

On the number of bins in a rank histogram

Claudio Heinrich *

*Norwegian Computing Center Oslo,
P.O. Box 114 Blindern, NO-0314 Oslo, Norway*

October 6, 2020

Abstract

Rank histograms are popular tools for assessing the reliability of meteorological ensemble forecast systems. A reliable forecast system leads to a uniform rank histogram, and deviations from uniformity can indicate miscalibrations. However, the ability to identify such deviations by visual inspection of rank histogram plots crucially depends on the number of bins chosen for the histogram. If too few bins are chosen, the rank histogram is likely to miss miscalibrations; if too many are chosen, even perfectly calibrated forecast systems can yield rank histograms that do not appear uniform. In this paper we address this trade-off and propose a method for choosing the number of bins for a rank histogram. The goal of our method is to select

*email: claudio.heinrich@nr.no

The author would like to thank Thordis Thorarinsdottir for helpful discussions, and two anonymous reviewers for their suggestions that helped to substantially improve the paper. He thanks his colleagues from the Norwegian Computing Center for labeling many histograms and is grateful to the Norwegian Computing Center for its financial support.

16 a number of bins such that the intuitive decision whether a histogram is uniform or
17 not is as close as possible to a formal statistical test. Our results indicate that it is
18 often appropriate to choose fewer bins than the usual choice of ensemble size plus
19 one, especially when the number of observations available for verification is small.

20 *Keywords: forecast verification, rank histograms, statistical testing*

21 **1 Introduction**

22 Rank histograms are widely used diagnostic tools for calibration assessment of forecasts
23 in meteorology. The underlying idea to consider the rank of the observation within a
24 predictive ensemble was proposed independently by Anderson (1996), Hamill and Colucci
25 (1997) and Talagrand et al. (1997). If the prediction system is well-calibrated (or reliable),
26 the rank of the observation within the ensemble is approximately uniformly distributed.
27 Deviations from uniformity indicate different types of miscalibration, for example, sloped
28 histograms indicate bias, and \cup - or \cap -shaped histograms indicate under- and overdispersion,
29 respectively. Rank histograms were originally applied to univariate forecasts, however,
30 several generalizations towards multivariate forecasts exist (Wilks, 2004; Thorarinsdottir
31 et al., 2016; Ziegel and Gneiting, 2014).

32 As pointed out by Wand (1997) in a different context, choosing the number of bins in
33 a histogram is generally a trade-off: More bins lead to a more detailed histogram while
34 also making it more susceptible to random fluctuations. In particular, when the available
35 number of forecast-observation pairs is small, the appearance can change quite dramatically
36 with different bin numbers, see Figure 1. The goal of this work is to address this trade-off
37 and provide guidance regarding the choice of a bin size in a rank histogram. We focus on the
38 case where only a relatively small number of forecast-observation pairs are available, say less
39 than 200. In this case, too many bins can lead to an over-interpretation of the histogram's
40 appearance. This situation occurs, for example, frequently in seasonal forecasting where

41 variables are averaged over long time-spans, leading to a drastically reduced number of
42 available observations, see Van Schaeybroeck and Vannitsem (2018).

43 When an ensemble forecast with m ensemble members is considered, the observation
44 rank can take values between 1 and $m + 1$. It is therefore intuitive and common practice
45 to use $m + 1$ bins for rank histograms, each bin corresponding to a single rank (e.g. Wilks
46 (2019)). We show how to construct rank histograms with any bin number such that every
47 bin accounts for the same number of ranks. This is necessary in order to address the above-
48 mentioned trade-off, and useful in its own right. It can, for example, be quite difficult to
49 compare histograms with different bin numbers. Therefore, when forecast systems with
50 different ensemble sizes are compared, it is useful to choose the same bin number for all of
51 them.

52 Our approach to finding ‘good’ bin numbers acknowledges that rank histograms are first
53 and foremost used for exploratory data analysis. They are typically generated and inspected
54 by scientists who then intuitively decide whether they look sufficiently uniform or not.
55 This implies, in particular, that good bin numbers are not an inherent statistical property
56 of the data, but require assumptions on scientists’ intuitive decisions. We will assume
57 that such decisions directly depend on the distance between the observed histogram and a
58 perfectly flat histogram, and that larger distances are more likely to lead to a rejection. This
59 constitutes a necessary oversimplification, which in particular does not take characteristic
60 shapes such as slopes or U-shapes into account. An empirical study is conducted where
61 several statisticians label more than 400 histograms as uniform or not, in order to assess
62 to what extent our assumption is justified.

63 Subject to this assumption, the bin number can be chosen to make the scientists’ de-
64 cision approximate the decision of a formal statistical test for uniformity. The underlying
65 intuition is that, when based on uniformly distributed data, histograms with fewer bins
66 tend to look flatter than those with many bins. Therefore, reducing the number of bins re-
67 duces the probability of an intuitive false reject (type I error). At the same time, it reduces

68 the amount of detail depicted by the histogram and therefore increases the probability of
69 a false accept (type II error). In this sense the trade-off in choosing the number of bins di-
70 rectly relates to the trade-off made in statistical testing when choosing a significance level,
71 which balances the probabilities of the two types of errors. We formalize this intuitive link,
72 which then allows us to associate a chosen number of bins with a probability for a false
73 reject. Establishing this link requires the selection of a subjective ‘acceptance threshold’,
74 indicating how large deviations from uniformity are deemed acceptable by the inspecting
75 scientist. We use the results from our empirical study to provide approximations for the
76 average scientists’ acceptance threshold.

77 There are several different tests for uniformity that have been applied in the context of
78 rank histograms. Besides the classical χ^2 -test, Delle Monache et al. (2006) considered a test
79 based on the so-called reliability index, and Taillardat et al. (2016) used a test based on an
80 entropy test statistic. These three tests have recently been compared by Wilks (2019). For
81 all three of them, the test statistic can be interpreted as a distance between the observed
82 and a perfectly flat histogram. This allows us to establish and analyze the above-mentioned
83 link between the choice of bin number and a statistical test for any of the three tests.

84 Given a significance level α and the number of available observations n , our methodology
85 selects a bin number k such that, when inspecting a histogram with k bins, a scientists’
86 intuitive decision closely approximates the test at significance level α . This bin number is
87 in most cases similar (and often identical) for the three different tests, which provides a
88 sanity check for our methodology: The selected bin number should lead to a false rejection
89 by the scientist with probability α , regardless of the test used in the derivation.

90 Our results show that when only few observations are available, even histograms with a
91 moderate number of bins lead to high probabilities of an intuitive false reject. For example,
92 when 100 observations are available, choosing more than 9 bins results in a probability of
93 more than 33% of a false reject; for 60 available observations, this probability is exceeded
94 when more than 6 bins are chosen.

95 Optimality criteria for histogram bin numbers and bin widths have been widely dis-
96 cussed in the literature, see e.g. Scott (1979); He and Meeden (1997); Muto et al. (2019)
97 and Knuth (2019). However, these criteria have generally been developed in a different
98 context and under assumptions that make them inappropriate for rank histograms. They
99 mostly focus on histograms as tools for estimating probability densities with the aim of
100 finding the number of bins that minimizes a distance (often the mean integrated squared
101 error) between the underlying density and the histogram of the data. In this context it is
102 commonly assumed that the density is continuous and sufficiently smooth over an inter-
103 val. Some early work even assumes approximately normally distributed data (Scott, 1979;
104 Sturges, 1926). These assumptions are not met for rank histograms based on discrete data.
105 Moreover, the vast majority of results derived in this strand of literature are of asymptotic
106 nature and therefore assume n to be large, in contrast to our assumptions. Thirdly, the
107 derived binnings are often data driven, i.e. the bin number depends on properties of the
108 data beyond the sample size n , such as for example the sample variance. In the context
109 of rank histograms, which are commonly used to compare different forecast systems this is
110 not desirable as all the histograms should have the same number of bins.

111 The remainder of the paper is organized as follows. In Section 2 we show how histograms
112 with any bin number can be derived from an m -member ensemble forecast. Section 3
113 describes the approach we take to relate the bin number to statistical tests. The optimal
114 bin number requires the choice of a subjective acceptance threshold. In Section 4 we present
115 an empirical study and use it to derive an approximation of this acceptance threshold. In
116 Section 5 we use the developed algorithm to find good bin numbers for a range of different
117 data sizes. Section 6 analyzes the rejection probability for histograms with the optimal bin
118 number under non-uniform distributions. Section 7 provides a discussion of the results and
119 Section 8 concludes.

