46	0.11	0.12	1.83	0	0	1.22	7.79	7.55	0.53	0.38	3.88	0.44
	100	3.2	0.09	4.96	0.07	0.09	2.6	0	0.03	1.76	8.58	8.12	0.44	0.38	3.2	0.41
	500	4.28	0.15	6.85	0.03	0.06	4.29	0	0.02	3.26	9.93	8.83	0.68	0.46	4.31	0.5
	1000	3.19	0.13	7.17	0.05	0.05	4.86	0	0.02	4.01	9.76	7.99	0.44	0.39	3.23	0.47
	5000	1.98	0	7.58	0	0	5.56	0	0	5.71	9.56	6.04	0.44	0.43	1.89	0.43
	10000	1.31	0	7.8	0	0	5.88	0	0	6.63	9.47	5.39	0.66	0.62	1.6	0.62
f ₂	50	4.41	0.49	2.24	0.5	0.46	0.11	0.73	0.68	0	5.76	5.47	0.62	0.17	1.04	0.22
	100	3.2	0.69	2.37	0.7	0.67	0.13	0.8	0.81	0	6.32	5.86	0.02	0.3	0.33	0.35
	500	1.99	1.02	3.74	0.99	1.01	0.32	1.05	1.1	0.13	6.5	5.29	0	0.52	0.53	0.52
	1000	1.69	1.17	3.73	1.12	1.16	0.43	1.15	1.22	0.2	6.02	4.13	0	0.61	0.61	0.61
	5000	1.39	1.39	3.46	1.31	1.39	0.71	1.3	1.39	0.41	4.71	2.34	0	0.69	0.69	0.69
	10000	1.43	1.51	3.28	1.38	1.51	0.81	1.37	1.49	0.6	4.16	2.03	0	0.78	0.78	0.78
f ₃	50	4.87	0.63	2.5	0.7	0.63	0.35	0.72	0.69	0	6.4	6.16	0.84	0.29	2.51	0.31
	100	4.04	0.51	2.8	0.57	0.5	0.43	0.65	0.58	0	7.12	6.51	0.27	0.13	0.18	0.16
	500	1.99	1.02	4.15	1.05	1.02	1.16	1.1	1.04	0.44	7.16	6.05	0	0.32	0.32	0.32
	1000	1.58	1.22	4.32	1.24	1.22	1.46	1.22	1.19	0.69	6.81	5.25	0	0.42	0.42	0.42
	5000	1.74	1.5	4.08	1.52	1.5	1.8	1.51	1.5	1.45	5.65	3.12	0	0.45	0.45	0.45
	10000	1.65	1.6	3.9	1.6	1.6	1.81	1.58	1.57	1.83	5.1	2.51	0	0.49	0.49	0.49
f ₄	50	5.23	0.5	3.02	0.78	0.79	0.53	0.71	0.72	0.37	6.23	5.99	0	0.45	5.22	0.46
	100	2.97	0.65	2.92	0.88	0.91	0.58	0.77	0.81	0.4	6.5	5.97	0	0.64	2.7	0.64
	500	0.61	0.56	3.18	0.67	0.75	0.45	0.63	0.74	0.32	5.81	4.14	0	0.67	0.61	0.58
	1000	0.19	0.26	2.87	0.69	0.76	0.5	0.67	0.76	0.37	4.88	2.53	0	0.74	0.24	0.4
	5000	1.01	1.38	3.29	2.26	2.39	2.13	1.89	2.09	1.59	4.22	1.8	0	2.77	1.01	1.38
	10000	1.07	1.53	3.1	2.17	2.34	2	1.82	2.08	1.42	3.66	1.55	0	2.69	1.07	1.53
f ₅	50	3	0	2.03	1.43	1.43	1.43	1.13	1.13	1.12	5.06	4.75	0.15	0.89	3	0
	100	2.5	0.34	1.94	1.94	1.94	1.94	1.56	1.56	1.53	5.41	4.8	0	0.57	2.5	0.34
	500	1.79	0.99	3.44	2.15	2.17	2.01	2.38	2.4	2.23	6.03	4.68	0	1.64	1.79	0.99
	1000	1.66	1.27	3.58	1.97	1.98	1.69	1.69	1.73	1.43	5.71	3.66	0	1.4	1.66	1.27
