

# SOME DISTRIBUTIONS ARISING IN MATCHING PROBLEMS

BY P. V. KRISHNA IYER

*Defence Science Laboratory, New Delhi*

## I. INTRODUCTION

SIMPLE matching is said to occur when elements of the same kind occur in the same order in two sets, each of  $n$  elements of  $k$  types, arranged at random, with fixed or varying probabilities. When two sets of  $n$  cards containing  $n_{11}, n_{12}, n_{21}, n_{22}, n_{31}, n_{32}, \dots, n_{k1}, n_{k2}$ ; cards of black, white, red, etc., colours are arranged at random in two rows and if we compare the cards of the two rows in order, we may find a number of instances where cards of the same kind occur together in the same order in both the rows. Such occurrences are usually termed 'matchings'. The idea of matchings between two sets can be extended for three or more sets. The matchings here may be between all the sets simultaneously or any two or more consecutive sets.

The distribution of the number of matchings between two or more sets of cards for varying compositions have been studied by Anderson (1943), Bartlett (1937), Battin (1942), Greenwood (1938), Greville (1941), Kaplansky and Riordon (1945), Wilks (1946) and many others. Most of this work relates to the number of matchings between cards of the same kind for 'finite sampling' with  $n_{11}, n_{12}, n_{21}, n_{22}, n_{31}, n_{32}, \dots, n_{k1}, n_{k2}$ ; cards of different colours. Obviously the probability for the occurrence of the cards changes as the order changes. It is possible that under certain circumstances this probability may be fixed for the various colours. Such a scheme wherein the probability for the occurrence of any of the cards is independent of the order of the card in the set may be termed 'infinite sampling'. (This appears to be justifiable because when sampling is carried out from an infinite population the probability for any of the cards can be assumed to be fixed.) Much work does not appear to have been done on the distribution of the number of matchings for infinite sampling. The object of this paper is to consider simultaneously, both for finite and infinite sampling, the distribution of the number of matching in two or more sets of cards by certain special methods developed by the author. It is also proposed to extend these methods for the discussion of distributions of a quantitative nature arising by assigning scores, say,  $\theta_1, \theta_2, \theta_3, \dots, \theta_k$  to the cards of various colours for the different sets.

For two sets we shall deal with the distributions of  $\Sigma(x_r - y_r)$ ,  $\Sigma|x_r - y_r|$  and  $\Sigma(x_r - y_r)^2$ , where  $x_r$  and  $y_r$  stand for the scores of the  $r$ -th cards in the first and the second set respectively. For three sets we shall discuss the distribution of  $\Sigma|x_r - 2y_r + z_r|$  and  $\Sigma(|x_r - y_r| + |y_r - z_r|)$ , where  $x_r, y_r$  and  $z_r$  are the  $r$ -th scores in the three sets of cards. In fact we can consider the distribution of any function  $f(x_r, y_r, z_r)$  of the scores.

## 2. SIMPLE MATCHINGS

Suppose there are two sets of cards,  $D_1$  and  $D_2$ , of  $k$  colours. Let there be  $n$  cards in each set arranged at random in a sequence. Assume that the probabilities for black, white, red, etc., cards in the two sets are fixed and are  $p_{11}, p_{21}, \dots, p_{k1}$  and  $p_{12}, p_{22}, \dots, p_{k2}$ . We shall consider the distribution of (i) the number of matched pairs of a given colour and (ii) the total number of matched pairs for all kinds of colours. For finite sampling, Battin (1942), Wilks (1946) and others have discussed this problem by obtaining the generating functions of the distributions. We shall investigate this problem by a general approach which is applicable both for finite and infinite sampling.

### (a) Number of matchings for a single colour in two sets

The probability generating function for the number of matched pairs for a given, say the  $r$ -th, colour is given by

$$[p_{r1} p_{r2} x + (1 - p_{r1} p_{r2})]^n \quad (2.1)$$

This is obvious, because the probability for a single matching of the  $r$ -th colour is  $p_{r1} p_{r2}$  and the P.G.F. for a single card is  $[p_{r1} p_{r2} x + (1 - p_{r1} p_{r2})]$ . If  $p_{r1}$  and  $p_{r2}$  are fixed, the probability for the  $s$ -th matching in the sequence is independent of the number of matchings noted earlier or later. Therefore the P.G.F. for  $n$  cards is (2.1).

The cumulants are the same as those for the binomial distribution, where  $P$  and  $Q$  are  $p_{r1} p_{r2}$  and  $(1 - p_{r1} p_{r2})$  respectively. Hence the first four cumulants are

$$\kappa_1 = nP, \quad \kappa_2 = nPQ, \quad \kappa_3 = nPQ(Q - P), \quad \kappa_4 = nPQ(1 - 6PQ) \quad (2.2)$$

The probability of  $s$  matchings for finite sampling follow from (2.1) by the following procedure: From (2.1) the probability for  $s$  matchings is

$$\binom{n}{s} (p_{r1} p_{r2})^s (1 - p_{r1} p_{r2})^{n-s} \quad (2.3)$$

Expanding (2.3) it becomes

$$\binom{n}{s} \left[ \left( p_{r_1} p_{r_2} \right)^s - \binom{n-s}{1} \left( p_{r_1} p_{r_2} \right)^{s+1} + \binom{n-s}{2} \left( p_{r_1} p_{r_2} \right)^{s+2} + (-1)^{n-s} \left( p_{r_1} p_{r_2} \right)^n \right] \quad (2.4)$$

Substitute

$$\frac{n_{r_1}^{(s)} n_{r_2}^{(t)}}{n^{(s)} n^{(t)}} = p_{r_1}^s p_{r_2}^t$$

in (2.4). Then (2.4) reduces to

$$\binom{n}{s} \left[ \frac{n_{r_1}^{(s)} n_{r_2}^{(s)}}{(n^{(s)})^2} - \binom{n-s}{1} \frac{n_{r_1}^{(s+1)} n_{r_2}^{(s+2)}}{(n^{(s+1)})^2} + \binom{n-s}{2} \frac{n_{r_1}^{(s+2)} n_{r_2}^{(s+2)}}{(n^{(s+2)})^2} \dots \right] \quad (2.5)$$

and is the probability for  $s$  matchings.

The cumulants of the distribution can be obtained from the first four factorial moments for infinite sampling. The factorial moments for the infinite sampling are

$$\left. \begin{aligned} \mu'_{[1]} &= n p_{r_1} p_{r_2} & ; & \quad \mu'_{[2]} = n^{(2)} p_{r_1}^2 p_{r_2}^2 & ; \\ \mu'_{[3]} &= n^{(3)} p_{r_1}^3 p_{r_2}^3 & ; & \quad \mu'_{[4]} = n^{(4)} p_{r_1}^4 p_{r_2}^4 & \end{aligned} \right\} \quad (2.6)$$

Making the substitution used in (2.5), the first four factorial moments for finite samplings reduce to

$$\left. \begin{aligned} \mu'_{[1]} &= \frac{n_{r_1} n_{r_2}}{n} & ; & \quad \mu'_{[2]} = \frac{n_{r_1}^{(2)} n_{r_2}^{(2)}}{n^{(2)}} & ; \\ \mu'_{[3]} &= \frac{n_{r_1}^{(3)} n_{r_2}^{(3)}}{n^{(3)}} & ; & \quad \mu'_{[4]} = \frac{n_{r_1}^{(4)} n_{r_2}^{(4)}}{n^{(4)}} & \end{aligned} \right\} \quad (2.7)$$

The second moment reduces to

$$\mu_2 = \frac{n_{r_1}^{(2)} n_{r_2}^{(2)}}{n^{(2)}} + \frac{n_{r_1} n_{r_2}}{n} - \left( \frac{n_{r_1} n_{r_2}}{n} \right)^2 \quad (2.8)$$

It may be observed that the distributions for both finite and infinite sampling tend to the normal form when  $p_{r_1} p_{r_2}$  is finite and  $n, n_{11}, n_{12}, \dots$ , etc., are large. When  $p_{r_1} p_{r_2}$  is very small, the distribution tends to the Poisson form.

(b) *Total number of matchings in two sets*

The P.G.F. for the distribution of the total number of matchings for two sets can be seen to be

$$[(\sum p_{r1} p_{r2}) x + (1 - \sum p_{r1} p_{r2})]^n \quad (2.9)$$

This follows by arguing on the same lines as in (a) above. Taking  $P'$  and  $Q'$  to be

$$P' = \sum_{r=1}^k p_{r1} p_{r2}; \quad Q' = (1 - \sum p_{r1} p_{r2}).$$

the cumulants are symbolically the same as given in (2.2).

