# Histogram binning revisited with a focus on human perception

Raphael Sahann\* Faculty of Computer Science, University of Vienna, Austria Torsten Möller<sup>†</sup> Faculty of Computer Science, Data Science @ Uni Vienna, University of Vienna, Austria Johanna Schmidt<sup>‡</sup> VRVis Zentrum für Virtual Reality und Visualisierung Forschungs-GmbH, Austria



Figure 1: Selection of histogram datasets used in our study. We evaluated how well human viewers can detect the underlying data distribution in a histogram when different sample sizes and bins are used. For this, we created datasets with a different number of samples (first row: few, last row: many) and a different number of bins (left column: 2, right column: 100). A bimodal distribution was used to create the datasets in this illustration.

# ABSTRACT

This paper presents a quantitative user study to evaluate how well users can visually perceive the underlying data distribution from a histogram representation. We used different sample and bin sizes and four different distributions (uniform, normal, bimodal, and gamma). The study results confirm that, in general, more bins correlate with fewer errors by the viewers. However, upon a certain number of bins, the error rate cannot be improved by adding more bins. By comparing our study results with the outcomes of existing mathematical models for histogram binning (e.g., Sturges' formula, Scott's normal reference rule, the Rice Rule, or Freedman–Diaconis' choice), we can see that most of them overestimate the number of bins necessary to make the distribution visible to a human viewer.

Keywords: empirical studies in visualization; histogram binning

# **1** INTRODUCTION

Histograms are a well-known and prevalent visualization technique [19] representing the distribution of univariate data by visualizing the tabulated frequency at certain intervals, represented as bars or *bins*. Several bins next to each other help human viewers to build a mental model of the data distribution. The most important parameter visualization designers have to set when creating a histogram is the *number of bins*.

Statisticians have developed several thumb rules to help researchers estimate the right number of bins when creating a histogram. For example, Sturge's formula [25] defines how to split the data into k bins based on the number of samples being available. Scott's normal reference rule [23] measures the discrepancy between the bin representation and the data distribution by employing mean integrated squared error. The Freedman–Diaconis choice [10] is based on minimizing the difference between the area under the data distribution and the area under the probability distribution defined by the binning. The so-called Rice's rule [26] can be applied for nonnormal distributed data. More advanced approaches, like the ones by Lolla and Hoberock [16] and Birge and Rozenholc [1], use mathematical concepts like Cumulative Distribution Functions (CDF) or penalty functions to calculate the number of bins, reflecting the idea that smooth distributions need fewer bins than rough distributions.

The mathematical models, on the forefront being Sturge's formula, Scott's normal reference rule, the Rice's rule, and the Freedman–Diaconis choice, are quite popular and are often used in current visualization systems and libraries. The mathematical models are fast and convenient binning estimations. Interestingly, the models and their suggested binnings have not been evaluated in perceptual user studies yet. The models in use today were statistically and mathematically evaluated. However, it is unclear how well the suggested numbers of bins match the human visual perception. The number of bins suggested by the mathematical models may be quite high (> 300). In case histograms need to be shown on small displays, e.g., smartphones, it would be interesting to know whether such a high number of bins is really needed when shown to a human viewer.

User studies [14] offer a scientifically sound method to measure how people read visualizations [11], and a number of studies have already been undertaken in an effort to assess these aspects [12]. When it comes to summary statistics, Lem et al. [15] and Kaplan et al. [13] noticed a general problem for students when trying to read and interpret aggregated information in histograms and box plots. Correll et al. [6] highlighted the importance of selecting the right number of bins for detecting missing values and outliers in a histogram. In general, the literacy of humans interpreting these visualizations seems to be decoupled from the statistical interpretation. According to Boels et al. [3], information reduction still seems to be an understudied topic in data visualization. This is also in line with the current need in research to understand how viewers can construct and interpret data visualizations [4] and with the need for further research on visualization guidelines [7].

<sup>\*</sup>e-mail: raphael.sahann@univie.ac.at

<sup>&</sup>lt;sup>†</sup>e-mail: torsten.moeller@univie.ac.at

<sup>&</sup>lt;sup>‡</sup>e-mail: johanna.schmidt@vrvis.at

# 2 QUANTITATIVE USER STUDY

In this paper we report about a user study, also sometimes called user evaluation study [9], that has been conducted to compare the number of bins suggested by mathematical models with human perception when analyzing histograms. We address the research question whether the numbers of bins suggested by statistical computations match the minimum number of bins required for human viewers to be able to detect the data's underlying distribution in a histogram. For this, we used datasets with four different distributions (*uniform, normal, gamma*, and *bimodal*), different sample sizes, and numbers of bins (see Figure 1), and asked participants to state which distributions they see in the different representations.