120 2 Changing the bin number for rank histograms

121 When computing rank histograms for an ensemble forecast with m members the observation
122 ranks r_1, \dots, r_n take values in $\{1, \dots, m + 1\}$. Therefore, the default is to use a histogram
123 with $m + 1$ bins, each bin containing the counts for one rank only. It is straightforward
124 to instead generate a rank histogram with $k < m + 1$ bins, as long as k divides $m + 1$.
125 Then, the first bin accounts for the first $(m + 1)/k$ ranks, and so on. However, this is quite
126 restrictive, especially as $m + 1$ is prime for some popular ensemble sizes such as 10, 30 and
127 100. As argued in the introduction, free choice of the bin number k is desirable and we
128 show in the following how this can be achieved.

129 The problem that arises when k does not divide $m + 1$ is that some bins get assigned
130 more ranks than others. Take the simple example of $m = 2$ where the observed ranks take
131 the values 1, 2, 3, and assume we want to plot a histogram with only two bins. Then, the
132 question arises whether the counts of rank 2 should be placed in the first or the second
133 bin. Both options lead to skewed histograms even if the ranks are perfectly uniformly
134 distributed. This issue can be resolved by randomization. For each count of rank 2 we
135 simply flip a coin and place it in the first bin if the coin shows tails, and in the second bin
136 otherwise. When moving beyond this simple example, the randomization becomes more
137 involved, as it needs to account for the fraction of overlap between bins and ranks: Say,
138 for example, we have ranks 1, ..., 5 and want to consider 4 bins, then the first bin should
139 account for all counts of rank 1 and $\frac{1}{4} - \frac{1}{5} = \frac{1}{20}$ th of the counts for the second bin. For each
140 count of rank 2 we should, therefore, flip a ‘skewed’ coin showing heads with probability
141 $1/20$, and place it in the first bin if heads comes up, and in the second bin otherwise.

This procedure can be simplified as follows. Consider ranks $r_1, \dots, r_n \in \{1, \dots, m + 1\}$
and compute the transformed ranks

$$\tilde{r}_i := \frac{r_i - 1 + U_i}{m + 1}, \quad (2.1)$$

142 where U_1, \dots, U_n are independent random variables, uniformly distributed on the interval

143 $[0,1]$. The transformed ranks can take any value between 0 and 1, and we can now generate
144 a histogram with any number of bins k in the usual way, i.e. the j th bin counts the number
145 of transformed ranks in the interval $[\frac{j-1}{k}, \frac{j}{k}]$. The random variables U_i take the roles of the
146 coinflips above, however, since they are uniformly distributed on $[0,1]$ they automatically
147 account for the fraction of overlap between the k bins and the $m + 1$ ranks.

148 The histogram of the modified ranks can be interpreted exactly as the original rank
149 histogram. In fact, the randomization only has an effect if a bin number that does not
150 divide $m + 1$ is chosen, otherwise the two histograms are identical. After this replacement,
151 histograms with any number of bins can be considered. Flatness is preserved and if the
152 original ranks are uniformly distributed so are the transformed ranks. Note that this also
153 allows us to consider histograms with more than $m + 1$ bins. If we, for example, consider
154 $k = 2(m + 1)$ bins, each count of rank 1 is simply assigned either to the first or to the
155 second bin with equal probability.

156 This randomization is closely related to randomized versions of the probability integral
157 transform (PIT), see e.g. Smith (1985). When a probability forecast with distribution
158 function F is issued and observation y materializes, the PIT simply considers $F(y)$. If
159 the forecast system is reliable and F is continuous, $F(y)$ follows a uniform distribution.
160 Therefore a histogram of $F_1(y_1), \dots, F_n(y_n)$, for a sequence of observations and associated
161 predictions, is a diagnostic tool for assessing the calibration of a probability forecast system,
162 very similar to rank histograms for ensemble forecast systems. If the probability forecast
163 F is not continuous, Smith (1985) suggested to modify the PIT by randomly filling in the
164 jumps: That is, whenever the observation y is at a discontinuity of F , the PIT value $F(y)$
165 is replaced by $F_-(y) + U(F_+(y) - F_-(y))$, where $F_-(y)$ and $F_+(y)$ are respectively the left
166 and right limits of F at y . This modification allows in particular to consider the PIT for
167 ensemble forecast systems by interpreting the ensemble forecast as its empirical distribution
168 (resulting in a discontinuous distribution function with m jumps). The resulting PIT
169 histogram is then identical to the modified rank histogram suggested above.

170 As mentioned in the introduction, having with (2.1) a simple way of changing the
171 number of bins in a histogram is useful in its own right. Especially when rank histograms
172 are calculated on the same observations for competing forecast systems (with potentially
173 different ensemble sizes), it is useful to make them comparable by creating histograms with
174 the same bin number for both systems. Such a direct comparison can for example reveal
175 if one of the two models is substantially more biased or underdispersed than the other.
176 However, it is important to recognize that rank histograms are diagnostic tools and not
177 designed for model comparison. As pointed out by Hamill (2001), flatness of histograms
178 may result from mutual compensation between situations where the ensemble system is
179 not reliable, and observed flatness must be interpreted with caution.

180 **3 Tests for uniformity depending on the bin number**

181 In this section we review three tests for uniformity of the distribution of observation ranks,
182 and consider the number of bins as an additional parameter in the test. This will allow
183 us to adjust the bin number such that the test is approximated by a scientists' intuitive
184 decision. It should be stressed that considering the bin number as a parameter is not
185 useful from a data-analytic point of view: Reducing the number of bins by aggregating
186 multiple observation ranks into the same bin constitutes a loss of information that generally
187 reduces the power of a test for uniformity. Therefore, for assessing whether the observation
188 ranks are uniformly distributed, statistical tests such as the χ^2 -test should be applied to
189 the observation ranks directly, without aggregating them into fewer bins. Adjusting the
190 number of bins used in a rank histogram is mostly relevant when histograms are used for
191 intuitive inspection, i.e. as tools for visual diagnostics.

192 The three tests we consider are the classical χ^2 -test, a test based on the so-called
193 reliability index (Delle Monache et al., 2006), and a test considered by Taillardat et al.
194 (2016) based on an entropy statistic. We will refer to the latter two as RI-test and entropy

195 test, respectively. For their formal definition, as well as a comparison of their performance,
 196 we refer to Wilks (2019). The tests are conceptually similar in that the test statistic is
 197 always a distance between the observed histogram and a perfectly flat histogram. The
 198 hypothesis of uniformity is rejected when this distance exceeds a threshold value, which is
 199 determined by the significance level of the test. However, the tests differ in their definition
 200 of distance: the χ^2 -test is based on the L^2 -distance, the RI-test is based on the L^1 -distance,
 201 and the entropy test is based on the Kullback-Leibler-divergence.

In our context, it is convenient to rescale histograms such that their domain is the interval $[0, 1]$ and integrate to a total area of one. In particular, we interpret rank histograms as histograms for data points distributed in the interval $[0, 1]$, with the transformation (2.1) in mind. This simplifies notation greatly when considering different bin numbers for the same underlying data. We generally denote the number of bins by k and the height of the bins by h_1, \dots, h_k . Consequently, the frequency of the observation falling into the j th bin is h_j/k , and for a perfectly flat histogram we have $h_1 = \dots = h_k = 1$. For a histogram H_k with k bins we then consider the three test statistics, or distances,

$$D_{L^2} := \frac{1}{k} \sum_{i=1}^k (h_i - 1)^2, \quad D_{L^1} := \frac{1}{k} \sum_{i=1}^k |h_i - 1|, \quad \text{and} \quad D_{KL} := \frac{1}{k} \sum_{i=1}^k h_i \log(h_i),$$

202 where for D_{KL} we follow the convention that $0 \log(0) = 0$. The first two are the L^2 -
 203 and L^1 -distance between H_k and a flat histogram, respectively. The third statistic is the
 204 Kullback-Leibler divergence from $P(H_k)$ to U , where $P(H_k)$ is the probability distribution
 205 defined by the bin frequencies of H_k , and U is the uniform distribution.

