

ASYMPTOTICALLY OPTIMAL CELLS FOR A HISTOGRAM

BY ATSUYUKI KOGURE

Fukushima University

The purpose of this paper is to examine the properties of the histogram when the cells are allowed to be arbitrary. Given a random sample from an unknown probability density f on I , we wish to construct a histogram. Any partition of I can be used as cells. The optimal partition minimizes the mean integrated squared error (MISE) of the histogram from f . An expression is found for the infimum of MISE over all partitions. It is proved that the infimum is attained asymptotically by minimizing MISE over a class of partitions of locally equisized cells.

1. Introduction. Let f be a continuous probability density on an interval I in the real line. I includes its endpoints if it is finite. Given a random sample X_1, X_2, \dots, X_n of size n from f , we wish to construct a histogram. Prior to the actual drawing of a histogram we need to choose cells into which the X_i 's are grouped. Any partition of I can be used as cells. (We call a collection $\{C_j\}$ of intervals of I a partition of I if $C_i \cap C_j = \emptyset$ whenever $i \neq j$ and if $I = \bigcup_j C_j$.) When a partition $Q = \{C_j\}$ is chosen, the histogram on C_j is

$$H_n(x|Q) := F_n(C_j)/|C_j|,$$

where $|C_j|$ denotes the width of cell C_j and

$$F_n(C_j) := n^{-1} \sum_{i=1}^n \{X_i \in C_j\}.$$

Considering the histogram an estimate of f , we adopt the mean integrated squared error,

$$(1.1) \quad \text{MISE}(n, Q) := E \left[\int_I (H_n(x|Q) - f(x))^2 dx \right],$$

as the risk is using Q . Here E denotes the expectation with respect to the X_i 's. The risk is decomposed into two components as

$$(1.2) \quad \text{MISE}(n, Q) := \frac{1}{n} \sum_j \frac{F(C_j)(1 - F(C_j))}{|C_j|} + \int_I (f_Q(x) - f(x))^2 dx,$$

where $f_Q(x) := E[H_n(x|Q)]$. The first component on the right of (1.2) represents the sampling variability and the second the bias. Shrinking cell sizes decreases the bias at the cost of increasing the sampling variability. The optimal partition is to minimize MISE by compromising the conflict between the two components.

Received April 1986; revised January 1987.

AMS 1980 subject classifications. Primary 62G05; secondary 62E20.

Key words and phrases. Histogram, cells, partition, mean integrated squared error.

In this paper we examine the properties of the histogram when the cells are allowed to be arbitrary. We ask how close, in terms of the MISE, the histogram comes to f . Several authors—including Tapia and Thompson (1978), Scott (1979) and Freedman and Diaconis (1981)—have investigated this problem. However, they deal with the subject under the restrictive setting that all cells in a partition have a common cell width, i.e., for all j , $|C_j| = h$ for some $h > 0$. We pursue this problem with no restrictions on cells. We derive an asymptotic expression for the infimum of MISE over all partitions and show that the infimum can be achieved by a “partition of locally equisized cells.”

We analyze the problem under the following conditions on f :

- (C.1) f is twice continuously differentiable;
- (C.2) f , f' and f'' all belong to L_2 ;
- (C.3) $\int |f'(x)f(x)|^{2/3} dx > 0$.

An important class of densities excluded by the conditions is that of step functions. For any such density, the optimal histogram is the density itself.

2. Asymptotic behaviour of MISE. How small can MISE be made as the sample size tends to infinity? Let \mathcal{P} denote the class of all partitions of I . We have

THEOREM 1. *Let (C.1) and (C.2) be satisfied. Then*

$$(2.1) \quad \inf_{Q \in \mathcal{P}} \text{MISE}(n, Q) = \left[\left(\frac{3}{2}\right)\left(\frac{1}{6}\right)^{1/3} \int_I |f'(x)f(x)|^{2/3} dx \right] n^{-2/3} + o(n^{-2/3}).$$