The probability for  $s$  matchings from (2.9) is

$$\binom{n}{s} P'^s (1 - P')^{n-s} \quad (2.10)$$

Expanding (2.10), it reduces to

$$\binom{n}{s} \left[ P'^s - \binom{n-s}{1} P'^{s+1} + \binom{n-s}{2} P'^{s+2} + \dots + (-1)^{n-s} P'^n \right] \quad (2.11)$$

Making the substitution used in (2.5), (2.11) after expansion in powers of  $p_{11}, p_{12}, \dots, p_{k2}$ , we get the probability for  $s$  matchings for finite sampling. The cumulants for infinite sampling are the same as shown in (2.2). For finite sampling

$$\left. \begin{aligned} \mu_1' &= \frac{\sum n_{r1} n_{r2}}{n}, \\ \mu_2 &= \frac{\sum n_{r1}^{(2)} n_{r2}^{(2)}}{n^{(2)}} + \frac{\sum n_{r1} n_{r2}}{n} - \left( \frac{\sum n_{r1} n_{r2}}{n} \right)^2. \end{aligned} \right\} \quad (2.12)$$

The distribution of the total number of matchings tends to the normal form because the cumulants can be put as a linear function of  $n$  plus a function  $O\left(\frac{1}{n}\right)$ .

(c) *Matchings in m sets*

The matchings may be between (i) all the sets, (ii) two adjacent sets and (iii) two or more adjacent sets. We shall consider the P.G.F.

for the distribution of the number of matchings between two or more adjacent sets.

The probability for a matching as defined above is

$$P'' = \sum_{l=1}^m \sum_{t=1}^m \sum_{r=1}^k p_{rl} p_{r+1} p_{r+2} \cdots p_{r+t}$$

where  $p_{rl}$  is the probability for a card of a particular colour in the  $l$ -th set. Therefore the P.G.F. for this distribution is

$$[P''x + (1 - P'')]^n \quad (2.13)$$

(d) *Matchings between cards of different kinds*

In sections (a), (b) and (c) above we discussed matching between cards of the same kind. We shall now obtain the distribution of the number of matchings between cards of different kinds. The P.G.F. for the distribution of the number of times that a black and a white card (including *vice versa*) come together is

$$[(p_{11}p_{22} + p_{21}p_{12})x + (1 - p_{11}p_{22} - p_{21}p_{12})]^n \quad (2.14)$$

where  $p_{11}, p_{21}, p_{31}, \dots, p_{k1}$  and  $p_{12}, p_{22}, \dots, p_{k2}$  are the probabilities for first and second set of cards being black, white, etc.

The expected number and the variance of the distribution for infinite sampling are

$$n(p_{11}p_{22} + p_{21}p_{12}) \text{ and } n(p_{11}p_{22} + p_{21}p_{12})(1 - p_{11}p_{22} - p_{21}p_{12}) \quad (2.15)$$

For finite sampling the corresponding values are

$$\begin{aligned} \mu_1' &= \frac{n_{11}n_{22} + n_{21}n_{12}}{n}, \\ \mu_2 &= \frac{1}{n^{(2)}} \left[ n_{11}^{(2)}n_{22}^{(2)} + 2n_{11}n_{22}n_{21}n_{12} + n_{21}^{(2)}n_{12}^{(2)} \right] \\ &\quad + \frac{n_{11}n_{22} + n_{21}n_{12}}{n} - \left( \frac{n_{11}n_{22} + n_{21}n_{12}}{n} \right)^2 \end{aligned} \quad (2.16)$$

### 3. DISTRIBUTIONS FROM MATCHINGS OF QUANTITATIVE MEASUREMENTS

Let the cards of the two sets be assigned the scores  $\theta_1, \theta_2, \dots, \theta_k$  for the different colours. Then the probabilities for  $\theta_1, \theta_2, \dots, \theta_k$  for

the two sets are  $p_{11}, p_{21}, \dots, p_{k1}$  and  $p_{12}, p_{22}, \dots, p_{k2}$  respectively. We shall consider the distribution of  $\Sigma(x_r - y_r)$ ,  $\Sigma |x_r - y_r|$  and  $\Sigma(x_r - y_r)^2$ , where  $x_r$  and  $y_r$  are the scores of  $r$ -th cards in the two sets. This aspect of the problem of matching has been mentioned by Wilks (1947) in his book. But no attempt appears to have been made so far to discuss the actual distributions.

The difference between the first and the second sets taken in order can be any of the values shown in the matrix  $[(\theta_r - \theta_s)]$ , where  $r$  and  $s$  take all values from 1 to  $k$ . The corresponding probabilities are given by  $[p_{r1} p_{s2}]$ . The generating function of the distribution of  $\Sigma |x_r - y_r|$  is given by

$$\left[ \sum_{r=1}^k p_{r1} \sum_{t=1}^k p_{t2} x^{| \theta_r - \theta_t |} \right]^n \tag{3.1}$$

Expanding (3.1) and taking the coefficients of  $x^s$  we get the probability for  $\Sigma |x_r - y_r|$  being equal to  $s$ .

If the number of  $\theta$ 's in the first set is fixed, *i.e.*, there are  $n_{11} \theta_1$ 's,  $n_{21} \theta_2$ 's,  $\dots$ ,  $n_{k1} \theta_{k1}$ 's in set 1, while the probability for the  $\theta$ 's in the second set is fixed, as in infinite sampling, then the probability for  $\Sigma |x_r - y_r|$  to be  $s$  is the coefficient of  $x^s$  in

$$\frac{n!}{n_{11}! n_{21}! \dots n_{k1}!} p_{11}^{n_{11}} p_{21}^{n_{21}} \dots p_{k1}^{n_{k1}} \prod_{r=1}^k \left[ \sum_{t=1}^k p_{t2} x^{| \theta_r - \theta_t |} \right]^{n_{r1}} \tag{3.2}$$

This follows by expanding (3.1) and substituting

$$p_{11}^r p_{21}^s p_{31}^t \dots = \frac{n_{11}^{(r)} n_{21}^{(s)} n_{31}^{(t)} \dots}{n^{(r+s+t \dots)}} \tag{3.3}$$

in the expansion. All the other terms will be zero. In fact (3.2) reduces to

$$\prod_{r=1}^k \left[ \sum_{t=1}^k p_{t2} x^{| \theta_r - \theta_t |} \right]^{n_{r1}} \tag{3.4}$$

If both sets refer to finite sampling then the probability for  $\Sigma |x_r - y_r|$  to be  $s$  is the coefficient of  $x^s$  in (3.4) subject to the substitution

$$p_{12}^r p_{22}^s p_{32}^t \dots = \frac{n_{12}^{(r)} n_{22}^{(s)} n_{32}^{(t)} \dots}{n^{(r+s+t \dots)}} \tag{3.5}$$

Using (3.5) it will be seen that (3.4) will reduce to a fairly simple form.

Suppose now  $\theta_1 = 0$ ,  $\theta_2 = 1$ ,  $\theta_3 = 2$ ,  $\theta_4 = 3$  and  $\theta_5 = 4$  and  $n_{11} = n_{21} = n_{31} = n_{41} = n_{51} = 5$ . Then (3.2) for infinite sampling of the second set reduces to

$$\frac{1}{5^{25}} (1 + x + x^2 + x^3 + x^4)^{10} (x + 1 + x + x^2 + x^3)^{10} \\ (x^2 + x + 1 + x + x^2)^5 \quad (3.6)$$

The coefficient of  $x^s$  in the expression (3.6) is the probability for  $\Sigma |x_r - y_r| = s$ . The above distribution has been calculated and given in Table I.

The generating functions for  $\Sigma (x_r - y_r)$  and  $\Sigma (x_r - y_r)^2$  can also be obtained from (3.2) and (3.4) by substituting  $(\theta_r - \theta_s)$  or  $(\theta_r - \theta_s)^2$  for  $|\theta_r - \theta_s|$ .