	5000	1.54	1.61	3.39	1.57	1.68	1.21	1.46	1.65	0.65	4.53	2.23	0	1.56	1.54	1.61
	10000	1.61	1.68	3.17	1.55	1.68	1.08	1.56	1.72	0.64	3.94	1.98	0	1.55	1.61	1.68
f ₆	50	3.41	0.33	2.23	0.4	0.32	0.47	0.33	0.29	0	6.19	5.87	0.69	0.33	3.27	0.33
	100	2.15	0.25	2.66	0.31	0.25	0.56	0.46	0.39	0	6.45	5.94	0.47	0.33	0.55	0.31
	500	1.15	0.51	3.59	0.54	0.51	0.72	0.55	0.5	0.05	6.46	5.25	0	0.1	0.1	0.11
	1000	1	0.57	3.6	0.6	0.57	0.85	0.62	0.59	0.14	6.05	4.46	0	0.06	0.06	0.07
	5000	1.28	1.19	3.84	1.21	1.19	1.58	1.17	1.15	1.28	5.45	2.83	0	0.55	0.55	0.55
	10000	1.48	1.42	3.71	1.44	1.42	1.67	1.44	1.42	1.72	4.94	2.35	0	0.62	0.62	0.62
f ₇	50	1.72	0.13	1.78	0.21	0.12	0.11	0.34	0.31	0	5.04	4.7	0.64	0.19	0.96	0.14
	100	1.89	0.03	1.97	0.02	0	0.28	0.26	0.24	0.2	5.72	5.27	0.26	0.23	1.3	0.15
	500	1.21	0.58	2.46	0.57	0.55	0	0.61	0.6	0	5.33	4.19	0.16	0.06	0.08	0.09
	1000	1.29	0.81	2.72	0.78	0.8	0.32	1.18	1.18	0	5.1	3.59	0.05	0.18	0.18	0.18
	5000	1.23	1.18	2.71	1.17	1.18	0.65	1.15	1.15	0.41	4.01	2.03	0	0.51	0.51	0.51
	10000	1.37	1.34	2.59	1.31	1.34	0.79	1.36	1.39	0.62	3.53	1.74	0	0.59	0.59	0.59
f ₈	50	2.77	0.31	1.28	0.3	0.27	0.08	0.32	0.29	0	4.73	4.47	0.45	0.13	0.73	0.14
	100	3.35	0.36	1.47	0.39	0.31	0	0.61	0.53	0.07	5.34	4.94	0.46	0.04	0.17	0.06
	500	1.65	0.45	2.3	0.55	0.51	0.33	0.72	0.65	0	5.1	3.99	0.21	0.07	0.07	0.07
	1000	1.25	0.37	2.25	0.47	0.42	0.28	0.57	0.52	0	4.6	3.22	0	0.06	0.06	0.06
	5000	1.02	0.85	2.13	0.95	0.92	0.67	0.94	0.91	0.42	3.61	1.95	0	0.49	0.49	0.49
	10000	1.22	1.12	2.06	1.17	1.14	0.82	1.16	1.12	0.74	3.09	1.59	0	0.82	0.82	0.82
f ₉	50	3.56	0.23	1.23	0.18	0.17	0.01	0.16	0.15	0	4.89	4.66	0.28	0.09	2.05	0.13
	100	3.8	0.41	2.14	0.47	0.32	0	1.07	1.03	0.73	5.66	5.22	0.39	0.37	2.66	0.31
	500	2.35	0.53	2.97	0.52	0.49	0.29	1.28	1.26	0.96	5.98	4.99	0	0.48	1.54	0.43
	1000	1.11	0.61	3.14	0.6	0.58	0.48	1.35	1.34	1.01	5.55	4.02	0	0.76	1.01	0.57
	5000	0.89	0.82	3	0.81	0.81	0.82	1.49	1.49	1.33	4.5	2.2	0	1	0.81	0.78
	10000	0.92	0.85	2.78	0.85	0.85	0.85	1.51	1.5	1.5	3.92	1.65	0	1.12	0.86	0.82
f ₁₀	50	2.07	0.28	1.82	1.39	1.39	1.39	1.39	1.39	3.92	3.51	0	1.32	2.07	1.32	
	100	2.23	0	2.82	3.15	3.15	3.15	3.12	3.12	3.12	5.4	4.49	0.84	2.99	2.23	0