Under the restrictive setting that all cells in a partition have a common width h , Freedman and Diaconis (1981) obtained a relation parallel to (2.1):

$$(2.2) \quad \inf_{h>0} \text{MISE}(n, h) = \left[\left(\frac{3}{2}\right)\left(\frac{1}{6}\right)^{1/3} \left\{ \int_I (f'(x))^2 dx \right\}^{1/3} \right] n^{-2/3} + o(n^{-2/3}),$$

where $\text{MISE}(n, h)$ is the MISE of the histogram with common cell width h . See also Scott (1979). The ratio of the leading term of (2.1) to that of (2.2) is $\int_I |f'(x)f(x)|^{2/3} dx / \left\{ \int_I (f'(x))^2 dx \right\}^{1/3}$, which is seen to be no more than 1 by Hölder's inequality. It is approximately 0.89 for the normal density and 0.24 for the Cauchy density.

Theorem 1 will be derived by proving the following two propositions. The first of these gives an asymptotic lower bound for the MISE.

PROPOSITION 1. *Under (C.1) and (C.2) we have*

$$(2.3) \quad \liminf_{n \rightarrow \infty} n^{2/3} \inf_{Q \in \mathcal{P}} \{\text{MISE}(n, Q)\} \geq \left(\frac{3}{2}\right) 6^{-1/3} \int_I |f'(x)f(x)|^{2/3} dx.$$

All we need then is to find a partition whose MISE behaves like the lower bound as n increases. Our scheme for getting such a partition is based on three

observations:

- (i) on a small interval f is well approximated by a linear function;
- (ii) the bias term $\int (f_Q(x) - f(x))^2 dx$ of $\text{MISE}(n, Q)$ is minimized for a partition of equisized cells if f is linear; and
- (iii) if f does not change much over the region, then the magnitude of the sampling variation $(1/n)\sum_j F(C_j)(1 - F(C_j))/|C_j|$ is mostly determined by the number of cells, regardless of the cell sizes being equal or unequal.

With these observations we might expect that the desired partition can be found in a class of partitions of locally equisized cells (POLEC). A POLEC is constructed by "dividing I twice," a process separated into two steps:

Step 1. Choose a reference point $x_0 \in I$ and a mesh width $w > 0$ to obtain an equally spaced net $\{\Delta_i | i \in Z\}$, where Z is a set of integers and

$$(2.4) \quad \begin{aligned} \Delta_i &:= [x_0 + (i - 1)w, x_0 + iw), \\ \bigcup_{i \in Z} \Delta_i &= I. \end{aligned}$$

Step 2. For each $i \in Z$ choose a positive integer k_i and divide Δ_i into k_i cells, allowing different k_i 's for different Δ_i 's. Put $\mathbf{k} := \{k_i | i \in Z\}$ and let $Q(w, \mathbf{k})$ denote the resultant POLEC.

Let $\{w_n\}$ and $\{U_n\}$ be sequences of positive numbers such that

$$(2.5) \quad n^{-1/3} \ll w_n \ll n^{-2/9}, \quad U_n \ll n^{1/3}, \quad (w_n/U_n) \ll n^{-1/3},$$

where $a_n \ll b_n$ means $\limsup_{n \rightarrow \infty} (a_n/b_n) = 0$. Define a class of POLECs as

$$(2.6) \quad \mathcal{P}(w_n, U_n) := \{Q(w_n, \mathbf{k}) | k_i \leq U_n \text{ for all } i \in Z\}.$$

Then we have

PROPOSITION 2. *Under conditions (C.1) and (C.2) we have*

$$(2.7) \quad \begin{aligned} &\inf\{\text{MISE}(n, Q) : Q \in \mathcal{P}(w_n, U_n)\} \\ &= \left(\frac{3}{2}\right)6^{-1/3} \int_I |f'(x)f(x)|^{2/3} dx n^{-2/3} + o(n^{-2/3}). \end{aligned}$$

Clearly, Propositions 1 and 2 together imply Theorem 1. The proofs of the propositions occupy the next section.

