

Manly (1984) and references therein; see also Cormack (1968) and the conference article of the same author (1973) who begins with the following remarks:

Many of the papers in this volume are concerned with the process of describing the development of an animal population by a mathematical model. The properties of such a model can then be derived, either by elegant mathematics or equally elegant computer simulation, in order to describe the future state of the population in terms of certain initial boundary conditions. The model becomes of scientific value when such predictions can be tested, which requires in turn that the mathematical symbols can be replaced by numbers. The parameters of the model must be estimated from data of a type that a biologist can collect about the population he is studying.

For an introductory treatment written for biologists, see Begon (1979).

3.3 THE POISSON DISTRIBUTION

We recall the definition and some elementary properties of Poisson random variables.

Definition A non-negative integer-valued random variable X has a Poisson distribution with parameter $\lambda > 0$ if

$$p_k = \Pr \{X = k\} = \frac{e^{-\lambda} \lambda^k}{k!}, \quad k = 0, 1, 2, \dots \quad (3.7)$$

From the definition of e^λ as $\sum_0^\infty \lambda^k/k!$ we find

$$\sum_{k=0}^\infty \Pr \{X = k\} = 1.$$

The mean and variance of X will easily be found to be

$$E(X) = \text{Var}(X) = \lambda.$$

The shape of the probability mass function depends on λ as Table 3.1 and the graphs of Fig. 3.2 illustrate.

Table 3.1 Probability mass functions for some Poisson random variables

	p_0	p_1	p_2	p_3	p_4	p_5	p_6	p_7	p_8
$\lambda = \frac{1}{2}$.607	.303	.076	.013	.002	<.001			
$\lambda = 1$.368	.368	.184	.061	.015	.003	<.001		
$\lambda = 2$.135	.271	.271	.180	.090	.036	.012	.003	<.001

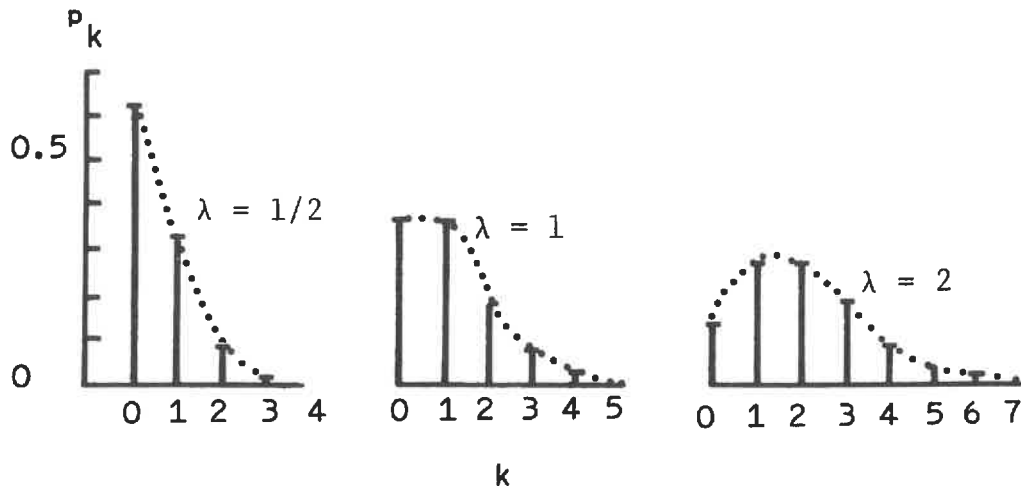


Figure 3.2 Probability mass functions for Poisson random variables with various parameter values.

There are two points which emerge just from looking at Fig. 3.2:

- (i) Poisson random variables with different parameters can have quite different looking mass functions.
- (ii) When λ gets large the mass function has the shape of a normal density (see Chapter 6).

Poisson random variables arise frequently in counting numbers of events. We will consider events which occur randomly in one-dimensional space or time and in two-dimensional space, the latter being of particular relevance in ecology. Generalizations to higher-dimensional spaces will also be briefly discussed.

