



Poisson Approximation for Some Statistics Based on Exchangeable Trials

A. D. Barbour; G. K. Eagleson

Advances in Applied Probability, Vol. 15, No. 3. (Sep., 1983), pp. 585-600.

Stable URL:

<http://links.jstor.org/sici?sici=0001-8678%28198309%2915%3A3%3C585%3APAFSSB%3E2.0.CO%3B2-5>

Advances in Applied Probability is currently published by Applied Probability Trust.

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Please contact the publisher regarding any further use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/apt.html>.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

JSTOR is an independent not-for-profit organization dedicated to creating and preserving a digital archive of scholarly journals. For more information regarding JSTOR, please contact support@jstor.org.

POISSON APPROXIMATION FOR SOME STATISTICS BASED ON EXCHANGEABLE TRIALS

A. D. BARBOUR,* *Gonville and Caius College, Cambridge*

G. K. EAGLESON,** *CSIRO Division of Mathematics and Statistics*

Abstract

Stein's (1970) method of proving limit theorems for sums of dependent random variables is used to derive Poisson approximations for a class of statistics, constructed from finitely exchangeable random variables.

Let $\{Y_i\}_{i=1}^M$ be exchangeable random elements of a space \mathcal{X} and, for I a k -subset of \mathcal{X} , let X_I be a 0–1 function. The statistics studied here are of the form

$$W = \sum_{I \in \mathcal{N}} X_I,$$

where \mathcal{N} is some collection of k -subsets of \mathcal{X} .

An estimate of the total variation distance between the distributions of W and an appropriate Poisson random variable is derived and is used to give conditions sufficient for W to be asymptotically Poisson. Two applications of these results are presented.

FINITELY EXCHANGEABLE RANDOM VARIABLES

1. Introduction

Let $\{Y_i\}_{i=1}^M$ be exchangeable random elements of a space \mathcal{X} , and let ϕ be a symmetric mapping from \mathcal{X}^k to $\{0, 1\}$. Let $D(k; r, s)$ denote the collection of k -subsets of $\{r, r+1, \dots, s\}$ and, for $I \in D(k; r, s)$, let $Y_I = \{Y_i : i \in I\}$ and $X_I = \phi(Y_I)$, where \mathbf{i} is any ordering of I : X_I is well defined because of the symmetry of ϕ .

Many statistics can be expressed in the form

$$W = \sum_{I \in \mathcal{N}} X_I,$$

where \mathcal{N} is some subset of $D \equiv D(k; 1, M)$, and we shall be concerned with studying their distribution.

Received 8 June 1982; revision received 5 January 1983.

* Present address: Institut für Angewandte Mathematik, Universität Zürich, Rämistrasse 74, CH-8001 Zürich, Switzerland.

** Postal address: CSIRO Division of Mathematics and Statistics, Bradfield Rd, West Lindfield, NSW 2070, Australia.

When the underlying $\{Y_j\}$ are independent and identically distributed and $\mathcal{N} = D$, W is asymptotically distributed (as $M \rightarrow \infty$) under rather natural conditions as a Poisson random variable (Silverman and Brown (1978)). In this case, W has been used to construct tests of randomness and of collinearities of points in the plane. Eagleson (1979) assumed only exchangeability, and took $\mathcal{N} = D(k; 1, r(M)) \subseteq D$; he showed that W is asymptotically distributed as a Poisson random variable (as $M \rightarrow \infty$) provided that some moment conditions are satisfied and that

$$\lim_{M \rightarrow \infty} r(M)/M = 0.$$

It seems natural to expect that similar results hold when \mathcal{N} is chosen to be some other relatively small subset of D . Such an extension is needed to provide a proof of the asymptotic distribution of Knox’s statistic, used in epidemiology, where the subset \mathcal{N} is chosen as those outbreaks of a disease which occur close in time (see Section 3 for more details). Unfortunately, the martingale techniques used in Eagleson (1979) do not readily give the necessary generalization.

An alternative approach to proving limit theorems for sums of dependent random variables was originated by Stein (1970) and developed in the Poisson context by Chen (1975). The basic idea is that, given a set $A \subset Z^+$, one defines a bounded function $x = x(\lambda, A): Z^+ \rightarrow R$ by

$$\begin{aligned} x(0) &= 0; \\ x(m+1) &= \lambda^{-m-1} e^\lambda m! [\mathcal{P}_\lambda(A \cap U_m) - \mathcal{P}_\lambda(A)\mathcal{P}_\lambda(U_m)], \quad m \geq 0, \end{aligned}$$

where $U_m = \{0, 1, \dots, m\}$ and where $\mathcal{P}_\lambda(B)$ denotes the probability that a Poisson random variable with mean λ belongs to B . The importance of such an x is that, for any non-negative integer-valued random variable T , it can be shown that

$$(1) \quad E\{\lambda x(T+1) - Tx(T)\} = P(T \in A) - \mathcal{P}_\lambda(A).$$

Chen used this identity to derive an estimate of the total variation distance between T and \mathcal{P}_λ when T is a sum of dependent Bernoulli random variables, satisfying a mixing condition. In the present paper, the symmetry properties of the joint distribution of $\{X_I; I \in D\}$ are used to calculate some conditional expectations explicitly and thus to approximate the left-hand side of (1).

It cannot be denied that the following analysis is far from simple, but the wide applicability of the results obtained and their potential importance justifies the detailed calculations. The main results are stated in Section 2, and their proofs given in Section 3. In Section 4 two applications are described; one a justification of the use of Knox’s statistic and the other a study of the asymptotic distribution of the number of ‘rare’ patterns occurring in a random arrangement of two types of points.