Another implication of the two propositions is

$$\inf\{\text{MISE}(n, Q) : Q \in \mathcal{P}(w_n, U_n)\} = \inf_{Q \in \mathcal{P}} \{\text{MISE}(n, Q)\} + o(n^{-2/3}).$$

This means that there is a POLEC whose associated MISE is nearly as small as that of the optimal partition as the sample size grows large. This may call forth some interest in developing a data-driven rule for choosing a POLEC. Some discussion on this can be found in Kogure (1986). For the related subject of the data-based choice of a cell width, see Scott (1979), Rudemo (1982), Chow, Geman and Wu (1983) and Stone (1985).

3. Proofs.

PROOF OF PROPOSITION 1. For each interval C let

$$\gamma(C) := \left(\frac{3}{2}\right)6^{-1/3} \int_C |f'(x)f(x)|^{2/3} dx$$

and for each $Q = \{C_j\} \in \mathcal{P}$ let $\alpha_n(Q) := n^{2/3} \text{MISE}(n, Q) - \gamma(I)$.

Put $h_n := (n^{2/9} \log n)^{-1}$ and define a subclass \mathcal{P}_n of \mathcal{P} as

$$\mathcal{P}_n := \left\{ Q = \{C_j\} \mid 0 < |C_j| \leq h_n \text{ for all } j \right\}.$$

For each interval C , $0 < |C| < \infty$, let

$$\delta_n(C) := \left(\frac{1}{12}\right)(n^{1/3}|C|)^2 + \int_C (f'(x))^2 dx + (n^{1/3}|C|)^{-1} F(C)$$

and for each $Q = \{C_j\} \in \mathcal{P}_n$ let $\beta_n(Q) := \sum_j (\delta_n(C_j) - \gamma(C_j))$. Observe that

$$(3.1) \quad \begin{aligned} \inf_{Q \in \mathcal{P}} \alpha_n(Q) &\geq \inf_{Q \in \mathcal{P}_n} \beta_n(Q) + \inf_{Q \in \mathcal{P}_n} (\alpha_n(Q) - \beta_n(Q)) \\ &+ \left(\inf_{Q \in \mathcal{P}} \alpha_n(Q) - \inf_{Q \in \mathcal{P}_n} \alpha_n(Q) \right). \end{aligned}$$

Then it suffices to show that each term on the right of (3.1) has a nonnegative lower limit. First, we will show

$$(3.2) \quad \liminf_{n \rightarrow \infty} \inf_{Q \in \mathcal{P}_n} \beta_n(Q) \geq 0.$$

Fix $Q = \{C_j\} \in \mathcal{P}_n$. For each $C_j \in Q$

$$\begin{aligned} \delta_n(C_j) &\geq \left(\frac{3}{2}\right)6^{-1/3} \left\{ \int_{C_j} (f'(x))^2 dx \right\}^{1/3} (F(C_j))^{2/3} \\ &\geq \left(\frac{3}{2}\right)6^{-1/3} \int_{C_j} |f'(x)f(x)|^{2/3} dx \quad (\text{by Hölder's inequality}) \\ &= \gamma(C_j). \end{aligned}$$

Thus, $\beta_n(Q) \geq 0$, which implies (3.2).

Second, we will show

$$(3.3) \quad \liminf_{n \rightarrow \infty} \inf_{Q \in \mathcal{P}_n} (\alpha_n(Q) - \beta_n(Q)) \geq 0.$$

Fix $Q = \{C_j\} \in \mathcal{P}_n$. Then

$$(3.4) \quad \begin{aligned} \alpha_n(Q) - \beta_n(Q) &= n^{2/3} \text{MISE}(n, Q) - \sum_j \delta_n(C_j) \\ &= n^{2/3} \left(\int_I (f_Q(x) - f(x))^2 dx - \left(\frac{1}{12}\right) \sum_j |C_j|^2 \int_{C_j} (f'(x))^2 dx \right) \\ &\quad - n^{1/3} \sum_j (F(C_j))^2 / |C_j|. \end{aligned}$$