We may now obtain the moments of the above distributions. They can be obtained either from the generating functions or by using certain results developed by the author for calculating factorial moments. As the latter method appears to be simpler it is proposed to obtain the first and the second moments by this method.

$$\mu_1' = E \left\{ \sum_{r,s=1}^k x_{rs} \mid \theta_r - \theta_s \right\}$$

Now  $|\theta_r - \theta_s|$  is fixed, while  $x_{rs}$ , the number of times that  $\theta_r$  and  $\theta_s$  will occur together, varies. Therefore

$$E \left\{ \sum_{r,s=1}^k x_{rs} \mid \theta_r - \theta_s \right\} = \sum_{r,s=1}^k |\theta_r - \theta_s| E(x_{rs}) \\ = n \sum_{r,s=1}^k |\theta_r - \theta_s| p_{r1} p_{s2} \quad (3.7)$$

The second moment can be obtained by taking the expectation of

$$E \left\{ \sum_{r,s=1}^k \delta x_{rs} \mid \theta_r - \theta_s \right\}^2 \quad (3.8)$$

Expanding (3.8), we get

$$\mu_2 = \Sigma (\theta_r - \theta_s)^2 E (\delta x_{rs})^2 + 2 \Sigma |\theta_r - \theta_s| |\theta_r - \theta_t| E (\delta x_{rs} \delta x_{rt}) \\ + 2 \Sigma |\theta_r - \theta_s| |\theta_t - \theta_s| E (\delta x_{rs} \delta x_{ts}) \\ + 2 \Sigma |\theta_r - \theta_s| |\theta_t - \theta_u| E (\delta x_{rs} \delta x_{tu}) \quad (3.9)$$





TABLE I—Contd.

| $\sum_1^{25} [1x_r - y_r]$ | Probability                              |
|----------------------------|--|
| 47                         | .0 3 3 2 8 8 0 5 2 5 8 0 7 1 5 1 2 1 9   |
| 48                         | .0 2 6 6 5 3 5 1 7 4 4 4 7 2 7 9 0 4 9   |
| 49                         | .0 2 0 7 0 0 9 7 5 9 5 3 8 1 5 9 6 9 7   |
| 50                         | .0 1 5 5 8 9 1 1 9 2 7 3 8 0 5 5 8 2 4   |
| 51                         | .0 1 1 3 7 7 4 0 4 8 1 8 7 1 7 0 6 1 2   |
| 52                         | .0 0 8 0 4 3 0 7 8 9 8 3 3 4 3 5 1 6     |
| 53                         | .0 0 5 5 0 4 1 9 1 7 8 0 5 2 9 0 6 9 1   |
| 54                         | .0 0 3 6 4 3 8 1 9 3 8 7 7 4 4 1 2 6 2   |
| 55                         | .0 0 2 3 3 1 7 0 9 1 1 0 9 5 4 5 6 6 1   |
| 56                         | .0 0 1 4 4 1 0 0 9 4 8 9 6 6 6 0 6 9 5   |
| 57                         | .0 0 0 8 5 9 2 2 9 3 1 4 6 9 6 9 1 6 0   |
| 58                         | .0 0 0 4 9 3 7 7 0 1 6 0 2 2 3 2 2 4 3   |
| 59                         | .0 0 0 2 7 3 1 3 6 9 8 7 2 8 5 6 2 3 9   |
| 60                         | .0 0 0 1 4 5 2 3 7 1 4 4 2 7 9 1 2 1 6   |
| 61                         | .0 0 0 0 7 4 1 2 1 4 8 8 8 1 3 8 0 6 1   |
| 62                         | .0 0 0 0 3 6 2 4 2 8 7 8 2 0 2 8 8 9 9   |
| 63                         | .0 0 0 0 1 6 9 4 5 6 3 3 7 9 1 2 6 1 5   |
| 64                         | .0 0 0 0 0 7 5 5 9 2 5 2 9 2 4 5 8 1 3   |
| 65                         | .0 0 0 0 0 3 2 0 9 0 6 8 7 9 1 1 1 1 6   |
| 66                         | .0 0 0 0 0 1 2 9 2 6 7 2 7 3 4 7 2 3 0   |
| 67                         | .0 0 0 0 0 0 4 9 2 4 2 5 6 5 4 0 8 4 6   |
| 68                         | .0 0 0 0 0 0 1 7 6 6 9 5 7 8 3 2 8 4 8   |
| 69                         | .0 0 0 0 0 0 0 5 9 4 4 7 5 3 7 8 4 4 6   |
| 70                         | .0 0 0 0 0 0 0 1 8 6 4 9 5 6 0 1 2 7 2   |
| 71                         | .0 0 0 0 0 0 0 0 5 4 1 9 2 2 5 8 8 7 8   |
| 72                         | .0 0 0 0 0 0 0 0 1 4 4 6 7 2 0 9 7 8 3   |
| 73                         | .0 0 0 0 0 0 0 0 0 0 3 5 1 1 9 6 2 5 4 6 |
| 74                         | .0 0 0 0 0 0 0 0 0 0 7 6 5 0 4 2 3 9 2   |
| 75                         | .0 0 0 0 0 0 0 0 0 0 1 4 6 9 3 8 8 4     |
| 76                         | .0 0 0 0 0 0 0 0 0 0 0 2 4 2 7 9 9 8 7   |
| 77                         | .0 0 0 0 0 0 0 0 0 0 0 0 3 3 2 8 6 0 0   |
| 78                         | .0 0 0 0 0 0 0 0 0 0 0 0 0 3 5 7 0 1 9   |
| 79                         | .0 0 0 0 0 0 0 0 0 0 0 0 0 2 6 8 4 4     |
| 80                         | .0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 7 4     |
|                            | 1.0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0  |

Now

$$\left. \begin{aligned}
 E(\delta x_{rs})^2 &= np_{r1} p_{s2} (1 - p_{r1} p_{s2}) \\
 E(\delta x_{rs} \delta x_{rt}) &= -np_{r1}^2 p_{s2} p_{t2} \\
 E(\delta x_{rs} \delta x_{ts}) &= -np_{r1} p_{t1} p_{s2}^2 \\
 E(\delta x_{rs} \delta x_{tu}) &= -np_{r1} p_{t1} p_{s2} p_{u2}
 \end{aligned} \right\} \quad (3.10)$$

Substituting the above results in (3.9), we get

$$\mu_2 = n [\Sigma (\theta_r - \theta_s)^2 p_{r1} p_{s2} - \{\Sigma | \theta_r - \theta_s | p_{r1} p_{s2}\}^2] \quad (3.11)$$

The formula (3.11) holds good only for infinite sampling.

When the first set refers to finite sampling and the second set to infinite sampling  $\mu_1'$  and  $\mu_2$  are obtained by the substitution used in (2.5) in the moments about the origin. Thus

$$\mu_1' = \Sigma | \theta_r - \theta_s | n_{r1} p_{s2} \quad (3.12)$$

$$\left. \begin{aligned} E(\delta x_{rs})^2 &= n_{r1} p_{s2} (1 - p_{s2}) \\ E(\delta x_{rs} \delta x_{rt}) &= -n_{r1} p_{s2} p_{t2} \\ E(\delta x_{rs} \delta x_{ts}) &= 0 \\ E(\delta x_{rs} \delta x_{tu}) &= 0 \end{aligned} \right\} \quad (3.13)$$

Therefore

$$\begin{aligned} \mu_2 &= [\Sigma (\theta_r - \theta_s)^2 n_{r1} p_{s2} (1 - p_{s2}) - 2 \Sigma | \theta_r - \theta_s | | \theta_r - \theta_t | \\ &\quad \times n_{r1} p_{s2} p_{t2}] \end{aligned} \quad (3.14)$$

When both the sets refer to finite sampling

$$\mu_1' = \frac{1}{n} \left[ \Sigma | \theta_r - \theta_s | n_{r1} n_{s2} \right] \quad (3.15)$$

$$\left. \begin{aligned} E(\delta x_{rs})^2 &= \frac{n^{(2)} n_{r1}^{(2)} n_{s2}^{(2)}}{(n^{(2)})^2} + \frac{n_{r1} n_{s2}}{n} - \left( \frac{n_{r1} n_{s2}}{n} \right)^2 \\ &= \frac{n_{r1} n_{s2}}{n} \left[ \frac{(n_{r1} - 1)(n_{s2} - 1)}{n - 1} - \frac{n_{r1} n_{s2}}{n} + 1 \right] \\ E(\delta x_{rs} \delta x_{rt}) &= \frac{n^{(2)} n_{r1}^{(2)} n_{s2} n_{t2}}{(n^{(2)})^2} - \frac{n^2_{r1} n_{s2} n_{t2}}{n^2} \\ E(\delta x_{rs} \delta x_{ts}) &= \frac{n_{r1} n_{t1} n_{s2}}{n} \left[ \frac{n_{s2} - 1}{n - 1} - \frac{n_{s2}}{n} \right] \\ E(\delta x_{rs} \delta x_{tu}) &= \frac{n_{r1} n_{t1} n_{s2} n_{u2}}{n} \left[ \frac{1}{n - 1} - \frac{1}{n} \right] \end{aligned} \right\} \quad (3.16)$$