For each of these distances, a statistical test is obtained for the null hypothesis that the underlying data is uniformly distributed. That is, the null hypothesis is rejected if

$$D(H_k) > c(\alpha, k, n), \tag{3.1}$$

206 where α is the significance level of the test. The threshold $c(\alpha, k, n)$ is defined as the smallest
 207 value c satisfying $P[D(H_k) > c] \leq \alpha$, when H_k is a histogram (with k bins) generated from

208 n independent uniformly distributed random variables. If we choose $D = D_{L^2}$, we recover
 209 the classical χ^2 -test, for $D = D_{L^1}$ we obtain the RI-test from Delle Monache et al. (2006),
 210 and, for $D = D_{KL}$, the entropy test from Taillardat et al. (2016).

211 We now aim to choose the bin number k such that a scientist’s intuitive decision ap-
 212 proximates such a formal statistical test. To this end we make the following assumption,
 213 for all three distances, i.e. $D \in \{D_{L^2}, D_{L^1}, D_{KL}\}$:

214 (A) There is an ‘acceptance threshold’ c_{acc} such that the scientist’s intuitive decision is
 215 well-approximated by rejecting whenever $D(H_k) > c_{\text{acc}}$. The acceptance threshold
 216 may depend on the chosen distance D .

217 Note that this assumption can be satisfied to different degrees for the different distances.
 218 It is, for example, possible that the use of an acceptance threshold constitutes a decent
 219 approximation to human behavior for $D = D_{L^2}$, but not for $D = D_{L^1}$. To what extent this
 220 assumption is satisfied by the different distances is assessed in the next section, where we
 221 also use the results of an empirical study to derive reasonable values for c_{acc} .

222 Subject to Assumption (A) being satisfied for one of the three distances D_0 , we can
 223 choose the bin number such that the scientist’s intuitive decision approximates the formal
 224 test based on D_0 . To this end, we choose a bin number k such that $c_{\text{acc}} \approx c(\alpha, k, n)$ from
 225 equation (3.1). Then, by Assumption (A), the scientist’s decision is close to the statistical
 226 test. The derived bin number then depends on the number of available observations n
 227 and on the significance level α of the test that is approximated. For a fixed number of
 228 observations n , the threshold $c(\alpha, k, n)$ is generally increasing in k and decreasing in α ,
 229 see Section 5. Consequently, if α is chosen small, k needs to be chosen small as well in
 230 order to achieve $c_{\text{acc}} \approx c(\alpha, k, n)$. This is intuitive, since for a small significance level only
 231 a small probability of a false reject is allowed. Reducing the bin number generally leads
 232 to flatter histograms if the underlying data is uniformly distributed, and therefore reduces
 233 the chance of an intuitive false reject by the scientist.

234 To sum up, in our proposed framework the optimal bin number k_{opt} is the one that
 235 minimizes $|c(\alpha, k, n) - c_{acc}|$. It depends on the number of available observations n , the
 236 selected significance level α , and the acceptance threshold c_{acc} . Such an optimal bin number
 237 can be derived for each of the three distances D_{L^1}, D_{L^2} and D_{KL} . Subject to Assumption
 238 (A), selecting this number of bins ensures that scientists' intuitive decisions are as close as
 239 possible to the statistical test associated with the corresponding distance.

240 4 The acceptance threshold

241 In this section we present the results of an empirical study assessing the validity of As-
 242 sumption (A) for the three different distances and derive approximations of the acceptance
 243 threshold. In this study several statisticians labeled histograms according to whether they
 244 believe them to be generated from uniform data or not. The histograms were in fact not
 245 based on underlying data at all, but were designed to have varying distances from uni-
 246 formity. Further details of the study design are given in the appendix. More than 15
 247 statisticians participated and 432 histograms were labeled.

248 For $D \in \{D_{L^2}, D_{L^1}, D_{KL}\}$ we consider the binary classifier

$$C_c(D(H_k)) = \begin{cases} \text{accept if } D(H_k) \leq c, \\ \text{reject if } D(H_k) > c \end{cases}$$

249 and compare the decision of this classifier to the intuitive decisions made by the statisticians.
 250 For a range of different c , we compute the misclassification rate of C_c , i.e. the proportion
 251 of cases where C_c decided differently than the statistician. The value c minimizing the
 252 misclassification rate then constitutes a good choice for c_{acc} , and the misclassification rate
 253 at this value provides a measure for the extent to which Assumption (A) is satisfied. The
 254 results for all three distances are shown in Figure 2. The lowest overall misclassification
 255 rate of 0.2 is achieved for $D = D_{L^2}$ and $c = 0.1$. In other words, rejecting a histogram

	c_{acc}	mcr	c_-	mcr	c_+	mcr
D_{L^2}	0.1	0.20	0.05	0.25	0.2	0.24
D_{L^1}	0.25	0.24	0.15	0.31	0.35	0.30
D_{KL}	0.05	0.21	0.02	0.27	0.09	0.26

Table 1: The three different values c_{acc} , c_- and c_+ considered as acceptance thresholds in Section 5, and their corresponding misclassification rates. The value c_{acc} is chosen to minimize the misclassification rate (mcr).

256 whenever its L^2 -distance exceeded 0.1 led to the same decision as the intuitive labeling
 257 for 4 out of 5 histograms. For D_{KL} a similarly small misclassification rate was achieved,
 258 whereas the misclassification rate for D_{L^1} was slightly higher, see Table 1 for details.

259 Different scientists have different preferences, and a histogram considered uniform by
 260 an optimist might be rejected by a pessimist. For the analysis in our next section we
 261 will therefore consider three different acceptance thresholds. The threshold minimizing the
 262 misclassification rate c_{acc} , which provides the best fit to the results of our empirical study,
 263 as well as thresholds c_- and c_+ , representing a pessimist and an optimist, respectively. For
 264 all three distances, c_- and c_+ were chosen such that the misclassification rate of C_c with
 265 respect to our study results was approximately 5% higher than for c_{acc} . The acceptance
 266 thresholds for the different distances and their corresponding misclassification rates are
 267 given in Table 1.

268 In practice, the decision of an expert to accept or reject can depend on an interplay
 269 between a distance from uniformity and the number of bins k . For example, an L^1 -distance
 270 of 0.25 for a histogram with 2 bins may be perceived as uniform, while the same distance
 271 of a histogram with 10 bins may be perceived as unacceptable. Such effects are unwanted
 272 in our context, since they are not accounted for by Assumption (A). In order to control for
 273 this effect, the 432 histograms labeled in the study had different bin numbers, namely 5,6,8,

274 or 10 bins. Figure 3 shows the acceptance rate of the scientists as a function of $D(H_k)$,
 275 for all three distances, and for each bin number k separately. The figures suggest that,
 276 at the same distance from uniformity, histograms with fewer bins tend to have a slightly
 277 higher acceptance rate. This is also supported by the correlation between bin number and
 278 scientist's decision, which was -0.16 if acceptance by the scientist got assigned the value 1
 279 and rejection got assigned the value 0. This effect is particularly clear for large values of
 280 D_{L^2} and D_{KL} and for 5 bins. An explanation for this could be that both D_{L^2} and D_{KL} put
 281 a higher penalty on outlier-bins than D_{L^1} , which could indicate that the labeling scientists
 282 found outlier-bins more likely to occur when few bins were used. Overall, however, the
 283 effect of the bin number on the decision is small compared to the effect of the distance.

284 5 Results

285 Here we present optimal bin numbers for a range of significance levels α and sample sizes
 286 n . As argued in the introduction, the results are mostly relevant for small data sizes n ,
 287 and we restrict our analysis to $n \leq 200$. We compute the optimal bin number for all three
 288 distances and the acceptance thresholds c_- , c_{acc} and c_+ given in Table 1. For α we consider
 289 the classical choice of 5%, as well as the more relaxed choices $\alpha = 10\%$ and $\alpha = 33\%$. While
 290 in most scenarios a statistical test with a false rejection probability of 33% is rather useless,
 291 such a threshold is not unreasonable in our informal setting where the test is approximated
 292 by scientists' intuitive decisions.

293 For given values of n, α, c and any of the distances D_{L^2}, D_{L^1}, D_{KL} , the optimal number k
 294 is then derived as follows. For all k in the range $k = 2, \dots, 12$ we compute $c(\alpha, k, n)$ from (3.1)
 295 and choose k such that $|c(\alpha, k, n) - c|$ is minimized. For the derivation of $c(\alpha, k, n)$ we do not
 296 rely on closed-form formulas (as in the original formulations of the tests), but use Monte-
 297 Carlo approximation with $N = 1.000.000$ samples. To be precise, we generate histograms
 298 H_1, \dots, H_N with k bins, each of which is based on n independent uniformly distributed data

299 points on $[0, 1]$. For each histogram we compute $D(H_k)$ and obtain $c(\alpha, k, n)$ as the minimal
300 value such that the fraction of histograms with $D(H_k) > c(\alpha, k, n)$ is smaller or equal to α .