It is easy to see that for each j there is a point x_j inside c_j such that

$$\int_{C_j} (f_Q(x) - f(x))^2 dx \geq \left(\frac{1}{12}\right) (f'(x_j))^2 |C_j|^3.$$

Then, by Lemma 2.22 of Freedman and Diaconis (1981), we have

$$\left| (f'(x_j))^2 |C_j| - \int_{C_j} (f'(x))^2 dx \right| \leq 2|C_j| \int_{C_j} |f'(x)f''(x)| dx.$$

Thus,

$$(3.5) \quad \int_I (f_Q(x) - f(x))^2 dx - \left(\frac{1}{12}\right) \sum_j |C_j|^2 \int_{C_j} (f'(x))^2 dx \\ \geq -2h_n^3 \int_I |f'(x)f''(x)| dx.$$

Note that

$$(3.6) \quad \sum_j (F(C_j))^2 / |C_j| = \int_I (f_Q(x))^2 dx \leq \int_I (f(x))^2 dx.$$

Combining (3.4), (3.5) and (3.6) and recalling that $h_n = (n^{2/9} \log n)^{-1}$ we have (3.3).

Lastly, we will show

$$(3.7) \quad \liminf_{n \rightarrow \infty} \left\{ \inf_{Q \in \mathcal{P}} \alpha_n(Q) - \inf_{Q \in \mathcal{P}_n} \alpha_n(Q) \right\} \geq 0.$$

It would suffice to show that for each $Q \in \mathcal{P}$, there is $Q^0 \in \mathcal{P}_n$ such that for all sufficiently large n

$$(3.8) \quad \alpha_n(Q) - \alpha_n(Q^0) \geq r_n,$$

where $\{r_n\}$ is a sequence of real numbers such that $\liminf_{n \rightarrow \infty} r_n = 0$. Because then for all sufficiently large n : $\inf_{Q \in \mathcal{P}} \alpha_n(Q) - \inf_{Q \in \mathcal{P}_n} \alpha_n(Q) \geq r_n - n^{-1}$. We construct such Q^0 as follows. Fix $C \in Q$. If $|C| \leq h_n$, then let C be included in Q^0 . If $|C| > h_n$, then divide C equally into subintervals of width h_n^* . Set h_n^* equal to h_n if $|C| = \infty$. If $|C| < \infty$, then there is a positive integer $m \geq 2$ such that $h_n(m-1) \leq |C| < h_n m$. Set h_n^* equal to $|C|/m$ in this case. Let the subintervals thus made be included in Q^0 . Repeat this for each $C \in Q$. Then we have

$$(3.9) \quad \alpha_n(Q) - \alpha_n(Q^0) \\ = n^{2/3} \left\{ \int_I (f_Q(x) - f(x))^2 dx - \int_I (f_{Q^0}(x) - f(x))^2 dx \right. \\ \left. + (1/n) \left(\sum_{C \in Q, C \in Q^0} F(C)(1 - F(C))/|C| \right. \right. \\ \left. \left. - \sum_{C \in Q, C \in Q^0} F(C)(1 - F(C))/|C| \right) \right\}.$$

Note that Q^0 is finer than Q in the sense that for each $C \in Q$ there is $C^0 \in Q^0$ such that $C^0 \subset C$. Thus

$$(3.10) \quad \int_I (f_Q(x) - f(x))^2 dx \geq \int_I (f_{Q^0}(x) - f(x))^2 dx.$$

Observe that

$$(3.11) \quad \sum_{C \neq Q, C \in Q^0} F(c)/|C| \leq \left(\inf_{C \in Q^0} |C| \right)^{-1} \leq Qh_n^{-1}.$$

Combining (3.6), (3.9), (3.10) and (3.11) we have

$$\alpha_n(Q) - \alpha_n(Q^0) \geq -n^{-1/3} \left(2h_n^{-1} + \int_I (f(x))^2 dx \right).$$

This in turn implies the existence of $\{r_n\}$ in (3.8). The proof of the proposition is now complete. \square