#### 2.1 Hypotheses generation

**Task definition**: Histograms as summary statistics provide the possibility to perform several tasks related to distribution analysis (e.g., identifying the mean and the median or comparing quartiles). One task related to distribution analysis is to identify the data's underlying distribution, which has been classified as the task to "describe and identify the shape and type of one distribution" in the literature [2]. The identification of the underlying distribution is the task we evaluated in our study.

**Distributions**: To get an overview of the distributions currently used in practice, we looked into literature targeted towards data scientists to learn more about the use of data distributions. Some examples are: *Doing Data Science* [22] lists 17 density functions that data scientists should be familiar with. The *Data Scientist's Crib Sheet* [20] describes 15 density functions that are important and highlights their relationships. In the *KDnuggets* tutorials [24] five density functions are explained that data scientists should be aware of. Based on this literature research and based on our own experience when working with data, we decided to classify the available density distributions based on their main shape characteristics. We defined four main classes:

- uniform: uniform distributions
- *unimodal*: distributions with one peak, similar to a Gaussian kernel
- *bimodal*: distributions with two peaks
- *skewed*: distributions with one peak which are skewed to one side of the distribution

This classification is also confirmed by Walker [17], who describes the most common shapes of distributions as *bell shaped*, *left skewed*, *right skewed*, *bimodal*, and *uniform*. A quick pre-test among students in the study preparation phase did not show any differences between the detection of left and right skewness. We, therefore, only considered distributions which are skewed to the left in our study. After deciding on the four classes, we identified one mathematical density function representing each class best:

- For the class uniform, a uniform distribution fits best.
- For the class *unimodal*, we selected the normal density function to represent this class.
- For the class *bimodal* we joined two normal density functions to form a bimodal distribution with two peaks.
- For the class *skewed* we selected the gamma density function to represent this class.

**Number of samples and bins**: The number of samples and bins to be tested were chosen based on a mathematical analysis of the distributions. We calculated the four moments (mean, variance, skew, and kurtosis) for each of the distributions. The combination of these four moments can uniquely identify a distribution's shape. For five different sample sizes (100, 1000, 10000, 100000, and 1000000)



Figure 2: Calculated average errors of the four moments based on 100,000 draws from binned samples (n = 1,000,000) compared to the actual moments of the underlying distribution. The x-axis shows the number of bins, the y-axis shows the error in percent. Note that the y-axes are scaled differently.

we drew 1,000 times from the four distributions specified above and recorded the actual moments from each draw. We then created evenly spaced binnings in steps of 1 from 2 up to 100 bins for each draw. Using the bins' centers as the outline of a new shape, we calculated its moments and compared them to the actual moments. Figure 2 shows the calculated errors in the four moments when using different bin sizes. The error is large for small numbers of bins for all moments but reaches an almost constant rate above 10 bins. Adding more bins does not affect the error rate anymore. We, therefore, decided that the range up to 10 bins would be of the greatest interest for our study. Since some moments (especially kurtosis) take a little longer to settle completely, we chose to include some values between 10 and 40 as well. We also included 100 bins as an upper boundary since it is equal to our study's smallest sample size.

**Hypotheses generation**: We agreed upon testing the following hypotheses:

- **Hypothesis H1**: The number of bins influences how well humans can perceive the underlying data distribution in a histogram.
- **Hypothesis H2**: Upon a certain number of bins, adding new bins does not improve the perception of underlying data distributions in a histogram.

### 2.2 Study design

To summarize, we tested

- four distributions (uniform, normal, bimodal, and gamma), with
- ten different bin counts (2,3,4,5,7,10,15,20,40, and 100), and
- four sample sizes (100, 1000, 10000, and 1000000).

Question 9 of 23



Figure 3: Study question. Participants were shown a histogram depicting the data's underlying distribution and asked to click the appropriate icon. Participants were also asked to state how confident they are about their answers. In this example *normal* and *confident* have been selected.

We used the approach of a web-based questionnaire to be able to reach a large group of participants [21]. A Cross-Site Request Forgery (CSRF) token was generated whenever a participant decided to start the survey. Since we then used only this token to identify the participant, the study was fully anonymized without any possibility to track the results back to the participants.