3.4 HOMOGENEOUS POISSON POINT PROCESS IN ONE DIMENSION

Let t represent a time variable. Suppose an experiment begins at $t = 0$. Events of a particular kind occur randomly, the first being at T_1 , the second at T_2 , etc., where T_1, T_2 , etc., are random variables. The values t_i of $T_i, i = 1, 2, \dots$ will be called **points of occurrence** or just *events* (see Fig. 3.3).

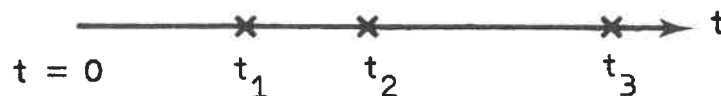


Figure 3.3 Events of a particular kind occur at t_1, t_2 , etc.

Let $(s_1, s_2]$ be a subinterval of the interval $[0, s]$ where $s < \infty$. Denote by $N(s_1, s_2)$ the number of points of occurrence in $(s_1, s_2]$. Then $N(s_1, s_2)$ is a random variable and the collection of all such random variables, abbreviated to N , for various subintervals (or actually any subsets of $[0, s]$) is called a **point process** on $[0, s]$.

Definition N is an homogeneous Poisson point process with rate λ if:

- (i) for any $0 \leq s_1 < s_2 \leq s$, $N(s_1, s_2)$ is a Poisson random variable with parameter $\lambda(s_2 - s_1)$;
- (ii) for any collection of times $0 \leq s_0 < s_1 < s_2 \dots < s_n \leq s$, where $n \geq 2$, the random variables $\{N(s_{k-1}, s_k), k = 1, 2, \dots, n\}$ are mutually independent.

We see therefore that the number of points of occurrence in $(0, t]$, which we denote by just $N(t)$, is a Poisson random variable with parameter λt . Also, the numbers of points falling in disjoint intervals are independent.

Now, the expected value of $N(t)$ is λt and this is also the expected number of points in $(0, t]$. Thus the expected number of points in the unit interval $(0, 1]$, or any other interval of unit length, is just λ . Hence the description of λ as the rate parameter, or as it is often called, the **intensity** of the process. The process is called **homogeneous** because the probability law of the number of points in any interval depends only on the length of the interval, not on its location.

We will now derive some properties of interest in connection with the distances (or time intervals) between points of occurrence (events) when the Poisson point process is defined on subsets of $[0, \infty)$. The role of s will now change.

The waiting time to the next event

Consider any fixed time point $s > 0$. Let T_1 be the time which elapses before the first event after s . Then we have the following result.

Theorem 3.3 The waiting time, T_1 , for an event is exponentially distributed with mean $1/\lambda$.

Note that the distribution of T_1 does not depend on s . We say the process has no memory, a fact which is traced to the definition since the numbers of events in $(s_1, s_2]$ and $(s_2, s_3]$ are independent.

Proof First we note that the probability of one event in any small interval of length Δt is

$$e^{-\lambda \Delta t}(\lambda \Delta t) = \lambda \Delta t + o(\Delta t), \quad (3.8)$$

where $o(\Delta t)$ here stands for terms which vanish faster than Δt as Δt goes to

40 Applications of hypergeometric and Poisson distributions

zero. We will have $T_1 \in (t, t + \Delta t]$ if there are no events in $(s, s + t]$ and one event in $(s + t, s + t + \Delta t]$. By independence, the probability of both of these is the product of the probabilities of either occurring separately. Hence

$$\Pr \{T_1 \in (t, t + \Delta t]\} = e^{-\lambda t}[\lambda \Delta t + o(\Delta t)].$$

It follows that the density of T_1 is given by

$$f_{T_1}(t) = \lambda e^{-\lambda t}, \quad t > 0$$

Alternatively, this result may be obtained by noting that

$$\Pr \{T_1 > t\} = \Pr \{N(s, s + t) = 0\} = e^{-\lambda t}.$$

We find that not only is the waiting time to an event exponentially distributed but also the following is true.

Theorem 3.4 The time interval between events is exponentially distributed with mean $1/\lambda$.

Proof This is Exercise 8.

In fact it can be shown that if the distances between consecutive points of occurrence are independent and identically exponentially distributed, then the point process is a homogeneous Poisson point process. This statement provides one basis for statistical tests for a Poisson process (Cox and Lewis, 1966).