2. Statement of results

The principal results are expressed in terms of convergence theorems for a sequence of statistics $W_n, n \geq 1$, each of which is derived in the manner described in the introduction. Thus, for example, the underlying probabilistic framework is that of a triangular array $(\{Y_{jn}\}_{j=1}^{M(n)}; n \geq 1)$ of random variables exchangeable within rows. The value of k is kept the same for all n , whereas all other quantities are allowed to vary: however, to simplify notation, the explicit dependence of the other quantities upon n is suppressed.

With notation as before, let

$$p = EX_I, \quad N = |\mathcal{N}| \quad \text{and} \quad \mu = EW = Np.$$

Then, for $0 \leq j \leq k$, let C_j be the number of pairs of elements (I, J) of \mathcal{N} such that $|I \cap J| = j$: further, for such pairs (I, J) , set $q_j = EX_I X_J$ and $\eta_j = |q_j - p^2|^{\frac{1}{2}}$. Finally, for any subset $K \subseteq \{1, 2, \dots, M\}$, let $\text{In}(K)$ consist of those elements of \mathcal{N} which have non-empty intersection with K , set $i_K = |\text{In}(K)|$, and let $\text{No}(K) = \mathcal{N} \setminus \text{In}(K)$.

Theorem 1. The distribution of W converges (as $n \rightarrow \infty$) in distribution to \mathcal{P}_λ , the Poisson distribution with parameter λ , if the following six conditions are satisfied:

- (i) $Np \rightarrow \lambda$
- (ii) $N^{-1} \max_{J \in \mathcal{N}} i_J \rightarrow 0$
- (iii) $NM^{-k} \rightarrow 0$
- (iv) $C_j q_j \rightarrow 0, \quad 1 \leq j \leq k-1$
- (v) $N^2 q_j M^{-j} \rightarrow 0, \quad 1 \leq j \leq k-1$
- (vi) $N\eta_0 \rightarrow 0.$

The conditions of Theorem 1 simplify considerably in various special cases. For example, if $\mathcal{N} = D$ and each row of the Y 's is the same independent and identically distributed sequence, they can be reduced to those of Silverman and Brown (1978), Theorem A. In fact, Corollary 4 below exhibits an improvement in this case on the rate of convergence to the Poisson limit given in Brown and Silverman (1979). When $\mathcal{N} = D(k; 1, r(M))$, the theorem reduces to that of Eagleson (1979), and when \mathcal{N} is a balanced subset and the Y 's are independent and identically distributed a theorem of Berman and Eagleson (1983) is recovered.

The following further cases are useful in applications.

Corollary 1 (subdiagonal). Let $M = n$ and $\mathcal{N} = \{\{r, r+1, \dots, r+k-1\}\}_{r=1}^{M-k+1}$. Then $W \rightarrow^D \mathcal{P}_\lambda$ if, as $M \rightarrow \infty, Mp \rightarrow \lambda, Mq_j \rightarrow 0, 1 \leq j \leq k-1$, and $M\eta_0 \rightarrow 0$.

Proof. Observe that here $\max_{J \in \mathcal{N}} i_J \leq 2k-1$ and that, for $1 \leq j \leq k-1, C_j \sim 2M$. Now apply Theorem 1.

converges to 0 if, as $M \rightarrow \infty$, $\lambda = Np \rightarrow \infty$, $mp \rightarrow 0$, $mM^{-1}\lambda^2 \rightarrow 0$, $\lambda^{-1}q_1\{C_1 \vee \lambda N^2/M\} \rightarrow 0$ and $\lambda^{-\frac{1}{2}}N\eta_0 \rightarrow 0$.

Proof. Easily, $N \leq Mm$, $N > \frac{1}{2}m^2$ and, for $J \in \mathcal{N}$, $i_J \leq 4m$. Hence, under the conditions of Part (i) of the corollary, by an argument similar to that in the proof of Corollary 2, Conditions (ii) and (iii) of Theorem 1 are satisfied while Conditions (i), (iv), (v) and (vi) are immediate.

Part (ii) of the corollary is proved in similar fashion, by appealing to Theorem 2.

3. Proofs

The proofs are based on the following three lemmas. Let $S \in D(s; 1, M)$, and, for $J \in \mathcal{N}$, define

$$W_{S,J} = \sum_{I \in \text{No}(S \cup J)} X_I.$$

Then, for each bounded $x : Z^+ \rightarrow R$, let

$$\|x\| = \sup_{m \in Z^+} |x(m)|,$$

and let

$$\Delta x = \sup_{m \in Z^+} |x(m+1) - x(m)|.$$

Lemma 1. For each bounded $x : Z^+ \rightarrow R$ and each $S \in D(s; 1, M)$,

$$\begin{aligned} & \left| E\{\mu x(W+1) - Wx(W)\} - \sum_{I \in \text{No}(S)} E\{(p - X_I)x(W_{S,I} + 1)\} \right| \\ & \leq (\Delta x \wedge 2 \|x\|) \left\{ p^2 \left(Ni_S + \sum_{I \in \mathcal{N}} i_I \right) + \sum_{j=1}^{k-1} q_j C_j + Nq_0 i_s \right\} + 2p \|x\| i_s. \end{aligned}$$

Proof. Consider the expression

$$\begin{aligned} & E\{\mu x(W+1) - Wx(W)\} - \sum_{I \in \text{No}(S)} E\{(p - X_I)x(W_{S,I} + 1)\} \\ & = \sum_{I \in \text{No}(S)} E\{p[x(W+1) - x(W_{S,I} + 1)] + X_I[x(W_{S,I} + 1) - x(W)]\} \\ & \quad + \sum_{I \in \text{In}(S)} E\{px(W+1) - X_I x(W)\}. \end{aligned}$$

