# A NOTE ON OPTIMUM GROUPING AND THE RELATIVE DISCRIMINATING POWER OF QUALITATIVE TO CONTINUOUS NORMAL VARIATES

Yung-tai HUNG and Anant M. KSHIRSAGAR

*Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA*

*Abstract:* Ogawa (1951) considered the efficiency of estimation of the population mean from suitably chosen order statistics in large samples. Cox (1957) has considered the relative amount of information retained by grouping the normal curve. Cochran and Hopkins (1961) determined the discriminating power retained after partitioning normally distributed variates into qualitative ones in multivariate classification problems. And Connor (1972) discussed the asymptotic efficiencies of the test for the trend using $m$ groups formed from a continuous variable. The same expression appears in all these investigations. This note throws some more light on the occurrence of the same expression in these seemingly unrelated problems.

*Keywords:* optimum grouping, efficiency of estimation using order statistics, multivariate classification with qualitative data.

## 1. Introduction

Consider a normal variate $X$ with mean $\mu$ and variance $\sigma^2$, let the $\lambda_i$-quantile of the population be $\chi_i$, $i = 1, \ldots, m-1$, that is,

$$\lambda_i = \int_{-\infty}^{\chi_i} f(t)\, dt, \quad 0 < \lambda_1 < \cdots < \lambda_{m-1} < 1, \quad (1.1)$$

where $f$ denotes the density function of $X$. To estimate $\mu$ using the limiting distribution of suitably chosen order statistics, Ogawa (1951) derived the quantity

$$\varepsilon = \sum_{i=1}^{m} \frac{(\phi_i - \phi_{i-1})^2}{\Phi_i - \Phi_{i-1}} \quad (1.2)$$

where

$$\phi_0 = \Phi_0 = 0, \qquad \phi_i = \phi(U_i), \qquad \Phi_i = \Phi(U_i),$$

$$U_i = (\chi_i - \mu)/\sigma,$$

$\phi$ and $\Phi$ are the density and cumulative function respectively of an $N(0, 1)$ variate. $\varepsilon$ is called the efficiency of estimation. Thus if $\sigma^2$ is known and $m$, the number of groups, is specified then an optimum spacing or maximum efficiency can be

achieved by choosing the cutpoints $\chi_i$, $i = 1, \ldots, m-1$ so as to maximize (1.2).

On a somewhat different subject, however, in order to condense observations from $X$ into a small number of groups while retaining as much information as possible, Cox (1957) proposed that group boundaries be chosen in such a way that the quantity

$$E\left\{ (X - \xi_i)^2/\sigma^2 \right\} \quad (1.3)$$

is minimized. Note that $\xi_i$ denotes the mean of the $i$th group to which $\chi$ is assigned. (1.3) can also be written as

$$1 - \sum_{i=1}^{m} P_i(\xi_i - \mu)^2/\sigma^2 \quad (1.4)$$

where $P_i$ is the probability of an observation falling into the $i$th group. The problem reduces to maximizing the second term of (1.4). The criterion was also proposed by Connor (1972) in a problem of maximizing the asymptotic efficiency of the test for the trend using optimal grouping. He also gave a list of other objectives relating to grouping which resulted in using this maximization criterion. We

shall prove in this note that the quantity to be maximized here is mathematically identical to (1.2).

Yet on another different subject the same quantity appeared again in Cochran and Hopkins (1961). The standard discriminating procedure in the case of two $K$-variate normal populations $\Pi_1$ and $\Pi_2$ with different means but the same variance-covariance matrices gives the chance of misclassification $\Pr(\Pi_2 \mid \Pi_1)$ as

$$\alpha_1 = \Phi\left(-\tfrac{1}{2}\Delta_K\right) \tag{1.5}$$

where $\Delta_K$ is the Mahalanobis distance between the two populations. Suppose each continuous normal variable is now replaced by an $m$-category qualitative variable, and also suppose that all variables $X_1, \ldots, X_K$ are independent each with unit variance. Then let $P_{si}$, $P'_{si}$ be the probability that an observation from $\Pi_1$ and $\Pi_2$ respectively falls into the $i$th group for the $s$th variate and let $R = \sum_{s=1}^{K} \log(P_{si}/P'_{si})$, Cochran and Hopkins (1961) found that using $R$ as the classification criterion (that is, assigning an observation to $\Pi_1$ if $R \geqslant 0$, to $\Pi_2$ otherwise) the chance of misclassification, namely, $\Pr(R < 0 \mid \Pi_1)$ is asymptotically

$$\Phi\left(-\frac{E(R \mid \Pi_1)}{\sqrt{\mathrm{Var}(R \mid \Pi_1)}}\right) = \Phi\left(-\tfrac{1}{2}\Delta'_K\right), \tag{1.6}$$

say.

It turns out that minimizing (1.6) is the same as maximizing $\Delta'_K$ and $\Delta'^2_K/\Delta^2_K$ is nothing but (1.2).