It can be seen that (3.9) reduces to

$$\begin{aligned} \mu_2 &= \sum \frac{n_{r1} n_{s2}}{n} \left[ \frac{(n_{r1} - 1)(n_{s2} - 1)}{n - 1} - \frac{n_{r1} n_{s2}}{n} + 1 \right] (\theta_r - \theta_s)^2 \\ &+ 2 \sum \frac{n_{r1} n_{s2} n_{t2}}{n} \left[ \frac{n_{r1} - 1}{n - 1} - \frac{n_{r1}}{n} \right] | \theta_r - \theta_s | | \theta_r - \theta_t | \\ &+ 2 \sum \frac{n_{r1} n_{t1} n_{s2}}{n} \left[ \frac{n_{s2} - 1}{n - 1} - \frac{n_{s2}}{n} \right] | \theta_r - \theta_s | | \theta_t - \theta_s | \\ &+ 2 \sum \frac{n_{r1} n_{s2} n_{t1} n_{u2}}{n} \left[ \frac{1}{n - 1} - \frac{1}{n} \right] | \theta_r - \theta_s | | \theta_t - \theta_u | \end{aligned} \quad (3.17)$$

Using the results for  $E(\delta x_{rs})^2$ ,  $E(\delta x_{rs} \delta x_{rt})$ ,  $E(\delta x_{rs} \delta x_{ts})$ ,  $E(\delta x_{rs} \delta x_{tu})$  we can obtain the expected values and the variances for  $E(x_r - y_r)$  and  $\Sigma(x_r - y_r)^2$  for all the three situations discussed in this paper.

Considering the distribution of  $\Sigma(x_r - y_r)$ ,  $\mu_1' = 0$ ,  $\mu_2 = (3.11)$  or (3.14) or (3.17) with the moduli replaced by the algebraic values.

For the distribution of  $\Sigma(x_r - y_r)^2$ ,  $\mu_1'$  and  $\mu_2$  are given by the same expressions as for  $\Sigma|x_r - y_r|$  with the change that  $|\theta_r - \theta_s|$  and  $(\theta_r - \theta_s)^2$  are replaced by  $(\theta_r - \theta_s)^2$  and  $(\theta_r - \theta_s)^4$  respectively.

4. SOME PARTICULAR CASES

The first two moments of  $\Sigma(x_r - y_r)$ ,  $\Sigma|x_r - y_r|$  and  $\Sigma(x_r - y_r)^2$  for two sets of cards, each of four, five and six groups of four, five and six cards respectively with scores 0, 1, 2, 3, 4 and 5 for the different groups are given in Table II for the various situations mentioned in this paper.

TABLE II  
Mean and variance for finite and infinite sampling

| Type of sampling                   | Distribution of |                         |                     |         |                     |         |                  |                       |         |                  |
|------------------------------------|-----------------|-------------------------|---------------------|---------|---------------------|---------|------------------|-----------------------|---------|------------------|
|                                    | No. of groups   | No. of cards in a group | $\Sigma(x_r - y_r)$ |         | $\Sigma x_r - y_r $ |         |                  | $\Sigma(x_r - y_r)^2$ |         |                  |
|                                    |                 |                         | $\mu_1'$            | $\mu_2$ | $\mu_1'$            | $\mu_2$ | coef. of vari. % | $\mu_1'$              | $\mu_2$ | coef. of vari. % |
| Both sets infinite                 | 4               | 4                       | 0                   | 40      | 20                  | 15      | 19.37            | 40                    | 132     | 28.72            |
|                                    | 5               | 5                       | 0                   | 100     | 40                  | 36      | 15.00            | 100                   | 540     | 23.24            |
|                                    | 6               | 6                       | 0                   | 210     | 70                  | 73.8    | 12.28            | 210                   | 1673    | 19.48            |
| 1st set finite<br>2nd set infinite | 4               | 4                       | 0                   | 20      | 20                  | 14      | 18.71            | 40                    | 116     | 26.93            |
|                                    | 5               | 5                       | 0                   | 50      | 40                  | 33.2    | 14.41            | 100                   | 470     | 21.68            |
|                                    | 6               | 6                       | 0                   | 105     | 70                  | 67.3    | 11.75            | 210                   | 1449    | 18.13            |
| 1st set finite<br>2nd set finite   | 4               | 4                       | 0                   | 0       | 20                  | 13.1    | 18.62            | 40                    | 106.3   | 25.82            |
|                                    | 5               | 5                       | 0                   | 0       | 40                  | 31.3    | 14.07            | 100                   | 416.3   | 20.41            |
|                                    | 6               | 6                       | 0                   | 0       | 70                  | 63.2    | 11.36            | 210                   | 1260    | 16.90            |

In the above table the coefficient of variation is least for the distribution of  $\Sigma|x_r - y_r|$  and therefore between the three statistics

considered above, it may be desirable to take  $\Sigma |x_r - y_r|$  for any statistical test. This question is examined in greater detail in a later section.

### 5. DISTRIBUTIONS ARISING BY MATCHING THREE SETS OF CARDS

We shall discuss in this section some of the distributions arising from three sets of cards. Let the cards have the scores  $\theta_1, \theta_2, \dots, \theta_k$  in the three sets. Assume that the probabilities for their occurrence are  $p_1, p_2, \dots, p_k$  and remain the same for all the sets. (The more general case is not considered here because it is not likely to be of much use in practical applications.) As in the previous sections, distributions of various functions of the  $r$ -th score in the different sets can be considered. For the present attention is devoted to the following two linear functions:

$$\Sigma |x_r - 2y_r + z_r| \text{ and } (\Sigma |x_r - y_r| + \Sigma |y_r - z_r|)$$

where  $x_r, y_r$  and  $z_r$  are the scores of cards in the three sets at the  $r$ -th place. The generating functions of these distributions are complicated and therefore no attempt has been made in this paper to obtain them. It can be shown that all such distributions tend to the normal form as  $n$ , the number of cards in the set increases. The first and the second moments which are useful in examining the deviations from randomness of the scores of the three sets are given below:

#### (a) Distributions of $\Sigma |x_r - 2y_r + z_r|$

Assuming that the probabilities for the occurrence of  $\theta_1, \theta_2, \dots, \theta_k$  are fixed,

$$E\left(\sum_1^n |x_r - 2y_r + z_r|\right) = n E(|x_r - 2y_r + z_r|) \quad (5.1)$$

Now

$$E(|x_r - 2y_r + z_r|) = \sum_{r,s,t=1}^k |\theta_r - 2\theta_s + \theta_t| p_r p_s p_t \quad (5.2)$$

Therefore

$$\mu'_1 = n \sum_{r,s,t=1}^k |\theta_r - 2\theta_s + \theta_t| p_r p_s p_t \quad (5.3)$$

Since the probabilities,  $p$ 's, are fixed and the functions considered involve only the  $r$ -th scores of the three sets, the variance for  $n$  cards (or observations) is

$$n \times (\text{variance of the function of the scores in a single column}).$$

The variance of  $|x_r - 2y_r + z_r|$  for a single column is

$$E(|x_r - 2y_r + z_r|^2) - \{E|x_r - 2y_r + z_r|\}^2$$

$$= \sum_{r,s,t=1}^k (\theta_r - 2\theta_s + \theta_t)^2 p_r p_s p_t - \{\sum |\theta_r - 2\theta_s + \theta_t| p_r p_s p_t\}^2 \quad (5.4)$$

Hence  $n$  (5.4) is the variance of  $\sum |x_r - 2y_r + z_r|$  for fixed  $p$ 's. This result can be used for examining the significance of the difference between three samples.

For finite sampling of  $3n$  observations in which  $\theta_1, \theta_2, \dots, \theta_k$  occur  $n_1, n_2, \dots, n_k$  times respectively and are arranged in three rows, the variance can be obtained by substituting

$$\frac{n_1^{(r)} n_2^{(s)} n_3^{(t)}}{3n^{(r+s+t)}} \text{ for } p_1^r p_2^s p_3^t \dots$$

in the moments about the origin. The results so obtained are useful in deciding whether the distribution of the  $3n$  observations in three rows is random or not.