For the first term on the right-hand side, one has the estimates

$$|x(W+1) - x(W_{S,I} + 1)| \leq 2 \|x\| I[x(W+1) - x(W_{S,I} + 1) \neq 0] \leq 2 \|x\| \left\{ \sum_{J \in \text{In}(S \cup I)} X_J \right\}$$

and

$$|x(W + 1) - x(W_{S,I} + 1)| \leq \Delta x \left\{ \sum_{J \in \text{In}(S \cup I)} X_J \right\};$$

hence

$$\begin{aligned} \left| \sum_{I \in \text{No}(S)} p E\{x(W + 1) - x(W_{S,I} + 1)\} \right| &\leq (2 \|x\| \wedge \Delta x) \sum_{I \in \text{No}(S)} i_{S \cup I} p^2 \\ &\leq (2 \|x\| \wedge \Delta x) p^2 \left\{ N i_S + \sum_{I \in \mathcal{N}} i_I \right\}. \end{aligned}$$

For the next term, in a similar way,

$$\begin{aligned} |X_I \{x(W_{S,I} + 1) - x(W)\}| &\leq X_I (2 \|x\| \wedge \Delta x) \sum_{J \in \text{In}(S \cup I) \setminus I} X_J \\ &\leq (2 \|x\| \wedge \Delta x) X_I \left\{ \sum_{J \in \text{In}(I) \setminus I} X_J + \sum_{J \in \text{In}(S \cup I) \setminus \text{In}(I)} X_J \right\}. \end{aligned}$$

Adding over $I \in \text{No}(S)$ and taking expectations gives

$$\begin{aligned} \left| \sum_{I \in \text{No}(S)} E\{X_I [x(W_{S,I} + 1) - x(W)]\} \right| &\leq (2 \|x\| \wedge \Delta x) \left\{ \sum_{\substack{I \neq J \in \mathcal{N} \\ I \cap J \neq \emptyset}} E(X_I X_J) + \sum_{\substack{I \in \text{No}(S) \\ J \in \text{In}(S \cap I) \setminus \text{In}(I)}} E(X_I X_J) \right\} \\ &\leq \left\{ \sum_{j=1}^{k-1} q_j C_j + N q_0 i_S \right\} (2 \|x\| \wedge \Delta x). \end{aligned}$$

The last term is easily seen to be no larger than $2p \|x\| i_S$.

Remark. It is always possible to choose S in such a way that $i_S \leq ksN/M$. This is because

$$\sum_{S \in \mathcal{D}(s; 1, M)} i_S = \sum_{I \in \mathcal{N}} \left\{ \binom{M}{s} - \binom{M-k}{s} \right\},$$

and hence, for some S among the $\binom{M}{s}$ possibilities, we must have

$$i_S \leq N \left\{ 1 - \binom{M-k}{s} / \binom{M}{s} \right\} \leq N \left\{ 1 - \left(1 - \frac{k}{M} \right)^s \right\} \leq Nks/M.$$

Lemma 2.

$$\left| \sum_{I \in \text{No}(S)} E\{(X_I - p)x(W_{S,I} + 1)\} \right| \leq \|x\| (k!)^{\frac{1}{2}} N \sum_{j=0}^k \left\{ \binom{k}{j} / (k-j)! \right\}^{\frac{1}{2}} \eta_j s^{-\frac{1}{2}j}.$$

Proof. The σ -fields

$$\mathcal{Y}_{S,I} = \sigma\{Y_j; j \in I \cup S\}$$

and

$$\mathcal{L}_{S,I} = \sigma\{Y_j; j \in \{1, \dots, M\} \setminus (I \cup S)\}$$

are generated by disjoint subsets of the random variables $\{Y_{ij}^M; X_I \in \mathcal{Y}_{S,I}$ and $W_{S,I} \in \mathcal{L}_{S,I}$. Hence

$$\begin{aligned} E\{(X_I - p)x(W_{S,I} + 1)\} &= E\{E\{(X_I - p)x(W_{S,I} + 1) \mid \mathcal{L}_{S,I}\}\} \\ &= E\{x(W_{S,I} + 1)E\{(X_I - p) \mid \mathcal{L}_{S,I}\}\} \\ &= E\left\{x(W_{S,I} + 1) \sum_{J \in D_{S \cup I}} \binom{s+k}{k}^{-1} E\{(X_J - p) \mid \mathcal{L}_{S,I}\}\right\}, \end{aligned}$$

where D_K denotes the set of k -subsets of K , the last equality because the Y 's are exchangeable. Thus

$$\begin{aligned} |E\{(X_I - p)x(W_{S,I} + 1)\}| &= \left| E\left\{x(W_{S,I} + 1) \sum_{J \in D_{S \cup I}} \binom{s+k}{k}^{-1} (X_J - p)\right\} \right| \\ &\leq \left\{ \|x\| / \binom{s+k}{k} \right\} E\left[\left\{ \sum_{J \in D_{S \cup I}} (X_J - p) \right\}^2 \right]^{\frac{1}{2}}. \end{aligned}$$

Direct computation now shows that for each $I \in \text{No}(S)$

$$E\left\{ \sum_{J \in D_{S \cup I}} (X_J - p) \right\}^2 \leq \binom{s+k}{k} \sum_{j=0}^k \binom{k}{j} \binom{s}{k-j} \eta_j^2,$$

and hence the lemma follows.