The waiting time to the k th point of occurrence

Theorem 3.5 Let T_k be the waiting time until the k th event after s , $k = 1, 2, \dots$. Then T_k has a gamma density with parameters k and λ .

Proof The k th point of occurrence will be the only one in $(s + t, s + t + \Delta t]$ if and only if there are $k - 1$ points in $(s, s + t]$ and one point is in $(s + t, s + t + \Delta t]$. It follows from (3.7) and (3.8) that

$$\Pr \{T_k \in (t, t + \Delta t]\} = \frac{e^{-\lambda t}(\lambda t)^{k-1}[\lambda \Delta t + o(\Delta t)]}{(k-1)!}, \quad k = 1, 2, \dots$$

Hence the density of T_k is

$$f_{T_k}(t) = \frac{\lambda(\lambda t)^{k-1}e^{-\lambda t}}{(k-1)!}, \quad t > 0 \tag{3.9}$$

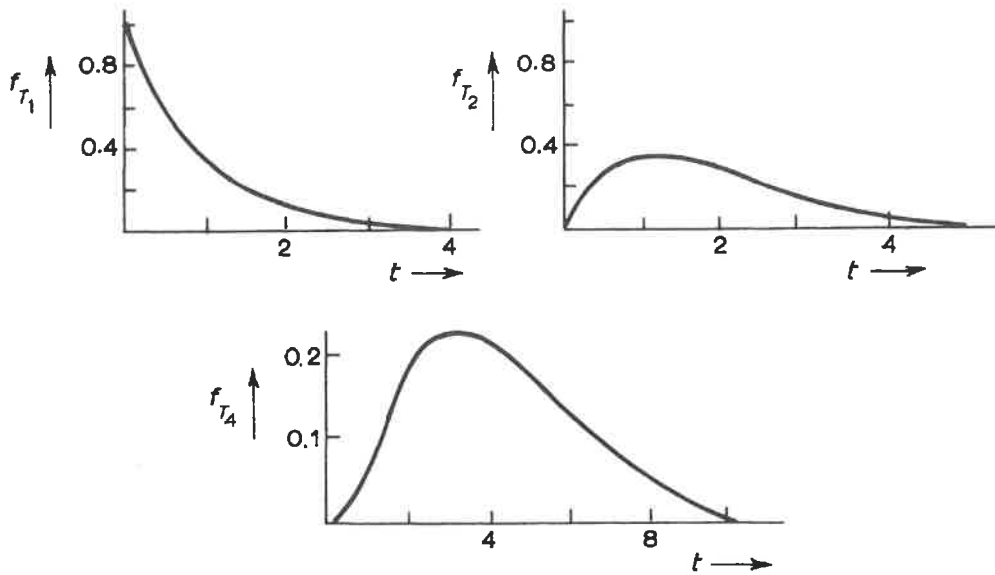


Figure 3.4 The densities of the waiting times for 1, 2 and 4 events in a homogeneous Poisson point process with $\lambda = 1$.

and the mean and variance of T_k are given by

$$E(T_k) = \frac{k}{\lambda}, \quad \text{Var}(T_k) = \frac{k}{\lambda^2}.$$

Note that this result can also be deduced from the fact that the sum of $k > 1$ independent exponentially distributed random variables, each with mean $1/\lambda$, has a gamma density as in (3.9) (prove this by using Theorem 2.4). Furthermore, it can be shown that the waiting time to the k th event after an event has a density given by (3.9).

The waiting times for $k = 1, 2$ and 4 events have densities as depicted in Fig. 3.4. Note that as k gets larger, the density approaches that of a normal random variable (see Chapter 6, when we discuss the central limit theorem).

3.5 OCCURRENCE OF POISSON PROCESSES IN NATURE

The following reasoning leads in a natural way to the Poisson point process. Points, representing the times of occurrence of an event, are sprinkled randomly on the interval $[0, s]$ under the assumptions:

- (i) the numbers of points in disjoint subintervals are independent;
- (ii) the probability of finding a point in a very small subinterval is

42 Applications of hypergeometric and Poisson distributions

proportional to its length, whereas the probability of finding more than one point is negligible.