When the different rows consist of equal number of observations and the probability for the  $\theta$ 's vary from row to row,  $\mu_2$  can be obtained from

$$n [\sum (\theta_{r1} - 2\theta_{s2} + \theta_{t3})^2 p_{r1} p_{s2} p_{t3} - (\sum |\theta_{r1} - 2\theta_{s2} + \theta_{t3}| p_{r1} p_{s2} p_{t3})^2]$$

$$+ \{n \sum |\theta_{r1} - 2\theta_{s2} + \theta_{t3}| p_{r1} p_{s2} p_{t3}\}^2 -$$

$$- \left[ n \sum \frac{|\theta_{r1} - 2\theta_{s2} + \theta_{t3}| n_{r1} n_{s2} n_{t3}}{n^{(3)}} \right]^2 \quad (5.5)$$

by substituting

$$\frac{n_{r1}^{(s)} n_{s2}^{(t)} n_{t3}^{(u)}}{n^{(s)} n^{(t)} n^{(u)}} \text{ for } p_{r1}^s p_{s2}^t p_{t3}^u,$$

where  $n_{r1}, n_{s2}, n_{t3}$  represent the number of  $\theta_r$ 's,  $\theta_s$ 's and  $\theta_t$ 's in the first, second and third rows.

(b) *Distribution of  $\sum |x_r - y_r| + \sum |y_r - z_r|$*

The cumulants of this distribution for fixed  $p$ 's is  $n$  times the cumulants for a single column. Now for a single column

$$\mu'_1 = E(|x_r - y_r| + |y_r - z_r|) = 2 \sum |\theta_r - \theta_s| p_r p_s \quad (5.6).$$

The variance for a single column is given by

$$\begin{aligned}
 \mu_2 &= E \{ |x_r - y_r| + |y_r - z_r| \}^2 - \{ E ( |x_r - y_r| + |y_r - z_r| ) \}^2 \\
 &= E ( |x_r - y_r| )^2 + E ( |y_r - z_r| )^2 + 2E ( |x_r - y_r| |y_r - z_r| ) \\
 &\quad - 4 [ \Sigma | \theta_r - \theta_s | p_r p_s ]^2 \\
 &= 2 \{ [ \Sigma | \theta_r - \theta_s |^2 p_r p_s - ( \Sigma | \theta_r - \theta_s | p_r p_s )^2 ] \\
 &+ [ \Sigma | \theta_r - \theta_s | | \theta_s - \theta_r | p_r^2 p_s + \Sigma | \theta_s - \theta_r | | \theta_r - \theta_s | p_r p_s^2 \\
 &+ ( \Sigma | \theta_r - \theta_s | | \theta_s - \theta_t | + \Sigma | \theta_s - \theta_r | | \theta_r - \theta_t | + \Sigma | \theta_r - \theta_t | | \theta_t - \theta_s | \\
 &+ \Sigma | \theta_s - \theta_t | | \theta_t - \theta_r | + \Sigma | \theta_t - \theta_r | | \theta_r - \theta_s | + \Sigma | \theta_t - \theta_s | | \theta_s - \theta_r | ) \\
 &\quad p_r p_s p_t - ( \Sigma | \theta_r - \theta_s | p_r p_s )^2 \} \tag{5.7}
 \end{aligned}$$

The results (5.6) and (5.7) are useful in deciding whether three samples belong to the same population or not.

By following the procedure described in (a) above, the expected value and the variance for finite sampling can be deduced.

It is obvious that the methods described in this section can be extended for four, five and more rows. The functions taken may be orthogonal linear functions of the  $r$ -th scores in the various sets. From Fisher and Yates Tables the orthogonal functions for four rows and five rows can be as following:

*Four rows*

$$\left. \begin{aligned}
 \text{(i)} & - 3x_r - y_r + z_r + 3w_r \\
 \text{(ii)} & x_r - y_r - z_r + w_r \\
 \text{(iii)} & - x_r + 3y_r - 3z_r + w_r
 \end{aligned} \right\} \tag{5.8}$$

*Five rows*

$$\left. \begin{aligned}
 \text{(i)} & - 2x_r - y_r + w_r + 2v_r \\
 \text{(ii)} & 2x_r - y_r - 2z_r - w_r + 2v_r \\
 \text{(iii)} & - x_r + 2y_r - 2w_r + v_r \\
 \text{(iv)} & x_r - 4y_r + 6z_r - 4w_r + v_r
 \end{aligned} \right\} \tag{5.9}$$

The expected values and the variances of (5.8) and (5.9) can be worked out on the same lines as for three rows. Here also all the

cumulants will be linear functions of the number of observations in each row.

## 6. TESTING OF TWO SAMPLES

We have considered in this paper the distributions of three statistics  $\Sigma(x_r - y_r)$ ,  $\Sigma |x_r - y_r|$  and  $\Sigma(x_r - y_r)^2$  for two samples. The distributions of all of them tend to the normal form as  $n$  tends to infinity. All the three statistics can, therefore, be used for testing the significance of the difference between two given samples. Before discussing the actual procedure for testing we may note the following:

- (1) The statistics  $\frac{1}{n} \Sigma(x_r - y_r)$ ,  $\frac{1}{n} \Sigma |x_r - y_r|$  and  $\frac{1}{n} \Sigma(x_r - y_r)^2$  are all consistent.
- (2) The single tail confidence intervals based on the above statistics are consistent in the sense defined by Wald and Wolfowitz (1940), when  $n$  tends to infinity. This means, the probability of rejecting the null hypothesis, when it is false, approaches one as the sample size increases.

The consistency of the statistics can be easily established by using Tchebycheff's inequality. This readily follows from the fact that the variances of the different statistics are inversely proportional to  $n$ , the size of the samples. The consistency of the single tail regions can be proved by using the technique employed by Man and Whitney (1947).

Let  $H_0$  be the hypothesis that the probabilities for the occurrence of  $\theta_1, \theta_2, \dots, \theta_k$  in the two samples are  $p_1, p_2, \dots, p_k$ . Let the probabilities for the alternate hypothesis  $H_1$  be  $p_1 + \delta_1, p_2 + \delta_2, \dots, p_k + \delta_k$  such that

$$\sum_1^k \delta_r = 0 \quad (6.1)$$

Taking the statistic  $\Sigma |x_r - y_r|$ , the expected value and the variances for the hypotheses  $H_0$  and  $H_1$  are as follows:

*Hypothesis  $H_0$*

$$\mu_1' = n \Sigma | \theta_r - \theta_s | p_r p_s = n z_1 \quad (6.2)$$

$$\begin{aligned} \mu_2 &= n \{ \Sigma (\theta_r - \theta_s)^2 p_r p_s - (\Sigma | \theta_r - \theta_s | p_r p_s)^2 \} \\ &= n w_1 \end{aligned} \quad (6.3)$$

TABLE III. Powers of different tests for various alternatives

In each cell the first figure represents the power of the test  $\Sigma(x_i - y_i)$ , the second  $\Sigma |x_i - y_i|$  and the third  $\Sigma(x_i - y_i)^2$