Lemma 3. For any $A \subset Z^+$ and $\lambda > 0$,

$$\begin{aligned} |P[W \in A] - \mathcal{P}_\lambda(A)| &\leq 2(1 \wedge 1 \cdot 4\lambda^{-\frac{1}{2}}) \left\{ \frac{1}{2}\varepsilon + p i_s + \sum_{j=0}^k N \eta_j s^{-j/2} \psi_{kj} \right\} \\ &\quad + (1 \wedge \lambda^{-1}) \left\{ p^2 \sum_{I \in \mathcal{N}} i_I + \sum_{j=1}^{k-1} q_j C_j + N(p^2 + q_0) i_s \right\}, \end{aligned}$$

where $\varepsilon = |\lambda - Np|$ and $\psi_{kj} = \frac{1}{2} \binom{k}{j} (j!)^{\frac{1}{2}}$.

Proof. Given $A \subset Z^+$, define $x = x(\lambda, A)$, $\mathcal{P}_\lambda(\cdot)$ and U_m as in the introduction. Recall that

$$E[\lambda x(W + 1) - Wx(W)] = P[W \in A] - \mathcal{P}_\lambda(A).$$

Now

$$|E\{(\lambda - Np)x(W + 1)\}| \leq \varepsilon \|x\|,$$

and, applying Lemmas 1 and 2,

$$\begin{aligned} |E\{Npx(W + 1) - Wx(W)\}| &\leq (2 \|x\| \wedge \Delta x) \left\{ p^2 \left(N i_s + \sum_{I \in \mathcal{N}} i_I \right) + \sum_{j=1}^{k-1} q_j C_j + N q_0 i_s \right\} \\ &\quad + 2 \|x\| \left\{ p i_s + \frac{1}{2} (k!)^{\frac{1}{2}} \sum_{j=0}^k \left\{ \binom{k}{j} / (k-j)! \right\}^{\frac{1}{2}} N \eta_j s^{-j/2} \right\}. \end{aligned}$$

The lemma now follows from the estimates for $\|x\|$ and Δx which are proved in the appendix.

Remark. Suppose that, as in Brown and Silverman (1979), $\mathcal{N} = D$ and the Y_j 's are independent and identically distributed. Then the random variables $\{Y_j\}_{j=1}^M$ are exchangeable amongst the larger collection $\{Y_j\}_{j=1}^\infty$, and so we may reformulate their problem by taking $M = \infty$ and $\mathcal{N} = D(k; 1, n)$. Now S may be chosen as large as we like from $\{n + 1, n + 2, \dots\}$ while keeping $i_S = 0$, and, because of independence of the Y_j 's, $\eta_0 = 0$. Hence the estimate in Lemma 3 is no larger than

$$\begin{aligned} (1 \wedge 1.4\lambda^{-\frac{1}{2}})\varepsilon + (1 \wedge \lambda^{-1}) \left\{ Np^2 \cdot Nk^2/n + \sum_{j=1}^{k-1} q_j N \binom{k}{j} \frac{n^{k-j}}{(k-j)!} \right\} \\ \leq (1 \wedge \lambda^{-1}) \left\{ \lambda^2 k^2/n + q_{k-1} N n^{k-1} \sum_{j=1}^{k-1} \binom{k}{j} \frac{n^{1-j}}{(k-j)!} \right\}, \end{aligned}$$

by Lemma 1 of Silverman and Brown (1978), if $\lambda = Np$. Setting $\rho = n^{2k}(q_{k-1} - p^2)$, this yields a bound on the total variation distance between W and \mathcal{P}_λ as follows.

Corollary 4. In the above circumstances, the total variation distance between W and \mathcal{P}_λ cannot exceed

$$(1 \wedge \lambda^{-1}) \{ n^{-1} \lambda^2 (k^2 + b_n) + \rho n^{-1} c_n / k! \},$$

where

$$\begin{aligned} c_n &= \sum_{j=1}^{k-1} \binom{k}{j} \frac{n^{1-j}}{(k-j)!} \rightarrow \frac{k}{(k-1)!}, \\ b_n &= \left\{ n^k / \binom{n}{k} \right\} c_n \rightarrow k^2, \end{aligned}$$

as $n \rightarrow \infty$.

This improves on the result of Brown and Silverman (1979) because $(\rho n^{-1})^{\frac{1}{2}}$ is replaced by ρn^{-1} in the convergence rate, because of the extra factor $1 \wedge \lambda^{-1}$, useful for large λ , and because all the constants are explicit. Of course, if better estimates of q_j , $1 \leq j \leq k - 2$, are known, there is a corresponding improvement in the total variation bound.

Proof of Theorem 1. Choose S so that $i_S \leq ksN/M$, and let $s \rightarrow \infty$ in such a way that $Ns^{-k} \rightarrow 0$ and $N^2 q_j s^{-j} \rightarrow 0$, $1 \leq j \leq k - 1$, whilst still ensuring that $s/M \rightarrow 0$; this is possible by (iii) and (v). Taking the estimate in Lemma 3 term

by term, it follows that:

$$\begin{aligned} \varepsilon &\rightarrow 0 \text{ by (i);} \\ p i_s &\leq ksNp/M \rightarrow 0 \text{ by (i) and the choice of } s; \\ \sum_{j=0}^k \psi_{kj} N \eta_j s^{-\frac{1}{2}j} &\rightarrow 0 \text{ by (i), (vi) and the choice of } s; \\ p^2 \sum_{i \in \mathcal{N}} i_i &\rightarrow 0 \text{ by (i) and (ii);} \\ \sum_{j=1}^{k-1} q_j C_j &\rightarrow 0 \text{ by (iv);} \end{aligned}$$

and

$$N(p^2 + q_0) i_s \leq ksM^{-1}(N^2 p^2 + N^2 q_0) \rightarrow 0 \text{ by (i), (vi) and the choice of } s.$$

Hence the distribution of W converges to \mathcal{P}_λ .