301 The results are presented in Figure 4. It is clear to see that the bin number tends
302 to increase in the sample size n which is intuitive, since larger values of n reduce the
303 sample variability and therefore allow for separating the data into more bins. This effect
304 is, nevertheless, remarkable since it is not obvious from the way the optimal bin number is
305 derived. Indeed, the occasional dips of the red curves in Figure 4 show that the increasing
306 behavior in n constitutes a tendency rather than a mathematical necessity. The increasing
307 behavior can be explained by properties of the three distances used in the derivation. When
308 the underlying data is uniformly distributed, the distance from uniformity of a histogram
309 with fixed bin number k tends to decrease when the number of data points n increases. On
310 the other hand, the distance from uniformity tends to increase if the number of bins k is
311 increased for a fixed sample size n . While this behavior is not directly shown in the figure,
312 it implies that larger sample size n is balanced by larger k , in order to keep the probability
313 that the distance from uniformity exceeds the acceptance threshold at approximately α ,
314 and therefore that the optimal bin number tends to increase in n .

315 The results differ strongly between the different acceptance thresholds c_- , c_{acc} and
316 c_+ , highlighting that the optimal bin number depends substantially on the preferences of
317 the inspecting scientist. We will focus on the results for c_{acc} , which provides the best
318 approximation to our empirical study. Moreover, the study suggests that D_{L^2} and D_{KL}
319 are better suited to approximate human behavior than D_{L^1} , which suggests to focus on the
320 results for these two distances. Furthermore, Wilks (2019) concludes from his comparative
321 analysis of the three tests that *‘the traditional χ^2 test is recommended as a consequence
322 of its generally superior power, particularly for the underdispersed ensembles that are most
323 commonly encountered, and the relative ease of obtaining the necessary critical values.’*
324 This suggests putting most emphasis on the bin numbers derived by using the L^2 -distance.
325 There is remarkable similarity between the optimal bin numbers for D_{L^2} and D_{KL} when

326 $c = c_{\text{acc}}$, which provides a sanity check for our approach: Even though the derivation of
327 the optimal bin number is based on different test statistics for different distances, the goal
328 remains the same. Namely, to find a bin number that leads to an intuitive rejection of
329 histograms of uniform data with probability α .

330 As we would expect, the bin number k increases not only in n but also in c and α .
331 The increase in α highlights that, if one is willing to accept large probabilities of a false
332 reject, one should consider rank histograms with many bins, since this also tends to increase
333 the probability of a correct reject (the power of the associated test) when the data is not
334 uniformly distributed. The variability in c mainly provides insight to what extent the
335 results depend on the personal preferences of the scientist, but it should be mentioned that
336 the selection of c_- and c_+ in Section 4 is rather arbitrary.

337 Overall, the bin numbers suggested by this approach are relatively small, especially
338 for small sample sizes n . For $n = 100$, our approach suggests to choose only 5 bins in
339 order to approximate a conservative test with significance level of 5% (focusing on c_{acc}
340 and either D_{L^2} or D_{KL}). If we relax the significance level to 10% (33%), the algorithm
341 selects 6 bins (9 bins) instead. In particular, if we have 100 forecast-observation pairs
342 available, and we choose to print a histogram with 9 bins, we need to expect a roughly
343 33% chance for an intuitive false reject if the ensemble forecast system is well-calibrated.
344 If only 50 observations are available, the bin numbers drop to 2 (5%), 3 (10%) and 5
345 (33%), respectively. Such bin numbers constitute a stark contrast to the common practice
346 of choosing $m + 1$ bins which typically results in 11 bins or more.

347 Instead of focusing on the theoretically optimal number of bins, we may analyze the
348 false rejection rate of the classifier C_c as a function of the bin number k . Figure 5 shows
349 the results for the bin numbers $k = 4, 6, 8$ and 10. Again, we observe that the differences
350 between the distances D_{L^2} , D_{L^1} and D_{KL} are small. Especially for the pessimistic threshold
351 c_- the false rejection probabilities are very large, even for small number of bins. This can be
352 interpreted as a warning not to be too pessimistic when visually inspecting rank histograms

353 based on few observations, but rather acknowledge that the natural variability is likely to
354 result in histograms that may not look approximately flat, even when the underlying data
355 is uniformly distributed.

356 **6 Rejection probabilities under non-uniform distribu-** 357 **tions**

358 In this section we analyze the rejection probability of the considered tests under non-
359 uniform distributions. We consider two distributions representing the most prominent
360 characteristic shapes that are important in rank histogram analysis. The first distribution
361 is sloped, with a density linearly increasing from $2/3$ at 0 to $4/3$ at 1 , representing rank
362 histograms based on a biased prediction system. The second distribution is U-shaped
363 representing rank histograms based on an underdispersed prediction system. The U-shaped
364 distribution has density $f(x) = 3(x - 1/2)^2 + 3/4$, which is symmetric around $1/2$ where
365 it reaches its minimum value of $3/4$. Figure 6 shows histograms of the two distributions
366 based on 200.000 samples.

367 We obtain rejection probabilities for the three distributions by generating, for a range
368 of n and k , 1000 histograms with k bins based on n data points with the corresponding
369 distribution, and computing the distances D_{L^1} , D_{L^2} and D_{KL} for these histograms. The
370 rejection probability for one of these distances and a given acceptance threshold c is then
371 the fraction of histograms for which the distance exceeds c . As acceptance thresholds we
372 consider the three values c_- , c_{acc} and c_+ specified in Table 1. Figure 7 shows the rejection
373 probabilities for these acceptance thresholds under the three distributions, for a range of
374 bin numbers and sample sizes. The figure only shows the results for the L^2 -distance, the
375 other distances lead to very similar results (not shown). Generally, the rejection probability
376 increases in the bin number, showing that histograms based on more bins tend to have a

377 higher distance from uniformity under all three considered distributions. The uniform
378 distribution gets rejected with the lowest probability, which indicates that the considered
379 tests are unbiased. However, when $k = 2$, the U-shaped histogram gets rejected with the
380 same probability. This highlights that histograms based on two bins are essentially useless
381 in practice, since they cannot indicate misspecified dispersion in the ensemble forecast
382 system.

383 The figure clearly visualizes the trade-off that is made in choosing the number of bins:
384 While a low rejection probability is desirable when the data is uniformly distributed, high
385 rejection probabilities are desirable for the two alternative distributions. Figure 7 shows
386 that using c_+ generally leads to very low rejection probabilities, even for non-uniform
387 data. The pessimistic threshold c_- , on the other hand, generally leads to much lower
388 rejection probabilities for uniformly distributed data than for data generated from the
389 alternative distributions. However, the probability for a false reject is generally very large
390 when c_- is used, for example it is more than 75% when 12 bins are chosen, even for
391 $n = 180$. The threshold c_{acc} suggested by our empirical study leads to a large difference
392 in acceptance probabilities between uniform and non-uniform distributions and, at the
393 same time, allows for reasonably small false rejection probabilities. It can generally be
394 observed that the differences in rejection probabilities between uniform and non-uniform
395 distribution are getting more clearly pronounced as n increases. This highlights the fact
396 that with more available data it becomes easier to differentiate between uniform and non-
397 uniform distributions. It is also worth mentioning that the optimist's acceptance threshold
398 c_+ performs reasonable well for $n < 100$. Consequently, for very small n , one should be
399 careful not to expect too uniform histograms.

400 Figure 8 shows the rejection probability for the three distributions when the optimal
401 bin number is used. Here, the optimal bin number is derived using the L^2 -distance, the
402 acceptance threshold c_{acc} and the significance level $\alpha = 5\%$. The significance level is shown
403 in the figure as dashed line. The plot in the middle shows that the bin number is selected in

404 order to align the blue line with the 5% significance level. Note that approximately $n = 40$
405 is required in order to achieve a false rejection rate of only 5%, even when only two bins
406 are used. The left hand side and right hand side plot show the rejection probabilities for
407 pessimist and optimist, respectively, when they inspect histograms based on the optimal
408 number of bins derived with the acceptance threshold c_{acc} .