PROOF OF PROPOSITION 2. The mean value theorem implies that for each $i \in Z$ there is $x_i \in \Delta_i$ such that

$$(3.12) \quad \int_{\Delta_i} |f'(x)f(x)|^{2/3} dx = w_n |f'(x_i)f(x_i)|^{2/3},$$

where Δ_i is defined as at (2.4). For each $i \in Z$ define a real-valued function $H_i(\cdot)$ on $(0, \infty)$ as

$$(3.13) \quad H_i(t) := \left(\frac{1}{12}\right)w_n^3 (f'(x_i))^2/t^2 + f(x_i)t/n.$$

Put $Z^0 := \{i \in Z | f(x_i) \neq 0 \text{ and } f'(x_i) \neq 0\}$. For each $i \in Z^0$ let

$$t_i^* := 6^{-1/3} \left\{ (f'(x_i))^2 / f(x_i) \right\}^{1/3} n^{1/3} w_n.$$

t_i^* minimizes $H_i(t)$ and $H_i(t_i^*) = \left(\frac{3}{2}\right)6^{-1/3} \int_{\Delta_i} |f'(x)f(x)|^{2/3} dx n^{-2/3}$. By (3.12) $\int_{\Delta_i} |f'(x)f(x)|^{2/3} dx = 0$ for $i \in Z^0$. Thus,

$$(3.14) \quad \sum_{i \in Z^0} H_i(t_i^*) = \left(\frac{3}{2}\right)6^{-1/3} \int_I |f'(x)f(x)|^{2/3} dx n^{-2/3}.$$

Let $Q(w_n, \mathbf{k})$ be a partition in $\mathcal{P}(w_n, U_n)$, where $\mathcal{P}(w_n, U_n)$ is defined as (2.6). It is easy to see that

$$(3.15) \quad \text{MISE}(n, Q(w_n, \mathbf{k})) = \sum_{i \in Z} H_i(k_i) + o(n^{-2/3}),$$

where the $o(n^{-2/3})$ term is bounded by a quantity which converges to zero faster than $n^{-2/3}$, uniformly in both $\{x_i\}$ and \mathbf{k} . In the light of (3.14) and (3.15) the conclusion follows if we verify the existence of a sequence of positive integers $\{k_i^* | i \in Z, k_i^* \leq U_n\}$ such that

$$(3.16) \quad \sum_{i \in Z} H_i(k_i^*) - \sum_{i \in Z^0} H_i(t_i^*) = o(n^{-2/3}).$$

If both $f'(x_i)$ and $f(x_i)$ are zero, then $H_i(t) \equiv 0$. Thus, assume that at least one of $f'(x_i)$ and $f(x_i)$ is not zero. For any $i \in Z^0$, let k_i^0 be the smallest positive

integer such that $H_i(k_i^0)$ is closest to $H_i(t_i^*)$. For each $i \in Z$ define k_i^* as

$$k_i^* = \begin{cases} k_i^0 & \text{if } i \in Z \text{ and } k_i^0 \leq U_n, \\ U_n - 1 & \text{if } i \in Z \text{ and } k_i^0 > U_n, \\ U_n & \text{if } f'(x_i) \neq 0 \text{ and } f(x_i) = 0, \\ 1 & \text{if } f'(x_i) = 0 \text{ and } f(x_i) \neq 0. \end{cases}$$

Now consider the four possible cases and the respective inequalities below. For the first two cases we utilize the following relation: If $i \in Z^0$, then

$$(3.17) \quad 2\left(\frac{1}{12}\right)(f'(x_i))^2 w_n^3 t^{-2} \geq (f(x_i)/n)t, \quad \text{iff } t_i^* \geq t.$$