It is convenient to divide $[0, s]$ into n subintervals of equal length $\Delta s = s/n$. Under the above assumptions the probability p that a given subinterval contains a point is $\lambda s/n$ where λ is a positive constant. Hence the chance of k occupied subintervals is

$$\begin{aligned}\Pr \{k \text{ points in } [0, s]\} &= b(k; n, p) \\ &= b\left(k; n, \frac{\lambda s}{n}\right).\end{aligned}$$

Now as $n \rightarrow \infty$, $\lambda s/n \rightarrow 0$ and we may invoke the Poisson approximation to the binomial probabilities (see also Chapter 6):

$$b(k; n, p) \xrightarrow{n \rightarrow \infty} \frac{\exp(-np)(np)^k}{k!}.$$

But $np = n(\lambda s)/n = \lambda s$. Hence in the limit as $n \rightarrow \infty$,

$$\Pr \{k \text{ points in } [0, s]\} = \frac{\exp(-\lambda s)(\lambda s)^k}{k!},$$

as required.

The above assumptions and limiting argument should help to make it understandable why approximations to Poisson point processes arise in the study of a broad range of natural random phenomena. The following examples provide evidence for this claim.

Examples

(i) Radioactive decay

The times at which a collection of atomic nuclei emit, for example, alpha-particles can be well approximated as a Poisson point process. Suppose there are N observation periods of duration T , say. In Exercise 18 it is shown that under the Poisson hypothesis, the expected value, n_k , of the number, N_k , of observation periods containing k emissions is

$$n_k = \frac{N \exp(-\bar{n})\bar{n}^k}{k!}, \quad k = 0, 1, 2, \dots \quad (3.10)$$

where $\bar{n} = \lambda T$ is the expected number of emissions per observation period. For an experimental data set, see Feller (1968, p. 160).

(ii) Arrival times

The times of arrival of customers at stores, banks, etc., can often be approximated by Poisson point processes. Similarly for the times at which

phone calls are made, appliances are switched on, accidents in factories or in traffic occur, etc. In queueing theory the Poisson assumption is usually made (see for example Blake, 1979), partly because of empirical evidence and partly because it leads to mathematical simplifications. In most of these situations the rate may vary so that $\lambda = \lambda(t)$. However, over short enough time periods, the assumption that λ is constant will often be valid.

(iii) *Mutations*

In cells changes in genetic (hereditary) material occur which are called mutations. These may be spontaneous or induced by external agents. If mutations occur in the reproductive cells (gametes) then the offspring inherits the mutant genes. In humans the rate at which spontaneous mutations occur per gene is about 4 per hundred thousand gametes (Strickberger, 1968). In the common bacterium *E. coli*, a mutant variety is resistant to the drug streptomycin. In one experiment, $N = 150$ petri dishes were plated with one million bacteria each. It was found that 98 petri dishes had no resistant colonies, 40 had one, 8 had two, 3 had three and 1 had four. The average number \bar{n} of mutants per million cells (bacteria) is therefore

$$\bar{n} = \frac{40 \times 1 + 8 \times 2 + 3 \times 3 + 1 \times 4}{150} = 0.46.$$

Under the Poisson hypothesis, the expected numbers n_k of dishes containing k mutants are as given in Table 3.2, as calculated using (3.10). The observed values N_k are also given and the agreement is reasonable. This can be demonstrated with a χ^2 test (see Chapter 1).

(iv) *Voltage changes at nerve–muscle junction*

The small voltage changes seen in a muscle cell attributable to spontaneous activity in neighbouring nerve cells occur at times which are well described as a Poisson point process. A further aspect of this will be elaborated on in Section 3.9. Figure 3.5 shows an experimental histogram of waiting times

Table 3.2 Bacterial mutation data*

k	n_k	$N_k(\text{Obs.})$
0	94.7	98
1	43.5	40
2	10.0	8
3	1.5	3
4	0.2	1

*From Strickberger (1968).