409 7 Discussion

410 Our study indicates that, when visually inspecting forecast calibration with rank his-
411 tograms, choosing a small number of bins can substantially lower the risk of wrongfully
412 rejecting the hypothesis that the underlying data is uniform.

413 In practice, rank histograms are applied to identify characteristic shapes indicating
414 certain miscalibrations of the ensemble forecast. This has several implications. The most
415 common characteristic shapes in the appearance of rank histograms are slopes (indicating
416 bias) as well as \cup - and \cap -shapes (indicating under- and overdispersion, respectively). In
417 particular, it is never advisable to only use two bins (as our approach suggests in some
418 cases for very small sample sizes), since such a histogram is unable to pick up on dispersion
419 misspecification. At the same time, these simple shapes are equally well captured by a
420 histogram with three bins than by histograms with many bins. More involved characteristic
421 shapes (e.g. S-shapes) can indicate misspecified skewness or combinations of bias and
422 misspecified dispersion. However, they often require a large sample size n to become clearly
423 visible, see Thorarinsdottir and Schuhen (2018). Such shapes are generally captured by
424 histograms with six or eight bins, and it is difficult to imagine any informative characteristic
425 shape that would require more than 10 bins in order to become visible. On the contrary, our
426 results indicate that increasing the bin number puts more emphasis on random fluctuations
427 in the data which can distract from characteristic shapes. Based on these considerations
428 we recommend to generally limit the number of bins in histograms to about 10. When the

429 number of available forecast-observation pairs is limited one should not hesitate to consider
430 histograms with fewer bins. Histograms with three bins might look somewhat unusual, but
431 may be more appropriate when n is very small in order to mitigate effects of sampling
432 uncertainty.

433 At the same time, choosing a very small number of bins increases the risk of not rec-
434 ognizing deviations from uniformity, as shown in Section 6. Moreover, in situations where
435 the size of the verification data set is not known to the inspector, a larger number of bins
436 can help the inspector to estimate how many forecast-observation-pairs were used and thus
437 to avoid false acceptance or rejection of uniformity.

438 We assumed throughout this paper that the ranks of the different forecast-observation-
439 pairs are independent. This assumption is commonly made when rank histograms are
440 constructed, but is violated in some applications, in particular when multiple spatial grid
441 points are considered as samples. Such complex dependence structure can make the his-
442 togram much harder to interpret and, in particular, prevent formal testing for uniformity.
443 See Hamill (2001) for an in-depth discussion of this topic.

444 8 Conclusion

445 We introduce a criterion for choosing the number of bins in a rank histogram. The crite-
446 rion attempts to make the intuitive decision of scientists regarding calibration close to a
447 statistical test. It addresses the trade-off that adding more bins leads to a more detailed
448 histogram but at the same time decreases statistical robustness, and attempts to optimize
449 intuitive decision making based on the histogram. Our results highlight that the probability
450 for intuitively rejecting a histogram tends to increase with the number of bins, even if the
451 underlying data is uniformly distributed. This generally questions the current practice of
452 choosing as many bins as possible. We showed that reducing the bin number can, to some
453 extent, be used to appropriately balance the probability of an intuitive false reject, which

454 also depends on the sample size n . This probability further depends on the preferences
455 and experience level of the inspecting scientist. The bin numbers derived in the previous
456 section are therefore merely suggestions based on our empirical study and do not constitute
457 theoretical optima that ought to be followed under all circumstances.

458 Our results indicate that, especially for small verification samples with less than 100
459 data points, histograms with five bins or fewer are preferable. If histograms with more bins
460 are considered, their appearance should not be over-interpreted, and rather large deviations
461 from flatness should be expected, even for histograms based on uniformly distributed data.
462 Moreover, for very small sample sizes of 50 or less, the probability for an intuitive false
463 reject is generally rather large (often 50% or higher), for any reasonable bin number ($k > 2$).
464 This highlights the large uncertainty associated with such small sample sizes and shows
465 that rank histograms should in such situations be interpreted very carefully. Generally, and
466 particularly in this case, rank histogram analysis should rely on the results of statistical tests
467 for uniformity rather than on intuitive inspection of the histogram plot. The importance of
468 this is highlighted by our study that showed that intuitive decisions are strongly dependent
469 on the selected number of bins, which is a property of the histogram plot only, not of the
470 distribution of observation ranks in the predictive ensemble.

471 This article is accompanied by the R-package `RankHistBins` which is available on the
472 authors github account github.com/ClaudioHeinrich/RankHistBins. The package in-
473 cludes functionality to generate histograms with any bin number from observed ranks
474 using the transformation (2.1), and to compute the optimal bin number for any sample
475 size n , acceptance threshold c and test size $1 - \alpha$. Moreover, it provides tools and guidance
476 that allow the reader to conduct the empirical study described in Section 4. By person-
477 ally labeling histograms you can derive your personal acceptance threshold c_{acc} , and derive
478 optimal bin numbers for histograms inspected by yourself.

Appendix: Details on the Empirical Study

Here we give more details about the design of the empirical study presented in Section 4. An early version of this paper only considered the L^1 -distance from uniformity. Therefore, the study originally focused on analyzing the effect of different L^1 -distances only. The analysis of L^2 -distance and Kullback-Leibler divergence was added later and not taken into account for study design. For the study, 1000 histograms were created with 5,6,8 or 10 bins, and with L^1 -distance in $\{0.1, 0.15, \dots, 0.45, 0.5, 0.6\}$. The histograms were not based on underlying data, but were sampled by an algorithm described below that allows to generate histograms with pre-specified number of bins and L^1 -distance. Considering 4 different bin numbers and 10 different L^1 -distances resulted in 40 categories, for each of which 25 histograms were created. The created histograms were shuffled, printed out and laid out in the break room of the statistics and data science group of the Norwegian Computing Center in Oslo, Norway, with a call to the group to label as many histograms as possible. The participants labeled the histograms according to whether they believe them to be based on uniform data or not, and were left unaware that the histograms were not based on underlying data at all. The labeling of histograms was anonymous and participants could label as many histograms as they wanted. More than 15 Statisticians confirmed that they participated, and 432 out of the 1000 printed histograms were labeled. In all 40 categories the number of labeled histograms was between 7 and 16 (out of 25), except for one category where only three histograms were labeled. A detailed key of how many histograms were labeled in which category is shown in Figure 9.

The following algorithm was used for creating a random histogram with pre-specified bin number k and L^1 -distance from uniformity D .

1. Choose a number of steps n (in the study $n = 50$) for the algorithm. Start out with a perfectly uniform histogram with k bins. Mark all the bins with a 0.
2. Randomly select one of the bins marked 0 or 1, and *increase* its height by $\frac{Dk}{2n}$. If the

505 bin was marked 0, change its mark to 1.

506 3. Randomly select one of the bins marked 0 or -1, and *decrease* its height by $\frac{Dk}{2n}$. If the
507 bin was marked 0, change its mark to -1.

508 4. Repeat steps 2 and 3 in total n times.

509 In this algorithm, both steps 2 and 3 increase the L^1 -distance from uniformity by $D/2n$, and,
510 since they are both repeated n times, the final histogram has L^1 -distance from uniformity D .
511 The marking is important to ensure that bins that have been increased (decreased) in height
512 will only ever be increased (decreased), which ensures that the distance in fact increases in
513 each step. The alternation between increasing and decreasing bin heights ensures that the
514 total integral of the histogram remains 1. The algorithm needs an additionally constraint
515 that prevents that bin heights are decreased beyond zero.

516 **References**

517 Anderson, J. L. (1996). A method for producing and evaluating probabilistic forecasts from
518 ensemble model integrations. *J. Climate* 9, 1518–1530.