Remarks. Let $M = n$, $n \geq k$, and take $\mathcal{N} = \{(1, 2, \dots, k-1, r)\}_{r=k}^M$. Then, for each $I \in \mathcal{N}$, $i_I = N$, showing that Condition (ii) need not always be satisfied. It is clear that none of the other conditions is automatic. Nor is (iv) necessarily implied by (v). For if $M = n$, $n \geq k$, and, for any j , $0 \leq j < k$, $\mathcal{N} = \{(1, 2, \dots, j, j+r+1, \dots, k+r)\}_{r=0}^{M-k}$, then

$$C_{j'} = 0, \quad j' < j, \quad C_j \sim n^2, \quad \text{and} \quad C_{j'} \sim 2n, \quad j' > j.$$

Proof of Theorem 2. Argue as in the proof of Theorem 1, using the λ -dependent estimates in Lemma 3, and now letting $s \rightarrow \infty$ in such a way that $Ns^{-k} \rightarrow 0$ and $\lambda^{-1} N^2 \eta_j^2 s^{-j} \rightarrow 0$, $1 \leq j \leq k-1$, while still ensuring that $\lambda s M^{-1} \rightarrow 0$.

4. Applications

(i) *Knox's statistic.* In a study of childhood leukaemia, Knox (1964) suggested that a suitable test statistic for contagion would be the number of reported outbreaks of leukaemia which are close in time and close in space. He conjectured that, under the hypothesis of no association between time and position of occurrence, this statistic could be regarded as a Poisson random variable. In support of this conjecture, Barton and David (1966) gave a graph-theoretic proof that, under the above hypothesis, the mean and variance of this statistic are equal, and that a Poisson limit holds under certain assumptions.

Let $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ be n points on the plane and $0 \leq T_1 \leq T_2 \leq \dots \leq T_n$ be n time points. Knox's statistic is

$$K_n = \sum_{i < j} I(|\mathbf{Z}_i - \mathbf{Z}_j| < d_n) I(|T_i - T_j| < t_n),$$

where the measures of closeness, d_n and t_n , are chosen to make the mean of K_n small. We are interested in the distribution of K_n when the Z 's and T 's are fixed but associated at random.

Suppose that π is a random permutation of $(1, \dots, n)$ and consider

$$\begin{aligned} K_n &= \sum_{i < j} I(|Z_{\pi(i)} - Z_{\pi(j)}| < d_n) I(|T_i - T_j| < t_n) \\ &= \sum_{(i,j) \in \mathcal{N}} \phi(Z_{\pi(i)}, Z_{\pi(j)}), \end{aligned}$$

where $\mathcal{N} = \{(i, j) : |T_i - T_j| < t_n\}$. As the $\{Z_{\pi(i)}\}_{i=1}^n$ are exchangeable, K_n is an example of the sort of statistic studied in the previous section. In fact there is rather more structure in K_n , due to the fact that time is measured on a linear scale. If $|T_i - T_j| < t_n$ for some $j > i$ then both $|T_i - T_{j-k}| < t_n$ and $|T_{i+k} - T_j| < t_n$ for $k = 1, \dots, j - i - 1$. Thus K_n is a statistic of the form considered in Corollary 3 and hence, under the conditions listed there, is asymptotically distributed as a Poisson random variable.

Given a set of observations, one cannot tell whether it is one of a possible sequence of such observations satisfying the conditions of Corollary 3. However, Lemma 3 can be used to bound the total variation distance of the distribution of K_n from that of a Poisson random variable. Unfortunately, like most general bounds, Lemma 3 is not useful in a practical situation. Data similar to Knox's were simulated and the criteria d_n and t_n chosen to ensure that a 2×2 contingency table of close/not close distances against close/not close times had the same marginals as Knox's table. More precisely, 96 points independently and uniformly distributed on the unit square and 96 time points independently and uniformly distributed on the unit interval were chosen. The criteria d_n and t_n were chosen so that 25 of the interpoint distances were less than d_n and 152 of the time differences were less than t_n . The statistic K_n was then calculated from these data, as was the bound in Lemma 3 for the special case of s_n being taken as large as possible while leaving $i_s = 0$. In the particular simulation, K_n was 0 and the bound 0.71.

While the bound calculated from Lemma 3 was bad, 10000 repetitions of the above simulation gave an empirical distribution of K_n which, when fitted to the Poisson distribution with mean $25 \times 152/4,560 = 5/6$, produced a total variation distance of 0.016 and a Pearson χ^2 of 13.404 with 7 degrees of freedom.

In order to obtain a general result for the distribution of K_n , a model for generating the Z 's and T 's has to be assumed. If these random variables are all independent, the Z 's and T 's being associated at random, then, when a sufficient number of points are observed, most realisations give a K_n which is close to Poisson in a sense which is made precise by the following corollary.

Suppose that $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ are independent, identically distributed points on the plane with a common density which is bounded and continuous. Suppose also that T_1, \dots, T_n are independent and identically distributed points on the positive half-line. For fixed n , choose random variables δ_n and τ_n such that

$$\frac{1}{\binom{n}{2}} \sum_{i < j} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \delta_n) \leq \lambda/n$$

and

$$\sum_{i < j} I(|T_i - T_j| \leq \tau_n) \approx n.$$

Define

$$K_n = \sum_{i < j} I(|T_i - T_j| < \tau_n) I(|\mathbf{Z}_{\pi(i)} - \mathbf{Z}_{\pi(j)}| < \delta_n),$$

where π is a random permutation of $(1, \dots, n)$.

Corollary 5. With the statistic K_n defined as above, for all $\epsilon > 0$ there exists an n_0 such that for all $n \geq n_0$

$$P\left\{ \sup_{A \subset \mathbb{Z}^+} |P(K_n \in A | \mathbf{Z}_1, \dots, \mathbf{Z}_n; T_1, \dots, T_n) - \mathcal{P}_\lambda(A)| < \epsilon \right\} > 1 - \epsilon$$

for all $A \subset \mathbb{Z}^+$.