| $\delta_1 \backslash \delta_2$ | -.175                | -.150                | -.125                | -.100                | -.075                | -.050                | -.025                | 0                    | .025                 | .050                 | .075                 | .100                 | .125                 | .150                 | .175                 |
|--------------------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|----------------------|
| -.175                          | .050<br>.050<br>.050 | .053<br>.085<br>.055 | .056<br>.082<br>.066 | .060<br>.071<br>.083 | .064<br>.081<br>.107 | .069<br>.095<br>.138 | .074<br>.113<br>.179 | .079<br>.136<br>.229 | .085<br>.165<br>.289 | .091<br>.202<br>.358 | .099<br>.247<br>.435 | .106<br>.300<br>.517 | .115<br>.362<br>.601 | .125<br>.432<br>.683 | .136<br>.508<br>.759 |
| -.150                          | .047<br>.045<br>.047 | .050<br>.050<br>.050 | .053<br>.056<br>.057 | .057<br>.063<br>.068 | .061<br>.072<br>.087 | .065<br>.084<br>.113 | .070<br>.100<br>.147 | .075<br>.121<br>.192 | .081<br>.149<br>.247 | .087<br>.183<br>.312 | .094<br>.226<br>.387 | .102<br>.278<br>.469 | .110<br>.340<br>.555 | .120<br>.410<br>.641 | .130<br>.488<br>.723 |
| -.125                          | .044<br>.042<br>.048 | .047<br>.045<br>.046 | .050<br>.050<br>.050 | .053<br>.056<br>.057 | .057<br>.063<br>.070 | .061<br>.073<br>.090 | .066<br>.087<br>.119 | .071<br>.106<br>.157 | .077<br>.131<br>.206 | .083<br>.163<br>.266 | .090<br>.204<br>.337 | .097<br>.254<br>.418 | .105<br>.315<br>.505 | .115<br>.385<br>.595 | .125<br>.463<br>.683 |
| -.100                          | .041<br>.040<br>.051 | .044<br>.042<br>.047 | .047<br>.046<br>.050 | .050<br>.050<br>.050 | .054<br>.054<br>.058 | .058<br>.064<br>.072 | .062<br>.075<br>.094 | .067<br>.091<br>.125 | .072<br>.113<br>.167 | .078<br>.142<br>.221 | .085<br>.180<br>.288 | .092<br>.228<br>.365 | .100<br>.286<br>.452 | .109<br>.356<br>.545 | .120<br>.434<br>.638 |
| -.075                          | .038<br>.039<br>.038 | .040<br>.041<br>.052 | .043<br>.043<br>.047 | .047<br>.046<br>.047 | .050<br>.050<br>.050 | .054<br>.056<br>.059 | .058<br>.065<br>.074 | .063<br>.078<br>.098 | .068<br>.096<br>.133 | .074<br>.121<br>.179 | .080<br>.155<br>.239 | .087<br>.200<br>.312 | .095<br>.255<br>.397 | .104<br>.322<br>.491 | .114<br>.400<br>.588 |
| -.050                          | .035<br>.041<br>.039 | .037<br>.042<br>.059 | .040<br>.043<br>.048 | .043<br>.044<br>.048 | .046<br>.046<br>.047 | .050<br>.050<br>.050 | .054<br>.056<br>.059 | .059<br>.066<br>.076 | .064<br>.080<br>.103 | .069<br>.101<br>.141 | .075<br>.131<br>.193 | .082<br>.171<br>.260 | .090<br>.222<br>.341 | .098<br>.286<br>.433 | .108<br>.362<br>.534 |
| -.025                          | .032<br>.045<br>.083 | .034<br>.045<br>.071 | .037<br>.045<br>.061 | .039<br>.045<br>.053 | .043<br>.046<br>.048 | .046<br>.047<br>.046 | .050<br>.050<br>.050 | .054<br>.056<br>.060 | .059<br>.067<br>.079 | .064<br>.083<br>.109 | .070<br>.108<br>.152 | .077<br>.142<br>.210 | .084<br>.188<br>.284 | .093<br>.247<br>.373 | .102<br>.320<br>.474 |



|      |      |      |      |      |      |      |      |      |      |      |      |      |      |       |      |
|------|------|------|------|------|------|------|------|------|------|------|------|------|------|-------|------|
| 0    | .029 | .031 | .033 | .036 | .039 | .042 | .046 | .050 | .055 | .060 | .065 | .072 | .079 | .087  | .096 |
|      | .052 | .052 | .051 | .049 | .048 | .047 | .048 | .050 | .050 | .050 | .050 | .050 | .050 | .050  | .050 |
|      | .102 | .087 | .078 | .062 | .053 | .048 | .046 | .046 | .061 | .082 | .115 | .164 | .229 | .313  | .411 |
| .025 | .026 | .028 | .030 | .033 | .035 | .038 | .042 | .046 | .050 | .055 | .060 | .066 | .073 | .081  | .089 |
|      | .062 | .062 | .060 | .057 | .054 | .051 | .049 | .048 | .050 | .057 | .070 | .091 | .124 | .170  | .231 |
|      | .126 | .107 | .090 | .076 | .063 | .054 | .048 | .046 | .050 | .062 | .086 | .123 | .178 | .252  | .346 |
| .050 | .023 | .025 | .027 | .029 | .032 | .035 | .038 | .041 | .045 | .050 | .055 | .061 | .067 | .075  | .083 |
|      | .077 | .076 | .074 | .070 | .065 | .060 | .054 | .050 | .048 | .050 | .057 | .072 | .096 | .134  | .187 |
|      | .157 | .134 | .113 | .094 | .078 | .065 | .054 | .048 | .045 | .050 | .064 | .090 | .132 | .195  | .280 |
| .075 | .020 | .022 | .024 | .026 | .028 | .031 | .034 | .037 | .041 | .045 | .050 | .055 | .062 | .069  | .077 |
|      | .098 | .097 | .093 | .088 | .082 | .074 | .066 | .058 | .052 | .049 | .050 | .058 | .074 | .102  | .145 |
|      | .194 | .167 | .141 | .121 | .098 | .080 | .066 | .055 | .047 | .045 | .050 | .065 | .095 | .144  | .216 |
| .100 | .017 | .019 | .021 | .023 | .025 | .027 | .030 | .033 | .037 | .040 | .045 | .050 | .056 | .062  | .070 |
|      | .126 | .125 | .121 | .114 | .106 | .095 | .084 | .073 | .062 | .054 | .049 | .050 | .058 | .077  | .109 |
|      | .241 | .208 | .178 | .150 | .124 | .102 | .083 | .067 | .055 | .047 | .045 | .050 | .057 | .101  | .158 |
| .125 | .015 | .016 | .018 | .019 | .021 | .024 | .026 | .029 | .032 | .036 | .040 | .045 | .050 | .056  | .063 |
|      | .164 | .163 | .158 | .149 | .139 | .125 | .111 | .096 | .081 | .067 | .056 | .050 | .050 | .059* | .080 |
|      | .300 | .261 | .224 | .190 | .159 | .131 | .107 | .086 | .069 | .056 | .047 | .044 | .050 | .070  | .109 |
| .150 | .012 | .014 | .015 | .016 | .018 | .020 | .022 | .025 | .028 | .031 | .035 | .039 | .044 | .050  | .057 |
|      | .214 | .212 | .206 | .196 | .183 | .167 | .148 | .128 | .108 | .089 | .072 | .059 | .051 | .050  | .060 |
|      | .372 | .327 | .283 | .242 | .204 | .170 | .139 | .113 | .090 | .072 | .057 | .047 | .044 | .050  | .073 |
| .175 | .010 | .011 | .012 | .014 | .015 | .017 | .019 | .021 | .024 | .027 | .030 | .034 | .039 | .044  | .050 |
|      | .279 | .276 | .269 | .257 | .241 | .222 | .199 | .174 | .148 | .123 | .099 | .078 | .061 | .051  | .050 |
|      | .456 | .409 | .358 | .310 | .264 | .221 | .183 | .149 | .119 | .095 | .074 | .059 | .048 | .044  | .050 |

SOME DISTRIBUTIONS ARISING IN MATCHING PROBLEMS

TABLE IV. *Region of bias and unbiasedness for different tests*

| Hypothesis $\delta_1$<br>(deviation from<br>.20) | Tests   | Region of unbiasedness $W_1$   | Region of bias $W_2$   | Remarks   |
|--|---|--|--|---|
| -.175 ..   | $\Sigma(x_r - y_r)$<br>$\Sigma x_r - y_r $<br>$\Sigma(x_r - y_r)^2$ | $W_1 \leq -.175$<br>$W_1 \leq -.175 \ \& \ \geq 0$<br>$W_1 \leq -.175 \ \& \ \geq -.100$     | $W_2 > -.175$<br>$-.175 < W_2 < 0$<br>$-.175 < W_2 < -.100$    | Most powerful   |
| -.150 ..   | $\Sigma(x_r - y_r)$<br>$\Sigma x_r - y_r $<br>$\Sigma(x_r - y_r)^2$ | $W_1 \leq -.150$<br>$W_1 \leq -.150 \ \& \ \geq 0$<br>$W_1 \leq -.150 \ \& \ \geq -.075$     | $W_2 > -.150$<br>$-.150 < W_2 < 0$<br>$-.150 < W_2 < -.075$    | Most powerful   |
| -.125 ..   | $\Sigma(x_r - y_r)$<br>$\Sigma x_r - y_r $<br>$\Sigma(x_r - y_r)^2$ | $W_1 \leq -.125$<br>$W_1 \leq -.125 \ \& \ \geq 0$<br>$W_1 \leq -.125 \ \& \ \geq -.050$     | $W_2 > -.125$<br>$-.125 < W_2 < 0$<br>$-.125 < W_2 < -.050$    | Most powerful   |
| -.100 ..   | $\Sigma(x_r - y_r)$<br>$\Sigma x_r - y_r $<br>$\Sigma(x_r - y_r)^2$ | $W_1 \leq -.100$<br>$W_1 \leq -.100 \ \& \ \geq -.025$<br>$W_1 \leq -.100 \ \& \ \geq -.025$ | $W_2 > -.100$<br>$-.100 < W_2 < .025$<br>$-.100 < W_2 < -.025$ | Most powerful   |
| -.075 ..   | $\Sigma(x_r - y_r)$<br>$\Sigma x_r - y_r $<br>$\Sigma(x_r - y_r)^2$ | $W_1 \leq -.075$<br>$W_1 \leq -.075 \ \& \ \geq .025$<br>$W_1 \leq -.075 \ \& \ \geq 0$      | $W_2 > -.075$<br>$-.075 < W_2 < .025$<br>$-.075 < W_2 < 0$     | Most powerful   |
| -.050 ..   | $\Sigma(x_r - y_r)$<br>$\Sigma x_r - y_r $<br>$\Sigma(x_r - y_r)^2$ | $W_1 \leq -.050$<br>$W_1 \leq -.050 \ \& \ \geq .025$<br>$W_1 \leq -.050 \ \& \ \geq .025$   | $W_2 > -.050$<br>$-.050 < W_2 < .025$<br>$-.050 < W_2 < .025$  | Most powerful excepting a<br>small region near $\delta_1 = -.050$   |
| -.025 .  | $\Sigma(x_r - y_r)$<br>$\Sigma x_r - y_r $<br>$\Sigma(x_r - y_r)^2$ | $W_1 \leq -.025$<br>$W_1 \leq -.025 \ \& \ \geq .050$<br>$W_1 \leq -.025 \ \& \ \geq .050$   | $W_2 > -.025$<br>$-.025 < W_2 < .050$<br>$-.025 < W_2 < .050$  | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |

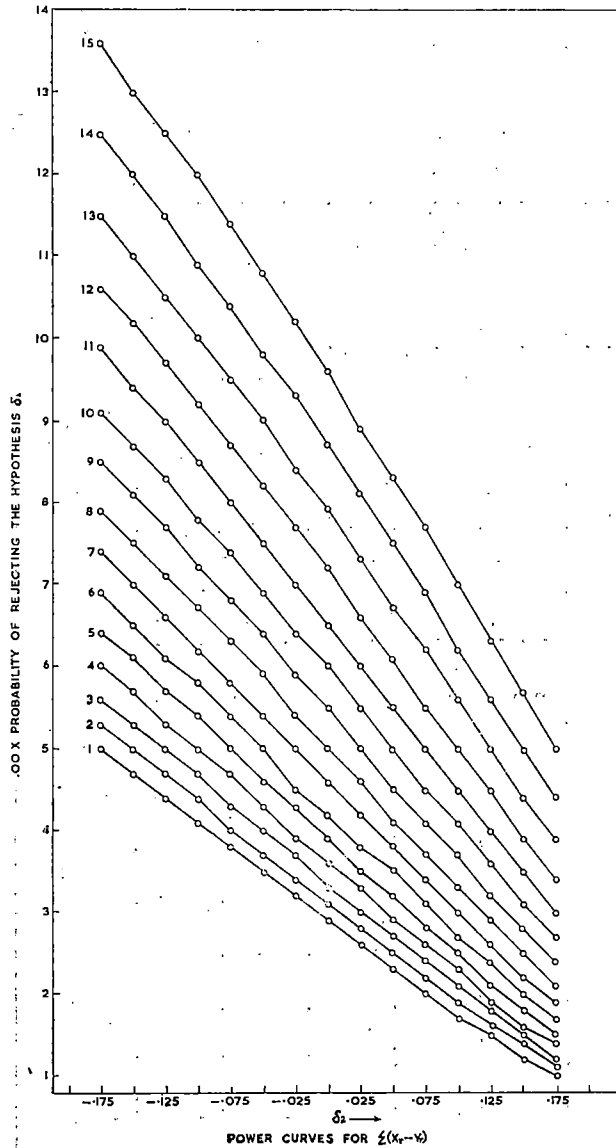
|      |    |   |   |   |   |
|------|----|---|---|---|---|
| 0    | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | $W_1 \leq 0$<br>$W_1 \leq 0$ & $\geq .050$<br>$W_1 \leq 0$ & $\geq .075$          | $W_2 > 0$<br>$0 < W_2 < .050$<br>$0 < W_2 < .075$           | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |
| .025 | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | $W_1 \leq .025$<br>$W_1 \leq .025$ & $\geq .075$<br>$W_1 \leq .025$ & $\geq .100$ | $W_2 > -.025$<br>$.025 < W_2 < .075$<br>$.025 < W_2 < .100$ | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |
| .050 | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | $W_1 \leq .050$<br>$W_1 \leq .050$ & $\geq .100$<br>$W_1 \leq .050$ & $\geq .125$ | $W_2 > .050$<br>$.050 < W_2 < .100$<br>$.050 < W_2 < .125$  | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |
| .075 | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | $W_1 \leq .075$<br>$W_1 \leq .075$ & $\geq .125$<br>$W_1 \leq .075$ & $\geq .150$ | $W_2 > .075$<br>$.075 < W_2 < .125$<br>$.075 < W_2 < .150$  | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |
| .100 | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | $W_1 \leq .100$<br>Unbiased for all alternatives<br>$W_1 \leq .100$ & $\geq .150$ | $W_2 > .100$<br>$.100 < W_2 < .150$                         | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |
| .125 | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | $W_1 \leq .125$<br>Unbiased for all alternatives<br>$W_1 \leq .125$               | $W_2 > .125$<br>$W_2 > .125$                                | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |
| .150 | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | $W_1 \leq .150$<br>Unbiased for all alternatives<br>$W_1 \leq .150$               | $W_2 > .150$<br>$W_2 > .150$                                | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |
| .175 | .. | $\frac{\sum(x_r - y_r)}{\sum(x_r - y_r)^2}$ | Unbiased for all alternatives less than .175<br>do<br>do                          |   | Most powerful for $\delta_2 > \delta_1$<br>do $\delta_2 < \delta_1$ |

Power was not actually calculated beyond the limits  $\delta = -.175$  to  $\delta = .175$ . But the regions of bias and unbiasedness have been decided from the trend of the power curves.

Alternate hypothesis  $H_1$

$$\left. \begin{aligned} \mu_1' &= nz_1 + n\psi_1(d_1, d_2, \dots, d_k) = nz_2 \\ \mu_2 &= nw_1 + n\psi_2(d_1, d_2, \dots, d_k) = nw_2 \end{aligned} \right\} \quad (6.4)$$

where  $\psi_1$  and  $\psi_2$  are functions of  $d_1, d_2, \dots, d_k$ .



GRAPH 1

Choose the tail region of rejection such that

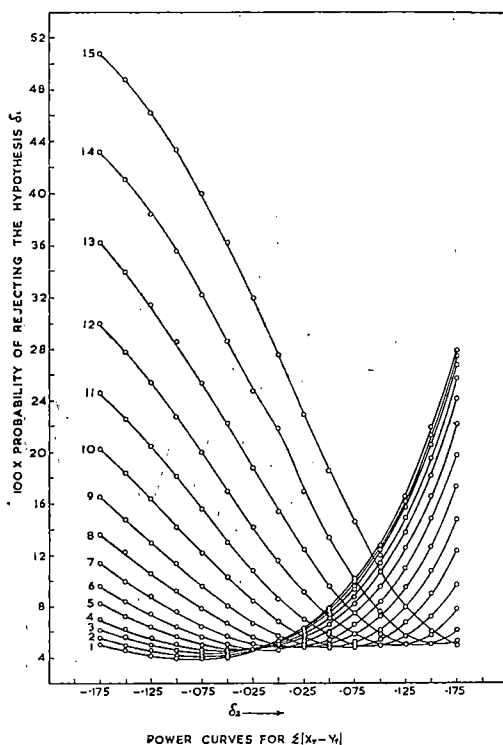
$$\frac{\sum |x_r - y_r|}{n} - z_1 \geq t_n \sqrt{\frac{w_1}{n}},$$

where  $t_n \rightarrow t$  as  $n$  tends to infinity. Then

$$P \left[ \frac{\sum |x_r - y_r|}{n} - z_1 \geq t_n \sqrt{\frac{w_1}{n}} \right] < \frac{1}{t_n^2} \quad (6.5)$$

Now the chance of accepting  $H_0$  when  $H_1$  is true is given by

$$\begin{aligned} & P \left\{ \frac{\sum |x_r - y_r|}{n} - z_1 < t_n \sqrt{\frac{w_1}{n}} \mid H_1 \right\} \\ &= P \left\{ \frac{\sum |x_r - y_r|}{n} - z_2 \geq k \sqrt{\frac{w_2}{n}} \right\} < \frac{1}{k^2} \end{aligned} \quad (6.6)$$