519 Delle Monache, L., J. P. Hacker, Y. Zhou, X. Deng, and R. B. Stull (2006). Probabilistic aspects of meteorological and ozone regional ensemble forecasts. *J. Geophys. Res.-Atmos.* 111, D24307.
520
521

522 Hamill, T. (2001). Interpretation of rank histograms for verifying ensemble forecasts. *Mon. Weather Rev.* 129(3), 550–560.
523

524 Hamill, T. and S. Colucci (1997). Verification of Eta-RSM Short-Range Ensemble Forecasts. *Mon. Weather Rev.* 125(6), 1312–1327.
525

- 526 He, K. and G. Meeden (1997). Selecting the number of bins in a histogram: A decision
527 theoretic approach. *J. Stat. Plan. Infer.* 61(1), 49–59.
- 528 Knuth, K. (2019). Optimal data-based binning for histograms and histogram-based prob-
529 ability density models. *Digital Signal Processing* 95, 102581.
- 530 Muto, K., H. Sakamoto, K. Matsuura, T. Arima, and M. Okada (2019). Multidimensional
531 bin-width optimization for histogram and its application to four-dimensional neutron
532 inelastic scattering data. *J. Phys. Soc. Jpn.* 88(4), 044002.
- 533 Scott, D. (1979). On optimal and data-based histograms. *Biometrika* 66(3), 605–610.
- 534 Smith, J. (1985). Diagnostic checks of non-standard time series models. *Journal of Fore-*
535 *casting* 4(3), 283–291.
- 536 Sturges, H. (1926). The choice of a class interval. *J. Am. Stat. Assoc.* 21(153), 65–66.
- 537 Taillardat, M., O. Mestre, M. Zamo, and P. Naveau (2016). Calibrated ensemble forecasts
538 using quantile regression forests and ensemble model output statistics. *Mon. Weather*
539 *Rev.* 144(6), 2375–2393.
- 540 Talagrand, O., R. Vautard, and B. Strauss (1997). Evaluation of probabilistic prediction
541 systems. In *Proceedings, Workshop on Predictability, European Centre for Medium-Range*
542 *Weather Forecasts*, pp. 1–25.
- 543 Thorarinsdottir, T. L., M. Scheuerer, and C. Heinz (2016). Assessing the calibration of
544 high-dimensional ensemble forecasts using rank histograms. *J. Comp. Graph. Stat.* 25(1),
545 105–122.
- 546 Thorarinsdottir, T. L. and N. Schuhen (2018). Verification: Assessment of calibration
547 and accuracy. In S. Vannitsem, D. S. Wilks, and J. W. Messner (Eds.), *Statistical*
548 *Postprocessing of Ensemble Forecasts*, Chapter 6, pp. 155–186. Elsevier.

- 549 Van Schaeybroeck, B. and S. Vannitsem (2018). Postprocessing of long-range forecasts. In
550 *Statistical Postprocessing of Ensemble Forecasts*, pp. 267–290. Elsevier.
- 551 Wand, M. (1997). Data-based choice of histogram bin width. *The American Statisti-*
552 *cian* 51(1), 59–64.
- 553 Wilks, D. S. (2004). The minimum spanning tree histogram as verification tool for multi-
554 dimensional ensemble forecasts. *Mon. Weather Rev.* 132, 1329–1340.
- 555 Wilks, D. S. (2019). Indices of rank histogram flatness and their sampling properties. *Mon.*
556 *Weather Rev.* 147(2), 763–769.
- 557 Ziegel, J. F. and T. Gneiting (2014). Copula calibration. *Electron. J. Stat.* 8(2), 2619–2638.

558 **List of Figures**

559 1 Three histograms based on the same data with different number of bins. The
560 data is a sample of 30 numbers, uniformly and independently distributed
561 on $[0,1]$. The middle plot shows an example for a distance from uniformity
562 considered in this paper: The size of the hatched area is the L^1 -distance D_{L^1}
563 between the considered histogram and a perfectly flat histogram. Clearly,
564 the distance varies with the number of bins. 27

565 2 The misclassification rate of the binary classifier C_c for the three distances
566 D_{L^2} , D_{L^1} and D_{KL} , as a function of the acceptance threshold c . The low
567 values all three curves attain at their minimum indicate that the classifier
568 C_c is a decent approximation for a scientist's intuitive decision, with D_{L^2}
569 and D_{KL} providing slightly better approximations than D_{L^1} . The misclas-
570 sification rate of D_{L^1} is a step function due to the design of the empirical
571 study: In a first version of this paper only the L^1 -distance was considered,
572 and the participants were therefore presented histograms that were gener-
573 ated to have a predefined L^1 -distance, namely $\{0.1,0.15,\dots\}$. The distances
574 D_{L^2} and D_{KL} of the labeled histograms were computed later on. 28

575 3 The acceptance rate of the statisticians as a function of the distance, sep-
576 arately for different bin numbers k . For D_{L^1} and D_{L^2} the histograms are
577 aggregated over intervals of length 0.1. As an example, the value shown at
578 $D = 0.2$ is the acceptance rate over all histograms with a distance in the
579 interval $(0.1, 0.2]$. For D_{KL} the same aggregation is applied over intervals of
580 length 0.05. 29

581 4 The optimal bin number as a function of the data size n , for three different
582 significance levels α , and the three choices of acceptance threshold c , specified
583 in Table 1. 30

584 5 The probability of a false rejection as a function of the data size n , for $k =$
585 4, 6, 8 and 10 bins. Increasing the bin number leads to a higher probability
586 for a false reject, but at the same time increases the probability for a correct
587 reject if the underlying data is not uniformly distributed, cf. Figure 6. . . 31

588 6 Histograms of the two non-uniform distributions considered in Section 8. . 32

589 7 Rejection probabilities for a range of n and k for the uniform distribution,
590 and the sloped and the U-shaped distribution described in Section 8. The
591 results are only shown for the test based on the L^2 -distance, and for the
592 three acceptance thresholds c_- , c_{acc} , c_+ given in Table 1. 33

593	8	Rejection probabilities for a range of n when the optimal bin number is used.	34
594	9	How many histograms (out of 25 possible) were labeled in each category in	
595		the empirical study.	35

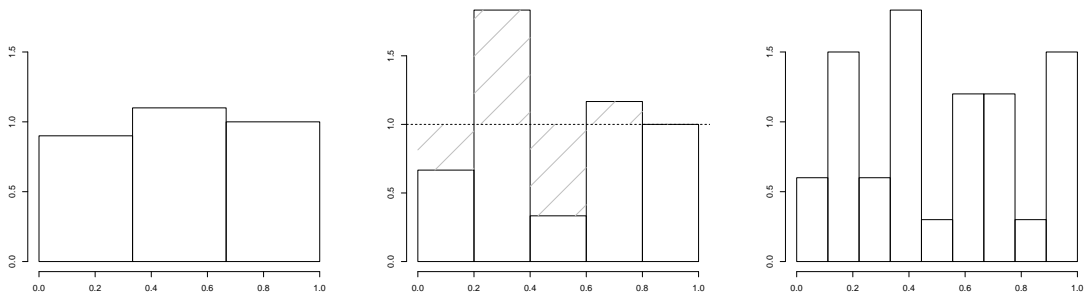


Figure 1: Three histograms based on the same data with different number of bins. The data is a sample of 30 numbers, uniformly and independently distributed on $[0,1]$. The middle plot shows an example for a distance from uniformity considered in this paper: The size of the hatched area is the L^1 -distance D_{L^1} between the considered histogram and a perfectly flat histogram. Clearly, the distance varies with the number of bins.

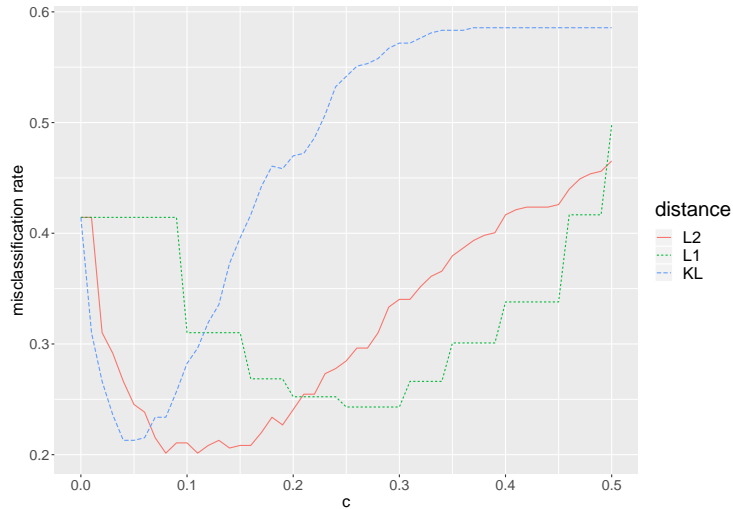


Figure 2: The misclassification rate of the binary classifier C_c for the three distances D_{L^2} , D_{L^1} and D_{KL} , as a function of the acceptance threshold c . The low values all three curves attain at their minimum indicate that the classifier C_c is a decent approximation for a scientist’s intuitive decision, with D_{L^2} and D_{KL} providing slightly better approximations than D_{L^1} . The misclassification rate of D_{L^1} is a step function due to the design of the empirical study: In a first version of this paper only the L^1 -distance was considered, and the participants were therefore presented histograms that were generated to have a predefined L^1 -distance, namely $\{0.1, 0.15, \dots\}$. The distances D_{L^2} and D_{KL} of the labeled histograms were computed later on.