	500	0.55	0	4.4	5.44	5.44	5.45	5.37	5.37	5.37	6.27	3.68	1.13	5.14	0.55	0
	1000	0.21	0	4.84	6.43	6.43	6.43	6.24	6.24	6.24	6.24	3.37	1.2	5.82	0.21	0
	5000	0.03	0	5.46	8.78	8.78	8.78	7.8	7.8	7.8	6.06	2.39	1.54	0.4	0.03	0
	10000	0.01	0	5.65	9.79	9.79	9.79	6.69	6.69	6.69	6.02	2.14	1.68	0.43	0.01	0
f ₁₁	50	1.6	0.06	2	0.02	0.03	0.38	0	0	0.23	4.72	4.43	1.84	0.06	1.59	0.08
	100	0.54	0.04	1.61	0	0.01	0.39	0	0	0.21	4.95	4.56	1.27	0.04	0.56	0.07
	500	2.01	0.68	1.95	0.9	0.73	0.34	1.03	1.03	0.14	4.72	3.76	0.62	0	0	0
	1000	1.64	0.66	2.03	0.67	0.65	0.48	0.88	0.72	0	4.35	3.14	0.39	0	0	0
	5000	1.24	0.97	1.68	1.01	0.96	0.56	1.01	0.96	0.32	3.04	2	0.08	0	0	0
	10000	1.22	1.08	1.46	1.12	1.07	0.69	1.14	1.08	0.51	2.32	1.49	0	0.01	0.01	0.01
f ₁₂	50	3.19	0.32	2.98	0.02	0.03	0.19	0	0	0.09	5.27	5.01	0.09	0.1	3.18	0.3
	100	2.22	0.4	2.26	0.24	0.24	0.21	0.22	0.23	0	5.39	4.82	0.27	0.18	2.1	0.33
	500	1.75	0.29	2.87	0.23	0.23	0.48	0.22	0.22	0	5.69	4.7	0.02	0.24	1.2	0.29
	1000	1.19	0.38	2.82	0.34	0.34	0.58	0.38	0.37	0.08	5.34	4.06	0	0.26	0.89	0.33
	5000	0.56	0.34	2.58	0.34	0.34	0.59	0.34	0.34	0.44	4.25	2.01	0	0.36	0.45	0.37
	10000	0.45	0.31	2.38	0.3	0.31	0.5	0.3	0.3	0.68	3.7	1.34	0	0.38	0.42	0.38
f ₁₃	50	2.27	0.2	2.44	0.14	0.15	0.08	0.09	0.09	0	5.94	5.73	0.48	0.06	2.19	0.08
	100	2.33	0.63	2.35	0.61	0.59	0.37	0.57	0.56	0.01	6.21	5.77	0.62	0	0.09	0.02
	500	3.23	1.9	5.13	1.94	1.88	2.56	2.11	2.07	2	8.31	7.39	1.6	0	0	0
	1000	3.13	1.46	5.55	1.5	1.45	3.19	1.2	1.14	2.3	8.17	6.67	1.6	0	0	0
	5000	2.2	1.37	6.13	1.37	1.37	4.15	1.81	1.81	4.16	7.9	4.71	1.77	0	0	0
	10000	1.93	1.4	6.3	1.38	1.39	4.43	2.05	2.05	4.95	7.79	4.38	1.93	0	0	0
f ₁₄	50	4.61	0.24	2.48	0.2	0.19	0.2	0.19	0.19	0	5.92	5.75	0.4	0.06	4.31	0.1
	100	3.14	0.32	1.9	0.28	0.27	0.29	0.4	0.36	0.15	5.65	5.21	0.38	0	3.07	0.07
	500	2.35	0.05	3.28	0.03	0	0.72	0.72	0.7	0.63	6.41	5.5	0.09	1	1.65	0.37
	1000	1.31	0.2	3.83	0.22	0.2	1.44	1.51	1.08	1.32	6.4	4.74	0	1.01	1.34	0.24
	5000	0.86	0	4.53	0	0	2.42	1.74	1.74	2.72	6.3	3.11	0.46	2.1	0.86	0
	10000	0.53	0.01	4.76	0	0	2.79	1.17	1.17	3.32	6.21	2.73	0.71	2.91	0.53	0.01
f ₁₅	50	2.84	0.02	1.1												