GRAPH 2

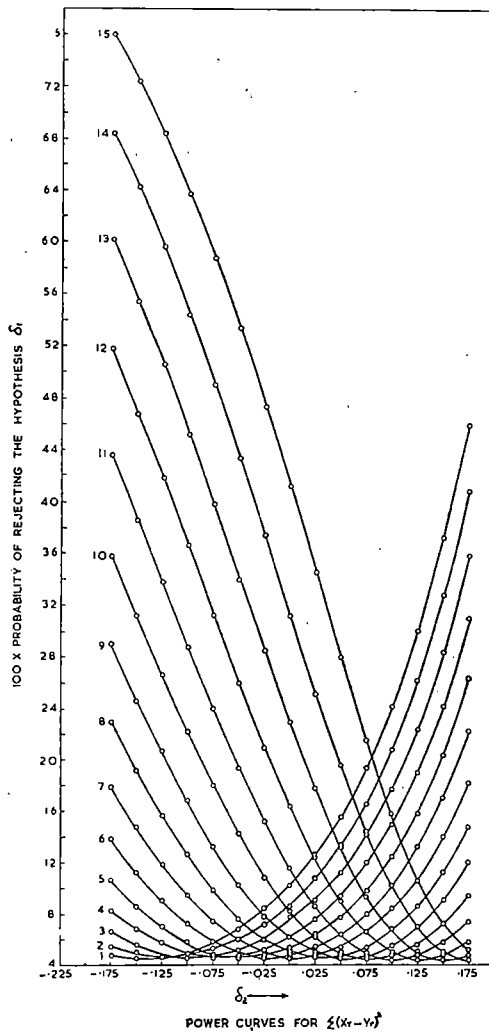
where

$$k = \left| \left[ t_n \sqrt{\frac{w_1}{n}} - \psi_1(d_1, d_2, \dots, d_k) \right] \right| / \sqrt{\frac{w_2}{n}},$$

provided  $n$  is large and  $\psi(d_1, d_2, \dots, d_k)$  is positive.

The probability of rejecting the null hypothesis is greater than

$$1 - \frac{w_2}{[t_n \sqrt{w_1} - \sqrt{n} \cdot \psi_1(d_1, d_2, \dots, d_k)]^2} \quad (6.7)$$



GRAPH 3

In the above Graphs curves 1, 2, 3, ..., 15 refer to hypotheses  $\delta_1 = -.175, -.150, -.125, \dots, +.175$  respectively.

If  $\psi_1$  is negative, we shall have to consider the other tail as the region of rejection.

We may now describe the actual procedure for testing two samples. Let two samples of sizes  $n$  be denoted by

Samples I.  $x_1, x_2, x_3, \dots, x_n$

II.  $y_1, y_2, y_3, \dots, y_n$

Let  $x_1, x_2, \dots, x_n; y_1, y_2, \dots, y_n$  assume any of the values  $\theta_1, \theta_2, \dots, \theta_k$  with probabilities  $p_1, p_2, \dots, p_k$ . The  $p$ 's are estimated by pooling the two samples together and finding the proportion of  $\theta_1, \theta_2, \dots, \theta_k$  in them. The expected values and variances of  $\Sigma(x_r - y_r)$ ,  $\Sigma|x_r - y_r|$  and  $\Sigma(x_r - y_r)^2$  can be calculated by using the results given in the previous sections. The observed values can be computed from the data. When  $n$  is not small the standardized deviate

$$z = \frac{x - m}{\sigma},$$

where  $x$  and  $m$  are the observed and the expected values of the *statistics*, and  $\sigma$  the standard deviation of the distribution, enables us to decide the significance of the difference between the two given samples. The probability that  $|z| \geq K$  is calculated on the assumption that the distribution of  $z$  is normal. As usual, if the probability that the observed  $|z|$  is less than or equal to 0.05, then we consider the two samples to be different from each other.

#### 7. POWER CURVES

We have already seen that of the three *statistics* considered in this paper,  $\Sigma|x_r - y_r|$  has the least coefficient of variation for the special cases examined and therefore  $\Sigma|x_r - y_r|$  may be preferred to the other *statistics* in actual practice. We shall examine this point in greater detail by finding the power curves of the three *statistics* for a particular value of  $n = 25$  for five values of  $\theta$ , i.e.,  $\theta_1 = 0, \theta_2 = 1, \theta_3 = 2, \theta_4 = 3$  and  $\theta_5 = 4$  with probabilities  $p_1 = p_2 = p_3 = .2$  and  $p_4 = .2 + \delta$  and  $p_5 = .2 - \delta$  for the hypothesis,  $\delta = \delta_1$  as compared to the alternate hypothesis  $\delta = \delta_2$  for varying values of  $\delta_1$  and  $\delta_2$  ranging from  $-0.175$  to  $.175$  for the first kind of errors 0.05. Tables III and IV, and graphs 1, 2 and 3 show the power of the different tests for various alternative hypotheses for the two equal tail regions of the normal curve for the probability 0.05. The unbiased region as determined from Table III is given in Table IV. It will be seen from this table that for the hypotheses  $\delta_1$ , ranging from  $-0.175$  to  $-0.050$ ,

$\Sigma(x_r - y_r)^2$  is the most powerful, while for  $\delta_1 > -0.50$ ,  $\Sigma|x_r - y_r|$  is the most powerful of the three tests when  $\delta_2 > \delta_1$ . When  $\delta_2 < \delta_1$ ,  $\Sigma(x_r - y_r)^2$  is the most powerful.  $\Sigma(x_r - y_r)$  is uniformly the least powerful of the three tests.

### 8. SUMMARY

This paper deals with a number of distributions arising from matching of two or more decks of cards of  $k$  characters which may be qualitative or quantitative. For two samples  $X$  and  $Y$  of size  $n$ , the sequence of observations being

$$X \dots x_1, x_2, x_3, \dots, x_n$$

$$Y \dots y_1, y_2, y_3, \dots, y_n$$

the distributions of  $\Sigma(x_r - y_r)$ ,  $\Sigma|x_r - y_r|$  and  $\Sigma(x_r - y_r)^2$  have been discussed for both finite and infinite sampling. All the cumulants are linear functions of the number of observations and therefore the distributions of the three *statistics* tend to the normal form as  $n$  tends to infinity. A table showing the probabilities for the distribution of  $\Sigma|x_r - y_r|$  for  $n = 25$ ,  $x$  and  $y$  taking any of the values 0, 1, 2, 3 and 4 with probabilities 0.2 has been given.

It has been shown that the three *statistics* are consistent and can be used for testing the significance of the difference between two given samples. Their powers in comparing two samples, each of 25 observations taking the values 0, 1, 2, 3 and 4 with equal probabilities, have been examined for different hypotheses and alternatives and it has been found that  $\Sigma|x_r - y_r|$  is the most powerful for certain regions, while for other regions  $\Sigma(x_r - y_r)^2$  is most powerful.

My sincere thanks are due to Prof. D. S. Kothari for suggesting to me to investigate the distribution of  $\Sigma|x_r - y_r|$  for two sets of samples taking the values 0, 1, 2, 3 and 4 with probabilities 0.2. I also wish to express my grateful thanks to Dr. F. C. Auluck, Messrs. S. P. Aggarwal and M. N. Bhattacharyya for helping me in preparing the tables presented in this paper.

### REFERENCES

- Anderson, T. W. .. "On Card Matching," *Ann. Math. Stat.*, 1943, **14**, 426.  
 Bartlett, M. S. .. "Properties of sufficiency and statistical tests," *Proc. Roy. Soc.*, 1937, **160 A**, 294.  
 Battin, I. L. .. "On the problem of multiple matching," *Ann. Math. Stat.*, 1942, **13**, 294.



- Greenwood, J. A. .. "The first four moments of a general matching problem," *Ann, Eug.*, 1940, **10**, 290.
- Greville, T. N. E. .. "The frequency distribution of a general matching problem," *Ann. Math. Stat.*, 1941, **12**, 350.
- Kaplansky, I. and Riordan, J. "Multiple matching and runs by the symbolic method," *ibid.*, 1945, **16**, 272.
- Mann, H. B. and Whitney, D. R. "On a test whether one of two random variables is stochastically larger than the other," *ibid.*, 1947, **18**, 50.
- Wald, A. and Wolfowitz, J. "On a test whether two samples are from the same population," *ibid.*, 1940, **11**, 147.
- Wilks, S. S. .. *Mathematical Statistics*, Princeton University Press, 1946.