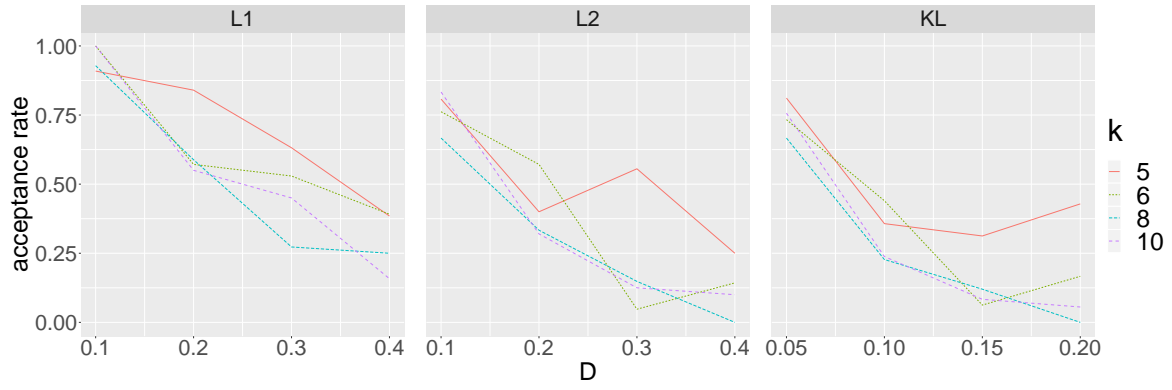


Figure 3: The acceptance rate of the statisticians as a function of the distance, separately for different bin numbers k . For D_{L^1} and D_{L^2} the histograms are aggregated over intervals of length 0.1. As an example, the value shown at $D = 0.2$ is the acceptance rate over all histograms with a distance in the interval $(0.1, 0.2]$. For D_{KL} the same aggregation is applied over intervals of length 0.05.

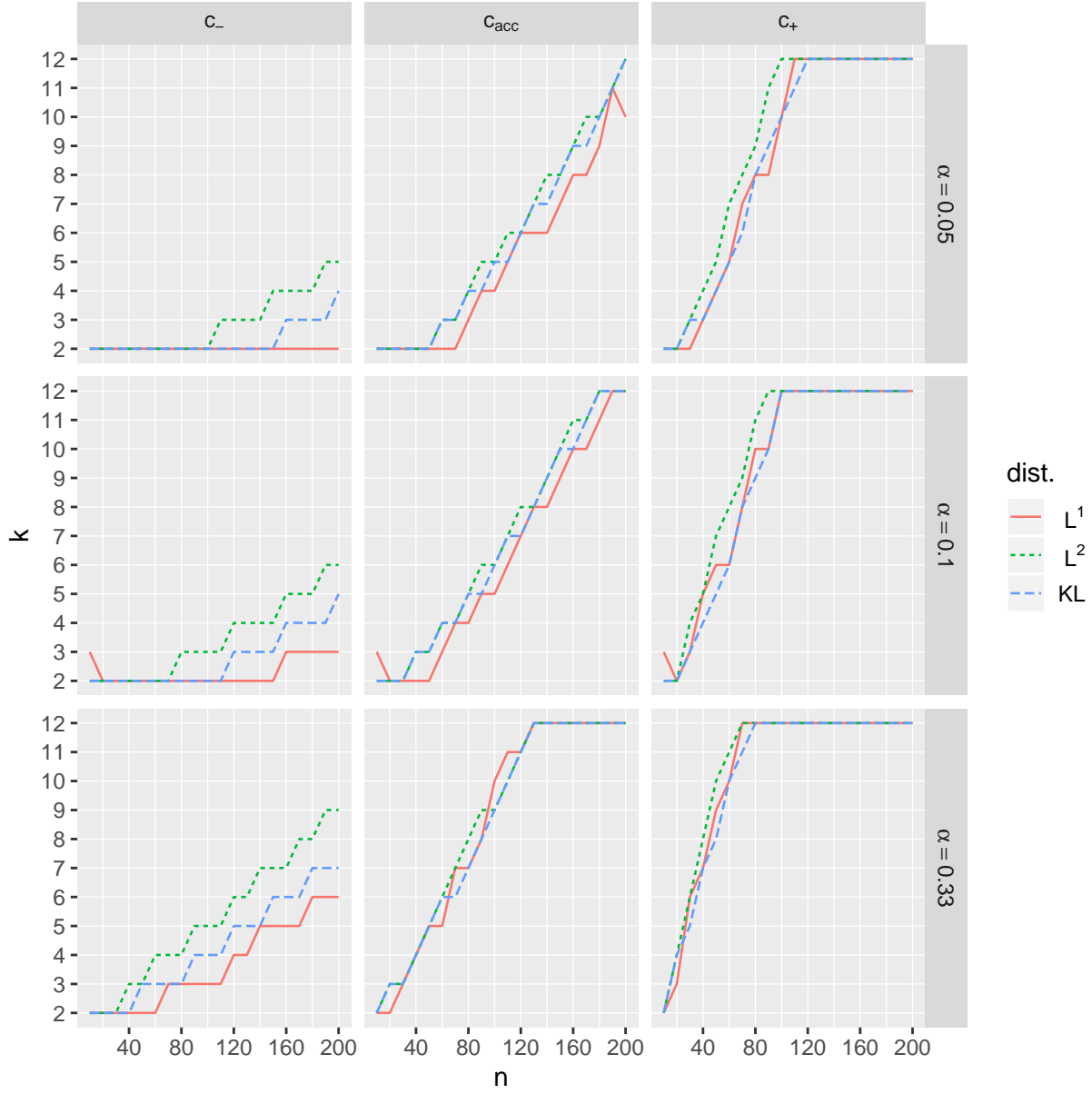


Figure 4: The optimal bin number as a function of the data size n , for three different significance levels α , and the three choices of acceptance threshold c , specified in Table 1.