Proof. The corollary will be proved if it can be shown that the conditions of Corollary 3 hold ‘in probability’. We have

$$N = n \rightarrow \infty;$$

$$m \leq \sqrt{2n} = o(n) \text{ always};$$

$$Np \rightarrow \lambda, \text{ by the choice of } \delta_n \text{ and } \tau_n;$$

$$C_1 \leq Cn^3 \text{ always.}$$

Furthermore,

$$\begin{aligned} q_1 &\leq Cn^{-3} \sum_{i < j < k} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \delta_n) I(|\mathbf{Z}_i - \mathbf{Z}_k| \leq \delta_n) \\ &\leq Cn^{-3} \sum_{i < j < k} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \delta_n) \{I(|\mathbf{Z}_i - \mathbf{Z}_k| \leq \epsilon_n) + I(\delta_n \geq \epsilon_n)\} \\ &\leq Cn^{-3} \sum_{i < j < k} \{I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \epsilon_n) + I(\delta_n \geq \epsilon_n)\} I(|\mathbf{Z}_i - \mathbf{Z}_k| \leq \epsilon_n) \\ &\quad + C\lambda n^{-1} I(\delta_n \geq \epsilon_n). \end{aligned}$$

Now on the set $[\delta_n \geq \epsilon_n]$,

$$\frac{1}{\binom{n}{2}} \sum_{i < k} I(|\mathbf{Z}_i - \mathbf{Z}_k| \leq \epsilon_n) \leq \lambda/n,$$

and so

$$q_1 \leq Cn^{-3} \sum \sum I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \epsilon_n) I(|\mathbf{Z}_i - \mathbf{Z}_k| \leq \epsilon_n) + 2c\lambda n^{-1} I(\delta_n \geq \epsilon_n).$$

Hence

$$Eq_1(\mathbf{Z}_1, \dots, \mathbf{Z}_n; T_1, \dots, T_n) \leq c_1 n^{-1} P(\delta_n \geq \epsilon_n) + c_2 P(|\mathbf{Z}_1 - \mathbf{Z}_2| \leq \epsilon_n, |\mathbf{Z}_1 - \mathbf{Z}_3| \leq \epsilon_n).$$

Since the \mathbf{Z} 's have a bounded density,

$$P(|\mathbf{Z}_1 - \mathbf{Z}_2| \leq \epsilon_n, |\mathbf{Z}_1 - \mathbf{Z}_3| \leq \epsilon_n) = E(P(|\mathbf{Z}_1 - \mathbf{Z}_2| \leq \epsilon_n | \mathbf{Z}_1)^2) \leq 4\epsilon_n^2 L^2,$$

where L is a bound on the density of the \mathbf{Z} 's. Choose $\epsilon_n \sim c/n$ so that $P[|\mathbf{Z}_i - \mathbf{Z}_j| \leq \epsilon_n] \geq 2\lambda/n$, possible because the \mathbf{Z} 's have a continuous density.

Then note that

$$I(\delta_n \geq \epsilon_n) = I\left[\binom{n}{2}^{-1} \sum_{i < j} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \epsilon_n) \leq \lambda/n\right],$$

so that if

$$U = \binom{n}{2}^{-1} \sum_{i < j} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \epsilon_n),$$

then

$$P(\delta_n \geq \epsilon_n) = P(U \leq \lambda/n) \leq n^2 \lambda^{-2} \text{Var } U.$$

It is easy to see that

$$\begin{aligned} \text{Var } U &= O(n^{-2})P(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \epsilon_n) \\ &\quad + O(n^{-1})P(|\mathbf{Z}_1 - \mathbf{Z}_2| \leq \epsilon_n, |\mathbf{Z}_1 - \mathbf{Z}_3| \leq \epsilon_n) \\ &= O(n^{-2})P(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \epsilon_n) \\ &\quad + O(n^{-1}\epsilon_n^2). \end{aligned}$$

Thus

$$P(\delta_n \geq \epsilon_n) = O(n^{-1}),$$

so that, with this choice of ϵ_n ,

$$Eq_1 \leq c_3 n^{-2}.$$

It follows that

$$C_1 q_1 \rightarrow^p 0 \text{ as } n \rightarrow \infty.$$

Finally, consider

$$\begin{aligned} N^2 |q_0 - p^2| &= n^2 \left\{ \binom{n-2}{2} \right\}^{-1} \left| \sum_{\substack{(i,j) \cap (k,l) = \emptyset \\ i < j, k < l}} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \delta_n) I(|\mathbf{Z}_k - \mathbf{Z}_l| \leq \delta_n) \right. \\ &\quad \left. - \left(\sum_{i < j} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \delta_n) \right)^2 \right| + O(n^{-1}) \\ &\leq \frac{c_4}{n^2} \left\{ \sum_{i < j} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \delta_n) \right. \\ &\quad \left. + \sum_i \sum_{j < k} I(|\mathbf{Z}_i - \mathbf{Z}_j| \leq \delta_n) I(|\mathbf{Z}_i - \mathbf{Z}_k| \leq \delta_n) \right\} + O(n^{-1}) \\ &= O(n^{-1}) + c_5 n q_1. \end{aligned}$$

It follows from the above calculations that

$$EN^2 |q_0 - p^2| \rightarrow 0$$

and hence that

$$N^2 |q_0 - p^2| \rightarrow^p 0$$

also.