For every participant, we started with an initial explanation what the study will be about. Afterwards we included a sanity check to filter out careless participants [18]. Afterwards the actual study questions started, where we showed 20 histograms to every participant, one after the other. The histograms were randomly selected from the pool of datasets. The histogram plots had two axes with ticks, but we did not show any numbers or scales. Participants were asked to answer the question

"1. Choose the distribution which resembles the image above most closely"

by clicking on one of the icons below the histogram showing different possible distributions. Participants were also asked to state

"2. How confident are you about your answer?"

on a four-point Likert scale. Participants could only proceed with the next histogram if they answered both questions. An example of how the web-based implementation of such a histogram question looked like is shown in Figure 3. The study setting can best be described as a *judgement study*, where the study's purpose is to gather a person's response to a set of stimuli [5]. According to the literature, judgement studies are a commonly used approach for perceptual studies.

# 3 RESULTS

In total, 82 participants finished the user study within a 14 day time frame. We only counted complete submissions and did not record the dropouts. Based on the sanity check questions, we had to exclude 10 data records from the evaluation, which led to a final number of 72 valid submissions.

The majority of the participants were between 20 and 49 years old. One-third of the participants (33%) were bachelor's, master's, or PhD students. The other participants were working part-time (14%), full-time (23%), or, more specifically, in research and education (26%). 27% of the participants had some experience in reading charts and plots in the media. 27% classified themselves as being experienced in reading data visualizations, and 41% stated that they are also creating data visualizations themselves. Only four participants stated that they do not have any experience with data visualization.

#### 3.1 Study results

The evaluation of the quantitative results led to the following results:

Insight 1: Small sample sizes generally make it harder to detect the underlying data distribution, which can only slightly be mitigated by using a higher number of bins.

For datasets with 100 samples, 35.4% of the answers were wrong. With 1000, 10000, and 1000000 samples being available, the detection error rate could be halved to 16.1%, 18.3%, and 18.8%. A Mann–Whitney U test [8] resulted in p-values p < 0.001 when comparing the results for all sample sizes, which confirms the statistical significance of the results. Participants stated to be less confident when judging the distribution with a sample size of 100. The amount of participants being *very confident* about their answers constantly increases with a rate of about 10% for larger samples sizes.

# Insight 2: Beyond a certain number of bins the error rate stays constant and is not improved by adding more bins.

Like the sample size, the number of bins affected participants' ability to recognize the underlying distribution correctly. More bins result in fewer errors being made by the participants (see Figure 4). This effect is different, depending on how many samples are available. In the case of 100 samples, the error rate stays rather high, also in cases where a higher number of bins was used. For other sample sizes, the error rate decreases in case a larger number of bins is used. However, for larger sample sizes, it can be seen that more bins do not improve the visual perception of the underlying data distribution. While the error rate is significantly better when comparing the bin size 2 to other parameters (p < 0.001), the difference between a larger number of bins is not significant any more (*bins/bins: p-value* -15/20: p = 0.072, 20/40: p = 0.442, 40/100: p = 0.121).



Figure 4: Correct (*blue*) and wrong (*orange*) answers based on the number of bins and sample size. With only 100 samples, the data distribution recognition is generally challenging. For larger sample sizes, a larger number of bins increases the recognition of the correct distribution. However, beyond 20 bins, the detection rate does not increase significantly anymore.



Figure 5: Recorded Answers by distribution and bin size. The chart title shows the actual shown distribution, the colored segments show the user responses. Users could mostly identify gamma and uniform distributions with two bins, needed three bins to distinguish a normal distribution, and they needed at least ten bins to recognize a bimodal distribution.

		~ .			
		Scott's		Freedman-	
	Sturge's	normal	Rice's	Diaconis	our
samples	formula	reference	rule	choice	results
100	8	9	14	18	20
1,000	11	20	29	38	20
10,000	15	43	62	80	20
1,000,000	21	200	287	371	20

Table 1: Suggested bin sizes from different mathematical models based on the number of samples in the data set. Apart from Sturge's formula almost all models overestimate the number of bins needed for 1,000 samples and above.

#### Insight 3: For bimodal distributions the number of bins is more important to recognize them correctly.

The percentages of correct and wrong answers for every distribution indicate that it is generally easier to detect gamma and uniform distributions (see Figure 4). The detection of bimodal distributions was harder and very strongly affected by the number of bins. In the case of normal distributions, apart from 2 bins, the detection worked quite well. Bimodal distributions are often confused with gamma distributions in case of a low number of bins (see Figure 5).

### Insight 4: Experience in reading data visualization had no impact on the error rate. More experience led to higher confidence when answering the questions.

Only minor, non-significant differences could be detected when analyzing the percentage of correct and wrong answers compared to the participants' stated experience with visualizations. Participants with no or mediocre experience were generally less confident when answering the questions than those who had extensive data visualization knowledge.