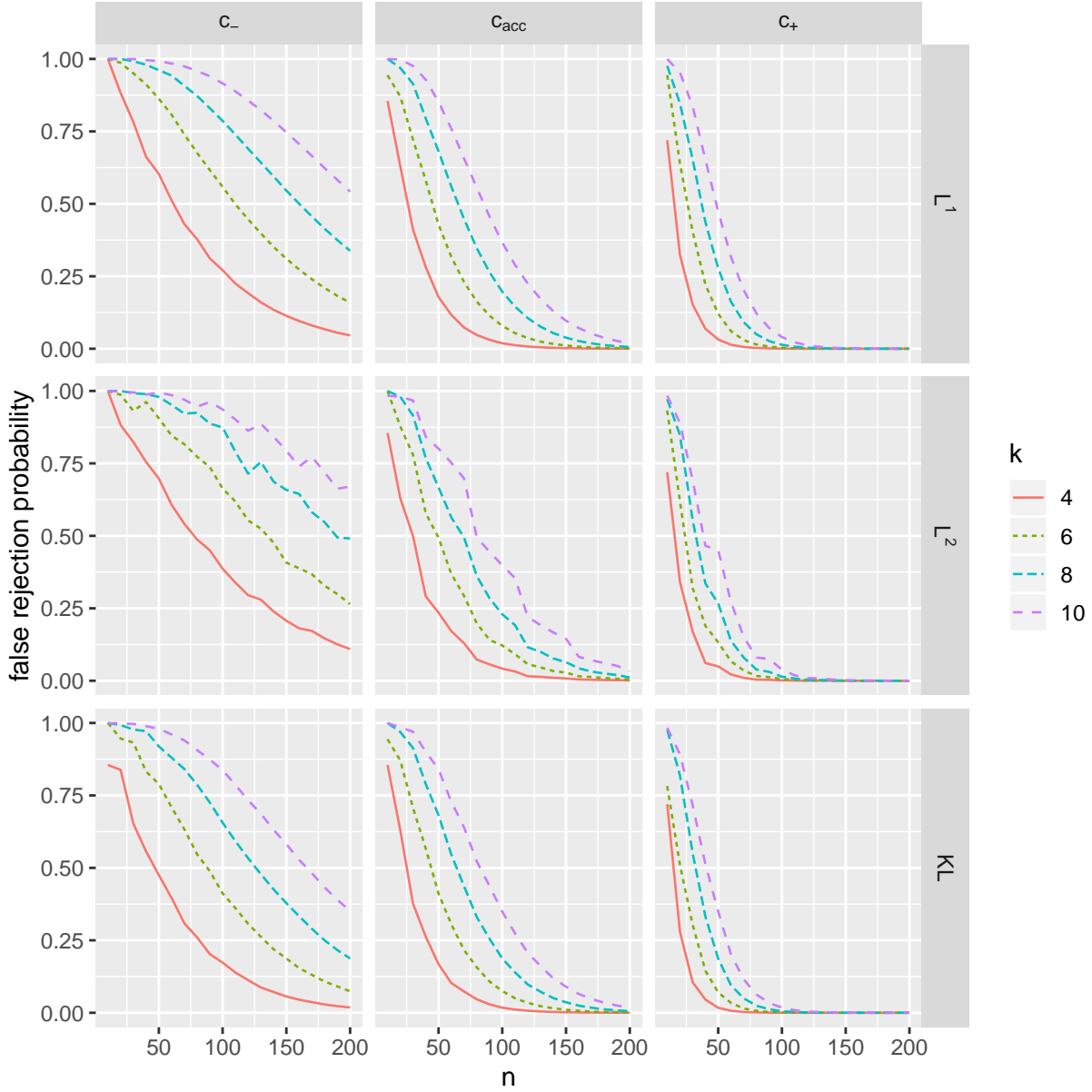


Figure 5: The probability of a false rejection as a function of the data size n , for $k = 4, 6, 8$ and 10 bins. Increasing the bin number leads to a higher probability for a false reject, but at the same time increases the probability for a correct reject if the underlying data is not uniformly distributed, cf. Figure 6.

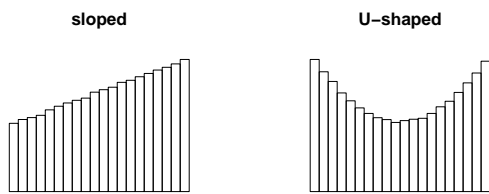


Figure 6: Histograms of the two non-uniform distributions considered in Section 6.

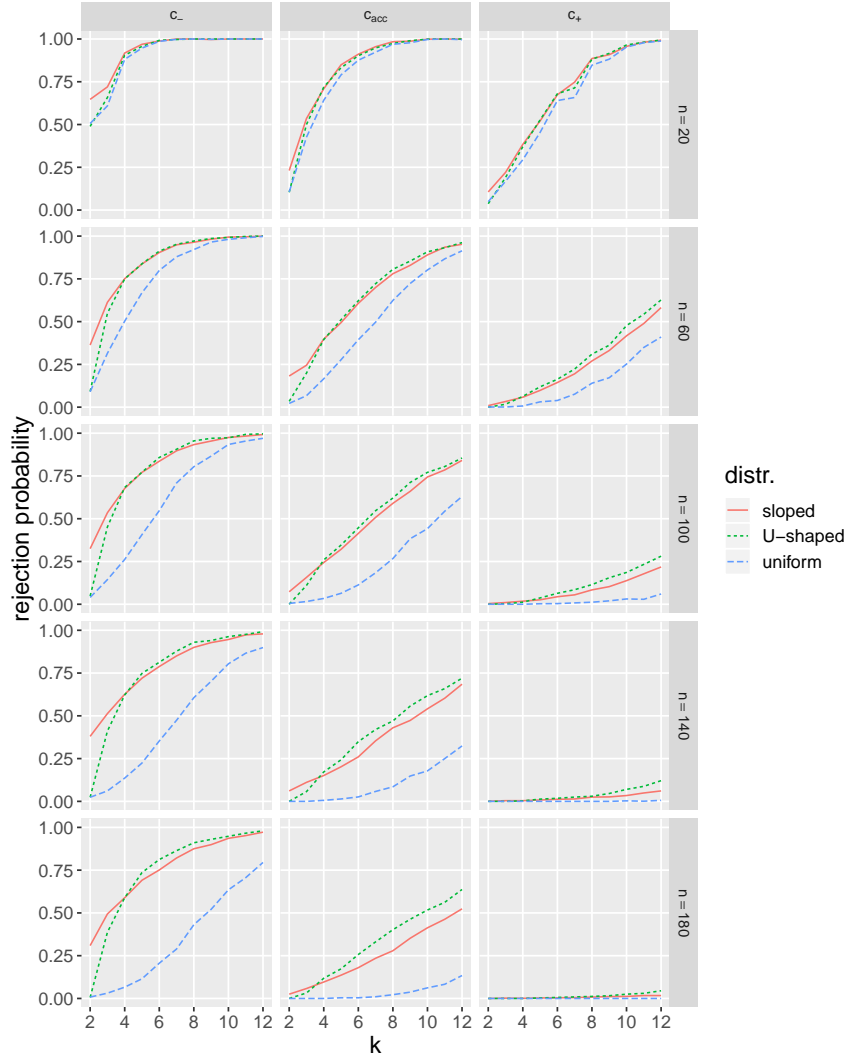


Figure 7: Rejection probabilities for a range of n and k for the uniform distribution, and the sloped and the U-shaped distribution described in Section 6. The results are only shown for the test based on the L^2 -distance, and for the three acceptance thresholds c_- , c_{acc} , c_+ given in Table 1.

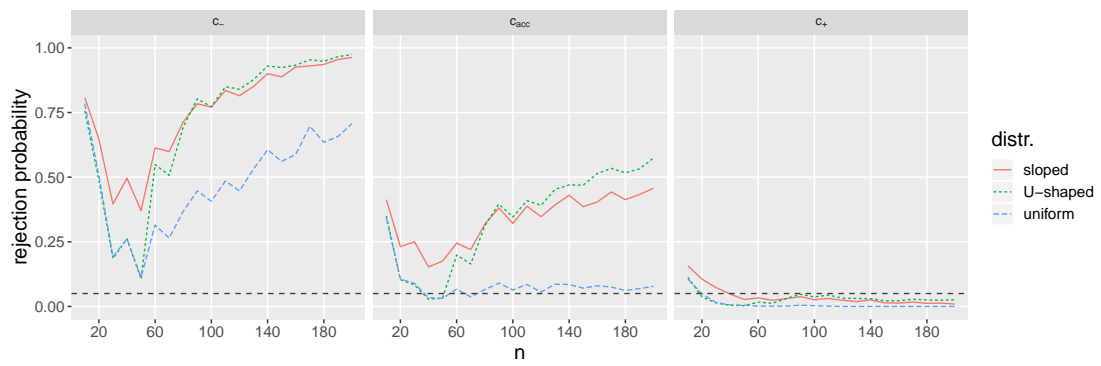


Figure 8: Rejection probabilities for a range of n when the optimal bin number is used.

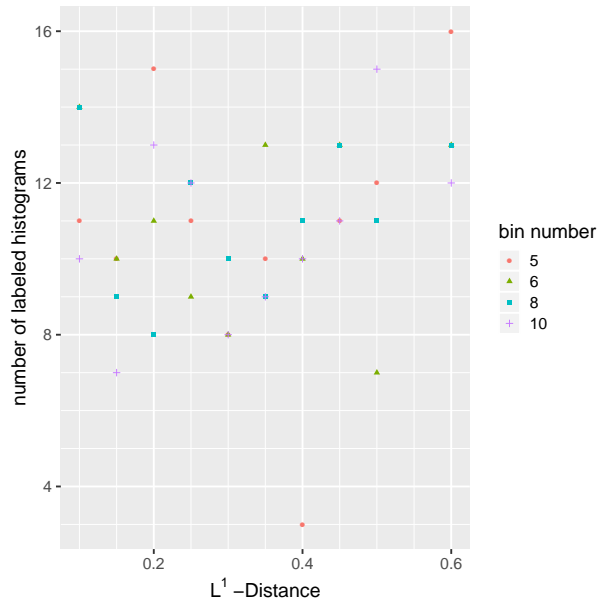


Figure 9: How many histograms (out of 25 possible) were labeled in each category in the empirical study.