(ii) *Rare patterns.* Suppose that one has n sites on a line, m of them coloured or distinguished in some way. Assume, as a null hypothesis, that the coloured sites have been chosen at random. To test an alternative of contagion, one might count the number of times k coloured sites are contiguous, rejecting the null hypothesis if the count were too large. In order to fix the significance points of the test, the distribution under the null hypothesis of the number of occurrences of k contiguous coloured sites needs to be found. If k is chosen such that the occurrence of k coloured sites together is suitably rare, one would expect the number of such occurrences to be close to Poisson. A similar theorem is proved for a Bernoulli model in Brown (1981), Corollary 2.

Let Y_1, \dots, Y_n be random variables, indicating the type of site:

$$Y_i = \begin{cases} 1 \\ 0 \end{cases} \text{ if site is } \begin{cases} \text{coloured} \\ \text{otherwise.} \end{cases}$$

The Y_1, \dots, Y_n are exchangeable. Consider k to be fixed and set

$$E_n = \sum_{i=1}^{n-k+1} I(Y_i Y_{i+1} \cdots Y_{i+k-1} = 1),$$

the number of occurrences of k coloured sites together. Under the conditions of Corollary 1, E_n will be asymptotically equivalent to the W defined there, so one needs only to check those conditions. Here $M = n$ and

$$p = (m)_k / (n)_k$$

so that

$$np = n(m)_k / (n)_k.$$

If $np \rightarrow \lambda$, $m = O(n^{1-k^{-1}})$ and

$$nq_j = n(m)_{2k-j} / (n)_{2k-j} = O(n^{j/k-1}) \rightarrow 0 \quad \text{if } 1 \leq j \leq k-1.$$

Finally,

$$\begin{aligned} n^2 |q_0 - p^2| &= n^2 \left| \frac{(m)_{2k}}{(n)_{2k}} - \left[\frac{(m)_k}{(n)_k} \right]^2 \right| \\ &= O(n^{k-1-1}) \rightarrow 0 \quad \text{if } k \geq 2. \end{aligned}$$

Thus we have proved the following.

Corollary 6. Let Y_{n1}, \dots, Y_{nm} be exchangeable 0–1 random variables and set

$$E_n = \sum_{i=1}^{n-k+1} I(Y_{ni} Y_{n,i+1} \cdots Y_{n,i+k-1} = 1).$$

Then E_n converges in distribution to a Poisson law as $n \rightarrow \infty$ if $m = \sum_{i=1}^n Y_{ni}$ is such that

$$\lim_{n \rightarrow \infty} n(m)_k / (n)_k = \lambda.$$

5. Appendix

Bounds are required for the test function, defined in the introduction, namely

$$x(0) = 0$$

$$x(m+1) = \lambda^{-m-1} e^\lambda m! [\mathcal{P}_\lambda(A \cap U_m) - \mathcal{P}_\lambda(A) \mathcal{P}_\lambda(U_m)], \quad m \geq 0,$$

where $A \subset Z^+$ and where $\mathcal{P}_\lambda(\cdot)$ and U_m have already been defined. It is easy to show that x is bounded and Chen (1975), Lemmas 3.3 and 3.4, shows that

$$\|x\| \leq 4(\lambda^{-\frac{1}{2}} \wedge 1) \quad \text{and} \quad \Delta x \leq 6(\lambda^{-\frac{1}{2}} \wedge 1).$$

In order to have reasonable constants in Lemma 3 and to improve the results for large λ , we here derive better bounds for $\|x\|$ and Δx .

Lemma 4. For x defined as above,

- (i) $\|x\| \leq 1 \wedge 1 \cdot 4\lambda^{-\frac{1}{2}}$;
- (ii) $\Delta x \leq \lambda^{-1}(1 - e^{-\lambda}) \leq (1 \wedge \lambda^{-1})$.

Proof. From the definition of x ,

$$x(m+1) = \lambda^{-m-1} m! e^\lambda \{ \mathcal{P}_\lambda(A \cap U_m) \mathcal{P}_\lambda(U_m^-) - \mathcal{P}_\lambda(A \cap U_m^-) \mathcal{P}_\lambda(U_m) \},$$

where B^- denotes the complement of B . Hence, since $0 \leq P(A \cap B)P(B^-) \leq P(B)P(B^-)$, for any $A \subset Z^+$,

$$(2) \quad |x(m+1)| \leq m! e^\lambda \lambda^{-m-1} \mathcal{P}_\lambda(U_m) \mathcal{P}_\lambda(U_m^-).$$

For $m \leq \lambda$,

$$(3) \quad |x(m+1)| \leq m! e^\lambda \lambda^{-m-1} \mathcal{P}_\lambda(U_m) = \frac{1}{\lambda} \sum_{j=0}^m \lambda^{-j} m! / (m-j)! \\ \leq \frac{1}{\lambda} \sum_{j=0}^m (m/\lambda)^j \leq (\lambda - m)^{-1};$$

similarly, for $m \geq \lambda - 3$,

$$(4) \quad |x(m+1)| \leq m! e^\lambda \lambda^{-m-1} \mathcal{P}_\lambda(U_m^-) = \sum_{j=0}^\infty \lambda^j m! / (m+1+j)! \\ \leq \frac{1}{m+1} \left\{ 1 + \frac{\lambda}{m+2} \sum_{j=0}^\infty (\lambda/(m+3))^j \right\} \\ = \frac{(m+2)(m+3) + \lambda}{(m+1)(m+2)(m+3-\lambda)}.$$

Hence, from (3), $|x(m+1)| \leq 1$ whenever $m \leq \lambda - 1$; from (4), $|x(m+1)| \leq 1$ whenever $m \geq \lambda - 1$ if $m \geq 2$, and whenever $\lambda \leq 12/7$ if $m = 1$. It therefore remains to consider $m = 0$ with $0 \leq \lambda \leq 1$ and $m = 1$ with $12/7 \leq \lambda \leq 2$. However, it follows from (2) that, for $m = 0$, $|x(1)| \leq \lambda^{-1}(1 - e^{-\lambda}) \leq 1$ for all λ : for $m = 1$,

$$|x(2)| \leq \{\lambda^{-2}(1 + \lambda)\} \{1 - e^{-\lambda}(1 + \lambda)\},$$

and, by observing that the two terms in braces are monotonic in λ , it follows that, for $12/7 \leq \lambda \leq 2$, $|x(2)| \leq (133/144)\{1 - 3e^{-2}\} < 1$. Thus $\|x\| \leq 1$.