# Insight 5: Most mathematical models overestimate the number of bins needed for characterizing distributions.

When comparing the suggested bin count of different mathematical models shown in Table 1 to the error rates from our study (as shown in Figure 4) it is clear that, apart from Sturge's formula, all models have a tendency suggest too many bins starting from 1,000 samples. Our results show that much fewer bins are needed for detecting distributions.

### 3.2 Hypotheses testing

We can summarize the results of the study based on our two hypotheses in the following way:

- **Hypothesis H1**: We can *partially confirm* that the visual perception of the underlying data distribution depends on the number of bins. At least with larger sample sizes (enough samples to resemble the underlying distribution), the recognition becomes better when using more bins.
- **Hypothesis H2**: We can *confirm* that upon a certain number of bins, adding new bins does not improve the perception of the underlying data distribution. After 20 bins, the error rate cannot be decreased significantly by adding additional bins.

#### **4** IMPACT AND DISCUSSION

The task that was tested in the study was on "*describing and identifying the shape and type of one distribution*" [2]. We, therefore, would like to emphasize, that all results are solely to be interpreted for this specific task.

We also asked ourselves how the study results for this task can probably also be explained by sampling theory. Binning can be seen as a way to sample the original distribution. As known from sampling theory, it is impossible to reconstruct the original function with too few samples. We, therefore, transferred the four density distributions we used (uniform, normal, bimodal, and gamma) into the frequency domain. None of these is band-limited, hence, there is no concrete Nyquist frequency. As a baseline we, therefore, used a representation of each density distribution with 1000000 uniform samples. We then compared this baseline to representations of the same density distribution with less (5, 10, 15, 20, 40, and 100) samples (= bins). For a comparison we measured the deviation (i.e., error) between the baseline and the binned representation in the frequency domain. The representation with only five samples (= bins) stands out to have the most significant error. From ten samples (= bins) on, the error starts to converge to zero. 100 samples (=bins) already ensure an error to the original representation very close to zero. It is important to note that the difference between five and ten samples (= bins) is much more significant than the difference between 40 and 100. The comparison of the representations in the frequency domain are, therefore, very similar to our study results.

A retrospective analysis of the study setting showed that, due to the fact the questions were randomly selected, an almost equal distribution of answers per question could be achieved (*bins: percentage of answers* -2: 9.02%, 3: 9.63%, 4: 10.79%, 5: 8.72%, 7: 10.49%, 10: 9.94%, 15: 9.57%, 20: 9.33%, 40: 11.89%, 100: 10.61%). As a drawback, participants mainly already had experience with data visualization. In the future, we would, therefore, like to access a broader range of users with a new study, and we would also like to test other tasks related to histograms.

# 5 CONCLUSION

We presented a quantitative evaluation and comparison of mathematically defined numbers of bins for histogram with human perception. The mathematical models (e.g., Scott's normal reference rule, the Rice Rule, Freedman-Diaconis' choice) mostly overestimate the number of bins necessary for a correct perception for human viewers. With around 20 bins, the error rate for human viewers to detect the data's underlying distribution becomes stable and does not improve by adding more bins.

#### ACKNOWLEDGMENTS

We thank the anonymous reviewers for their helpful comments and suggestions which helped to improve the paper. VRVis is funded by BMK, BMDW, Styria, SFG, Tyrol and Vienna Business Agency in the scope of COMET - Competence Centers for Excellent Technologies (879730) which is managed by FFG.