Case 1. $i \in Z^0$ and $k_i^0 \leq U_n$. By (3.17)

$$H_i(t_i^* + 1) \leq \left(\frac{3}{2}\right)(f(x_i)/n)(t_i^* + 1)$$

and

$$H_i(t_i^*) = \left(\frac{3}{2}\right)(f(x_i)/n)t_i^*.$$

Thus,

$$(3.18) \quad \begin{aligned} H_i(k_i^0) - H_i(t_i^*) &\leq H_i(t_i^* + 1) - H_i(t_i^*) \\ &= \left(\frac{3}{2}\right)(f(x_i)/n). \end{aligned}$$

Case 2. $i \in Z^0$ and $k_i^0 > U_n$. Then $t_i^* + 1 \geq U_n$. Thus,

$$(3.19) \quad H_i(U_n - 1) \leq \left(\frac{1}{4}\right)(f'(x_i))^2 w_n^3 (U_n - 1)^{-1}.$$

Case 3. $f'(x_i) \neq 0$ and $f(x_i) = 0$. Then

$$(3.20) \quad H_i(U_n) = \left(\frac{1}{12}\right)(f'(x_i))^2 w_n^3 U_n^{-2}.$$

Case 4. $f'(x_i) = 0$ and $f(x_i) \neq 0$. Then

$$(3.21) \quad H_i(1) = f(x_i)/n.$$

Combining (3.18)–(3.21) we have

$$(3.22) \quad \begin{aligned} &\sum_{i \in Z} H_i(k_i^*) - \sum_{i \in Z^0} H_i(t_i^*) \\ &\leq \left(\frac{3}{2}\right)\left(\sum_{i \in Z} f(x_i)\right)/n + \left(\frac{1}{4}\right)\left(\sum_{i \in Z} (f'(x_i))^2\right) w_n^3 (U_n - 1)^{-2}. \end{aligned}$$

By Corollary 2.24 of Freedman and Diaconis (1981) we have

$$(3.23) \quad \left| \sum_{i \in Z} f(x_i) w_n - \int_I f(x) dx \right| \leq w_n \int_I |f'(x)| dx$$

and

$$(3.24) \quad \left| \sum_{i \in Z} (f'(x_i))^2 w_n - \int_I (f'(x))^2 dx \right| \leq 2w_n \int_I |f'(x)f(x)| dx.$$

Equation (3.22) together with (3.23) and (3.24) implies

$$\left| \sum_{i \in Z} H_i(k_i^*) - \sum_{i \in Z^0} H_i(t_i^*) \right| = o((nw_n)^{-1} + (w_n/U_n)^2) \\ = o(n^{-2/3}).$$

The last relation holds because of (2.5). This completes the proof of Proposition 2. \square

Acknowledgments. This work is part of the author's Ph.D. thesis at Yale University, written under the supervision of Professor J. A. Hartigan, whose guidance and suggestions are gratefully acknowledged and appreciated. Thanks are also due to Professors K. Takeuchi and Y. Hosoya for helpful comments.

REFERENCES

- CHOW, Y.-S., GEMAN, S. and WU, L. D. (1983). Consistent crossvalidated density estimation. *Ann. Statist.* **11** 25–38.
- FREEDMAN, D. and DIACONIS, P. (1981). On the histogram as a density estimator: L_2 theory. *Z. Wahrsch. verw. Gebiete* **57** 453–476.
- KOGURE, A. (1986). Optimal cells for a histogram. Ph.D. thesis, Yale Univ.
- RUDEMO, M. (1982). Empirical choice of histograms and kernel density estimators. *Scand. J. Statist.* **9** 65–78.
- SCOTT, D. W. (1979). On optimal and data-based histograms. *Biometrika* **66** 605–610.
- STONE, C. (1985). An asymptotically optimal histogram selection rule. In *Proc. of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer* (L. M. Le Cam and R. A. Olshen, eds.) **2** 513–520. Wadsworth, Belmont, Calif.
- TAPIA, R. A. and THOMPSON, J. R. (1978). *Nonparametric Probability Density Estimation*. Johns Hopkins Univ. Press, Baltimore, Md.

FACULTY OF ECONOMICS
FUKUSHIMA UNIVERSITY
MATSUKAWA-MACHI, FUKUSHIMA 960-12
JAPAN