Returning to (2), it is always true that $\mathcal{P}_\lambda(U_m) \mathcal{P}_\lambda(U_m^-) \leq \frac{1}{4}$. Using Stirling's approximation that $m! \leq \sqrt{(2\pi)} m^{m+\frac{1}{2}} \exp(-m + 1/12m)$, $m \geq 1$, it follows that

$$(5) \quad |x(m+1)| \leq \frac{1}{4} \sqrt{(2\pi)} \lambda^{-\frac{1}{2}} (m/\lambda)^{m+\frac{1}{2}} \exp\left(\lambda - m + \frac{1}{12m}\right) \\ \leq \frac{1}{4} \sqrt{(2\pi)} \lambda^{-\frac{1}{2}} \exp\left\{\lambda^{-1}(m-\lambda)(m-\lambda+\frac{1}{2}) + \frac{1}{12m}\right\}, \quad m \geq 1.$$

As before, $|x(1)| \leq \lambda^{-1}(1 - e^{-\lambda}) \leq \lambda^{-1}$. Comparing (5) with (3) and (4), and taking the smaller, yields the estimate $\|x\| \leq c\lambda^{-\frac{1}{2}}$, for c with a value no greater than 1.4. This completes the proof of Part (i).

For Part (ii), let x_j denote $x_{\lambda, \{j\}}$, and note that $x_{\lambda, A} = \sum_{j \in A} x_j$. From its definition,

$$x_j(m+1) = \begin{cases} (\lambda^{j-1}/j!) \{m! \lambda^{-m} \mathcal{P}_\lambda(U_m^-)\} & m \geq j \\ -(\lambda^{j-1}/j!) \{m! \lambda^{-m} \mathcal{P}_\lambda(U_m)\} & m < j, \end{cases}$$

and, from the series expansion of the Poisson probabilities, it is easily seen that x_j is positive and decreasing in $m \geq j + 1$ and is negative and decreasing in $m \leq j$. Hence the only positive value taken by $x_j(m + 1) - x_j(m)$ is

$$\begin{aligned} x_j(j + 1) - x_j(j) &= \frac{e^{-\lambda}}{\lambda} \left\{ \sum_{r=j+1}^{\infty} (\lambda^r/r!) + \sum_{r=1}^j (\lambda^r/r!) \frac{r}{j} \right\} \\ &\leq \frac{e^{-\lambda}}{\lambda} \{e^\lambda - 1\} = \lambda^{-1} \{1 - e^{-\lambda}\}. \end{aligned}$$

Thus, for any $A \subset Z^+$,

$$\begin{aligned} x_{\lambda,A}(m + 1) - x_{\lambda,A}(m) &= I[m \in A] \{x_m(m + 1) - x_m(m)\} \\ &\quad + \sum_{\substack{j \in A \\ j \neq m}} \{x_j(m + 1) - x_j(m)\}, \end{aligned}$$

in which expression only the first term is positive, and so

$$\sup_{m \in Z^+} \max_{A \subset Z^+} \{x_{\lambda,A}(m + 1) - x_{\lambda,A}(m)\} = \sup_{m \in Z^+} \{x_m(m + 1) - x_m(m)\} \leq \lambda^{-1} (1 - e^{-\lambda}).$$

Since also $x_{\lambda,A} = -x_{\lambda,A^c}$, Part (ii) follows from the above estimate.

Acknowledgements

The authors are grateful for the help of J. Dawson in performing the simulations mentioned in Section 4, for fruitful discussions with T. C. Brown and D. J. Aldous, and for a referee’s pertinent comments on the presentation of the paper.

References

BARTON, D. E. AND DAVID, F. N. (1966) The random intersection of two graphs. In *Research Papers in Statistics*, ed. F. N. David, Wiley, New York, 445–459.

BERMAN, M. AND EAGLESON, G. K. (1983) A Poisson limit theorem for incomplete symmetric statistics. *J. Appl. Prob.* **20**, 47–60.

BROWN, T. C. (1981) Compensators and Cox convergence. *Math. Proc. Camb. Phil. Soc.* **90**, 305–319.

BROWN, T. C. AND SILVERMAN, B. S. (1979) Rates of Poisson convergence for U -statistics. *J. Appl. Prob.* **16**, 428–432.

CHEN, L. H. Y. (1975) Poisson approximation for dependent trials. *Ann. Prob.* **3**, 534–545.

EAGLESON, G. K. (1979) A Poisson limit theorem for weakly exchangeable arrays. *J. Appl. Prob.* **16**, 794–802.

KNOX, G. (1964) Epidemiology of childhood leukaemia in Northumberland and Durham. *Brit. J. Prev. Soc. Med.* **18**, 17–24.

SILVERMAN, B. S. AND BROWN, T. C. (1978) Short distances, flat triangles and Poisson limits. *J. Appl. Prob.* **15**, 815–825.

STEIN, C. (1970) A bound for the error in the normal approximation to the distribution of a sum of dependent random variables. *Proc. 6th Berkeley Symp. Math. Statist. Prob.* **2**, 583–602.