### REFERENCES

- L. Birge and Y. Rozenholc. How many bins should be put in a regular histogram. *ESAIM: Probability and Statistics*, 10:24–45, 2006. doi: 10. 1051/ps:2006001
- [2] M. Blumenschein, L. J. Debbeler, N. C. Lages, B. Renner, D. A. Keim, and M. El-Assady. v-plots: Designing Hybrid Charts for the Comparative Analysis of Data Distributions. *Computer Graphics Forum*, 39(3):565–577, 2020. doi: 10.1111/cgf.14002
- [3] L. Boels, A. Bakker, W. Van Dooren, and P. Drijvers. Conceptual difficulties when interpreting histograms: A review. *Educational Research Review*, 28:100291, 2019. doi: 10.1016/j.edurev.2019.100291
- [4] K. Börner, A. Bueckle, and M. Ginda. Data visualization literacy: Definitions, conceptual frameworks, exercises, and assessments. *Proceedings of the National Academy of Sciences*, 116(6):1857–1864, 2019. doi: 10.1073/pnas.1807180116
- [5] S. Carpendale. Evaluating Information Visualizations, pp. 19–45. Springer Berlin Heidelberg, 2008. doi: 10.1007/978-3-540-70956-5\_2
- [6] M. Correll, M. Li, G. Kindlmann, and C. Scheidegger. Looks Good To Me: Visualizations As Sanity Checks. *IEEE Transactions on Visualization and Computer Graphics*, 25(1):830–839, 2019. doi: 10. 1109/TVCG.2018.2864907
- [7] A. Diehl, A. Abdul-Rahman, M. El-Assady, B. Bach, D. Keim, and M. Chen. Visguides: A Forum for Discussing Visualization Guidelines. In *Proceedings of the Eurographics/IEEE VGTC Conference on Visualization: Short Papers*, EuroVis '18, pp. 61–65. Eurographics Association, Brno, Czech Republic, June 4–8 2018. doi: 10.2312/ eurovisshort.20181079
- [8] A. Field and G. Hole. *How to Design and Report Experiments*. SAGE Publications, 2003.
- [9] C. Forsell. A Guide to Scientific Evaluation in Information Visualization. In *Proceedings of the 14th International Conference Information Visualisation*, IV '10, pp. 162–169. London, UK, July 26–29 2010. doi: 10.1109/IV.2010.33
- [10] D. Freedman and P. Diaconis. On the histogram as a density estimator: L2 theory. Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete, 57(4):453–476, 1981. doi: 10.1007/BF01025868
- [11] W. Huang. Handbook of Human Centric Visualization. Springer New York, 2014.
- [12] T. Isenberg, P. Isenberg, J. Chen, M. Sedlmair, and T. Möller. A Systematic Review on the Practice of Evaluating Visualization. *IEEE Transactions on Visualization and Computer Graphics*, 19(12):2818–2827, 2013. doi: 10.1109/TVCG.2013.126
- [13] J. J. Kaplan, J. G. Gabrosek, P. Curtiss, and C. Malone. Investigating Student Understanding of Histograms. *Journal of Statistics Education*, 22(2), 2014. doi: 10.1080/10691898.2014.11889701
- [14] R. Kosara, C. G. Healey, V. Interrante, D. H. Laidlaw, and C. Ware. User Studies: Why, How, and When? *IEEE Computer Graphics and Applications*, 23(4):20–25, 2003. doi: 10.1109/MCG.2003.1210860
- [15] S. Lem, P. Onghena, L. Verschaffel, and W. Van Dooren. On the misinterpretation of histograms and box plots. *Educational Psychology*, 33(2):155–174, 2013. doi: 10.1080/01443410.2012.674006
- [16] S. V. G. Lolla and L. L. Hoberock. On Selecting the Number of Bins for a Histogram. In *Proceedings of the International Conference on Data Mining*, DMIN '11, 2011.
- [17] MathBootCamps. Common shapes of distributions. https://www.mathbootcamps.com/common-shapes-of-distributions/, 2017. last accessed [2020-11-29].
- [18] A. S. M. Niessen, R. R. Meijer, and J. N. Tendeiro. Detecting careless respondents in web-based questionnaires: Which method to use? *Journal of Research in Personality*, 63:1–11, 2016. doi: 10.1016/j.jrp.2016. 04.010
- [19] R. L. Nuzzo. Histograms: A Useful Data Analysis Visualization. *PM&R*, 11(3):309–312, 2019. doi: 10.1002/pmrj.12145
- [20] S. Owen. Common Probability Distributions: The Data Scientist's Crib Sheet. https://medium.com/@srowen/common-probabilitydistributions-347e6b945ce4, 2018. last accessed [2021-03-15].
- [21] U.-D. Reips. The methodology of Internet-based experiments, pp. 373–390. Oxford University Press, 2007. doi: 10.1093/oxfordhb/ 9780199561803.013.0024

- [22] R. Schutt and C. O'Neil. *Doing Data Science*. O'Reilly Media, Inc., 2013.
- [23] D. W. Scott. On optimal and data-based histograms. *Biometrika*, 66(3):605–610, 1979. doi: 10.1093/biomet/66.3.605
- [24] L. Strika. 5 Probability Distributions Every Data Scientist Should Know (KDnuggets, Tutorials, Overviews). https://www.kdnuggets.com/2019/07/5-probability-distributionsevery-data-scientist-should-know.html, 2019. last accessed [2021-03-17].
- [25] H. A. Sturges. The Choice of a Class Interval. *Journal of the American Statistical Association*, 21(153):65–66, 1926. doi: 10.1080/01621459. 1926.10502161
- [26] G. R. Terrell and D. W. Scott. Oversmoothed Nonparametric Density Estimates. *Journal of the American Statistical Association*, 80(389):209–214, 1985. doi: 10.2307/2288074