## 2. Proofs

We first show that the second term of (1.4) is in fact (1.2). Without loss of generality assume $\sigma^2 = 1$, $\mu = 0$ then

$$\sum_{i=1}^{m} P_i(\xi_i - \mu)^2/\sigma^2 = \sum_{i=1}^{m} P_i\xi_i^2. \tag{2.1}$$

Note that

$$P_i = \int_{U_{i-1}}^{U_i} \frac{1}{\sqrt{2\pi}} \exp\left(-\tfrac{1}{2}t^2\right) \mathrm{d}t = \Phi_i - \Phi_{i-1} \tag{2.2}$$

where $(u_i, u_{i-1})$ are the boundaries of the $i$th

group. Also,

$$\xi_i = E(X \mid U_{i-1} < X < U_i)$$
$$= \frac{1}{\Pr(U_{i-1} < X < U_i)} \int_{U_{i-1}}^{U_i} y \frac{1}{\sqrt{2\pi}} \exp\left(-\tfrac{1}{2}y^2\right) \mathrm{d}y$$
$$= (\phi_{i-1} - \phi_i)/(\Phi_i - \Phi_{i-1}). \tag{2.3}$$

Putting $\xi_i$ back into equation (2.1), the desired result follows.

We now show that $\Delta'^2_K/\Delta^2_K$ is the same as (1.2). Assume that the mean of $X_s$ ($s = 1, \ldots, K$) is $\delta_s$ in $\Pi_1$ and 0 in $\Pi_2$, then assuming

$$\Phi\left(x - \frac{\delta}{2}\right) \doteq \Phi(x) - \frac{\delta}{2}\phi(x), \tag{2.4}$$

where $\doteq$ stands for approximate equality, neglecting higher order terms, we obtain,

$$P_{si} = \int_{U_{i-1} + (1/2)\delta_s}^{U_i + (1/2)\delta_s} \frac{1}{\sqrt{2\pi}} \exp\left[-\tfrac{1}{2}(t - \delta_s)^2\right] \mathrm{d}t$$
$$\doteq (\Phi_i - \Phi_{i-1}) + \tfrac{1}{2}\delta_s(\phi_{i-1} - \phi_i). \tag{2.5}$$

Similarly,

$$P'_{si} \doteq (\Phi_i - \Phi_{i-1}) + \tfrac{1}{2}\delta_s(\phi_i - \phi_{i-1}). \tag{2.6}$$

so

$$\log\frac{P_{si}}{P'_{si}} = \log\left[1 + \frac{\delta_s(\phi_{i-1} - \phi_i)}{\Phi_i - \Phi_{i-1} - (1/2)\delta_s(\phi_{i-1} - \phi_i)}\right]$$
$$\doteq \frac{\delta_s(\phi_{i-1} - \phi_i)}{\Phi_i - \Phi_{i-1}}, \tag{2.7}$$

ignoring higher order terms, and

$$E(R \mid \Pi_1) = \sum_{s=1}^{K} \sum_{i=1}^{m} P_{si} \log(P_{si}/P'_{si})$$
$$\doteq \sum_{s=1}^{K} \tfrac{1}{2}\delta_s^2 \sum_{i=1}^{m} \frac{(\phi_{i-1} - \phi_i)^2}{\Phi_i - \Phi_{i-1}}, \tag{2.8}$$

$$\mathrm{Var}(R \mid \Pi_1) \doteq \sum_{s=1}^{K} \delta_s^2 \sum_{i=1}^{m} \frac{(\phi_{i-1} - \phi_i)^2}{\Phi_i - \Phi_{i-1}} \tag{2.9}$$

(for derivations, see Kshirsagar (1972, p. 238)).

As a result,

$$\frac{E(R \mid \Pi_1)}{\sqrt{\mathrm{Var}(R \mid \Pi_1)}} \doteq \sqrt{\tfrac{1}{4} \sum_{s=1}^{K} \delta_s^2 \sum_{i=1}^{m} \frac{(\phi_{i-1} - \phi_i)^2}{\Phi_i - \Phi_{i-1}}}.$$

$$\tag{2.10}$$

Putting (2.9) in (1.6), clearly,

$$\tfrac{1}{4}\Delta_K'^2 \doteq \tfrac{1}{4} \sum_{s=1}^{K} \delta_s^2 \sum_{i=1}^{m} \frac{(\phi_{i-1} - \phi_i)^2}{\Phi_i - \Phi_{i-1}}$$

$$= \tfrac{1}{4} \sum_{s=1}^{K} \delta_s^2 g(U), \qquad (2.11)$$

say.

Since $g(U)$ does not involve $\delta_s$, the values of cutpoints that maximize $g(U)$ are the same for each $s$. Denote the maximum value of $g(U)$ by $g^*(U)$, the maximum of $\tfrac{1}{4}\Delta_K'^2$ will be $\tfrac{1}{4}\Delta_K^2 g^*(U)$ where $\Delta_K^2 = \sum_{s=1}^{K} \delta_s^2 = K\delta_0^2$, say, is the square of the Mahalanobis distance between $\Pi_1$ and $\Pi_2$ based on the independent variables $X_1, \ldots, X_K$, all of which have unit variances. Note that assuming (2.4) and ignoring higher order terms in (2.7), $\Delta_K'^2$ is approximately equal to $\sum_{s=1}^{K} \delta_s^2 g(U)$. Therefore $\Delta_K'^2 / \Delta_K^2$ is approximately equal to $g(U)$ which is (1.2). Hence the proof is completed.

## References

Cochran, W.G. and C.E. Hopkins (1961), Some classification problems with multivariate qualitative data, *Biometrics* **17**, 10–32.

Connor, R.J. (1972), Grouping for testing trends in categorical data, *Journal of the American Statistical Association* **67**, 601–604.

Cox, D.R. (1957), Note on Grouping, *Journal of the American Statistical Association* **52**, 543–547.

Kshirsagar, A.M. (1972), *Multivariate Analysis*, Marcel Dekker, New York.

Ogawa, J. (1951), Contributions to the Theory of Systematic Statistics, I, *Osaka Mathematical Journal* **3**, 175–213.