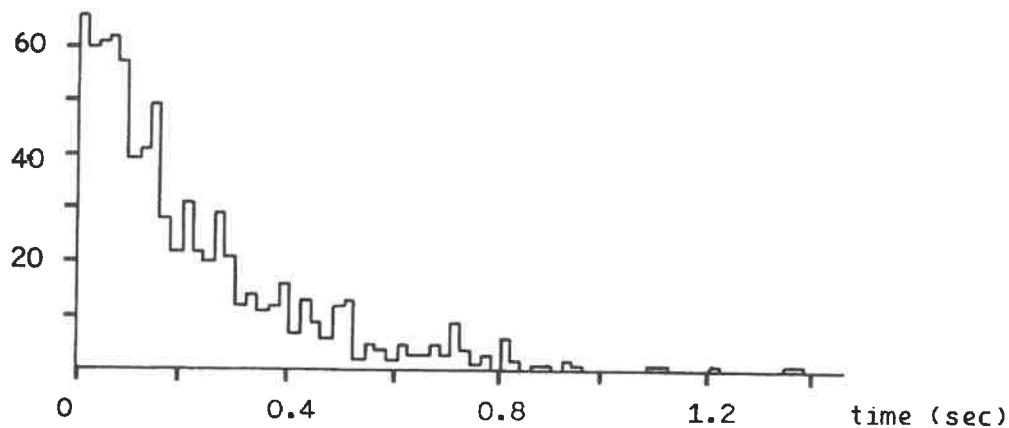


Figure 3.5 A histogram of waiting times between spontaneously occurring small voltage changes in a muscle cell due to activity in a neighbouring nerve cell. From Fatt and Katz (1952).

between such events. According to the Poisson assumption, the waiting time should have an exponential density which is seen to be a good approximation to the observed data. This may also be rendered more precise with a χ^2 goodness of fit test. For further details see Van der Kloot *et al.* (1975).

3.6 POISSON POINT PROCESSES IN TWO DIMENSIONS

Instead of considering random points on the line we may consider random points in the plane $\mathbb{R}^2 = \{(x, y) | -\infty < x < \infty, -\infty < y < \infty\}$, or subsets thereof.

Definition A point process N is an homogeneous Poisson point process in the plane with intensity λ if:

- (i) for any subset A of \mathbb{R}^2 , the number of points $N(A)$ occurring in A is a Poisson random variable with parameter $\lambda|A|$, where $|A|$ is the area of A ;
- (ii) for any collection of disjoint subsets of \mathbb{R}^2 , A_1, A_2, \dots, A_n , the random variables $\{N(A_k), k = 1, 2, \dots, n\}$ are mutually independent.

Note that the number of points in $[0, x] \times [0, y]$ is a Poisson random variable with parameter λxy . Putting $x = y = 1$ we find that the number of points in the unit square is Poisson with parameter λ . Hence λ is the expected number of points per unit area.

Application to ecological patterns

Ecologists are interested in the spatial distributions of plants and animals (see for example MacArthur and Connell, 1966). Three of the situations of interest are:

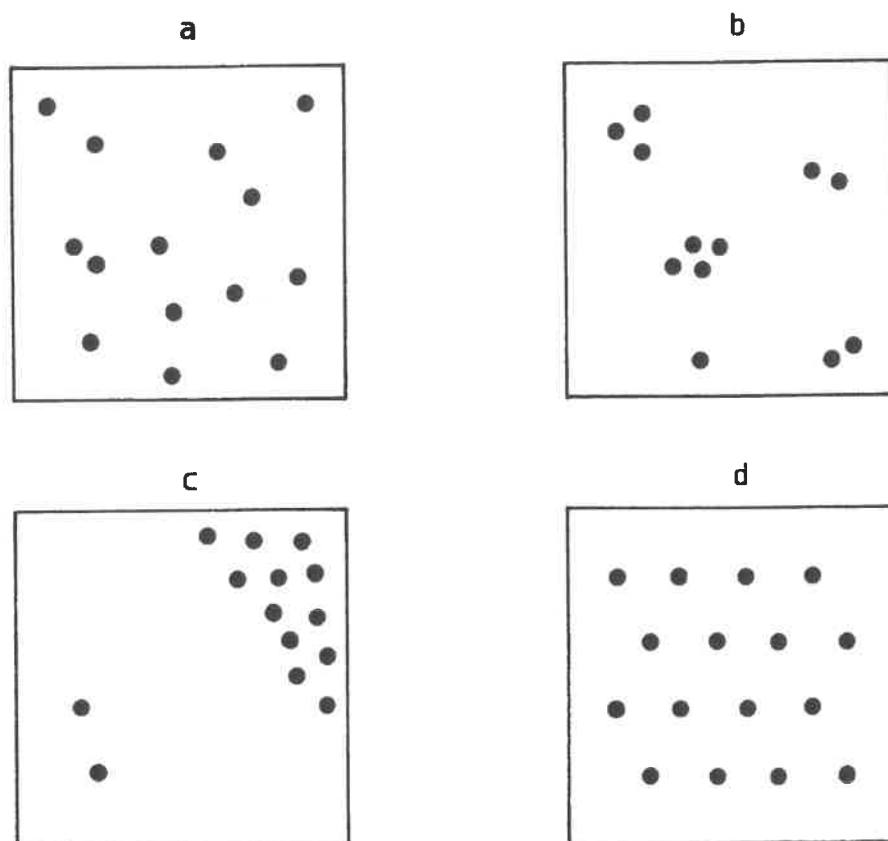


Figure 3.6 Some representative spatial patterns of organisms: (a) random, (b) clumping in groups, (c) preferred location, (d) regular.

- (i) the organisms are distributed randomly;
- (ii) the organisms have preferred locations in the sense that they tend to occur in groups (i.e. are clustered or clumped) or in some regions more frequently than others;
- (iii) the organisms are distributed in a regular fashion in the sense that the distances between them and their nearest neighbours tend to be constant.

These situations are illustrated in Fig. 3.6. We note that clumping indicates cooperation between organisms. The kind of spacing shown in Fig. 3.6(d) indicates competition as the organisms tend to maintain a certain distance between themselves and their neighbours.

An important reason for analysing the underlying pattern is that if it is known, the total population may be estimated from a study of the numbers in a small region. This is of particular importance in the forest industry.

The hypothesis of randomness leads naturally, by the same kind of argument as in Section 3.5, to a Poisson point process in the plane. Ecologists refer to this as a **Poisson forest**. Under the assumption of a Poisson forest we may derive the probability density function of the distance from one organism

(e.g. tree) to its nearest neighbour. We may use this density to test the hypothesis of randomness. We first note the following result.

Theorem 3.6 In a Poisson forest, the distance R_1 from an arbitrary fixed point to the nearest event has the probability density

$$f_{R_1}(r) = 2\lambda\pi r e^{-\lambda\pi r^2}, \quad r > 0. \quad (3.11)$$

Proof We will have $R_1 > r$ if and only if there are no events in the circle of radius r with centre at the fixed point under consideration. Such a circle has area πr^2 , so from the definition of a Poisson point process in the plane, the number of events inside the circle is a Poisson random variable with mean $\lambda\pi r^2$. This gives

$$\Pr\{R_1 > r\} = e^{-\lambda\pi r^2}.$$

We must then have

$$f_{R_1}(r) = \frac{d}{dr}(1 - e^{-\lambda\pi r^2})$$

which leads to (3.11) as required.

We may also prove that the distance from an event to its nearest neighbour in a Poisson forest has the density given by (3.11). It is left as an exercise to prove the following result.

Theorem 3.7 In a Poisson forest the distance R_k to the k th nearest event has the density

$$f_{R_k}(r) = \frac{2\pi\lambda r (\lambda\pi r^2)^{k-1} e^{-\lambda\pi r^2}}{(k-1)!}, \quad r > 0, \quad k = 1, 2, \dots$$

Estimating the number of trees in a forest

If one is going to estimate the number of trees in a forest, it must first be ensured that the assumed probability model is valid. The obvious hypothesis to begin with is that one is dealing with a Poisson process in the plane. A few methods of testing this hypothesis and a method of estimating λ are now outlined. For some further references see Patil *et al.* (1971) and Heltshe and Ritchey (1984). An actual data set is shown in Fig. 3.7.

Method 1 – Distance measurements

Under the assumption of a Poisson forest the point–nearest tree or tree–nearest tree distance has the density f_{R_1} given in (3.11). The actual measure–

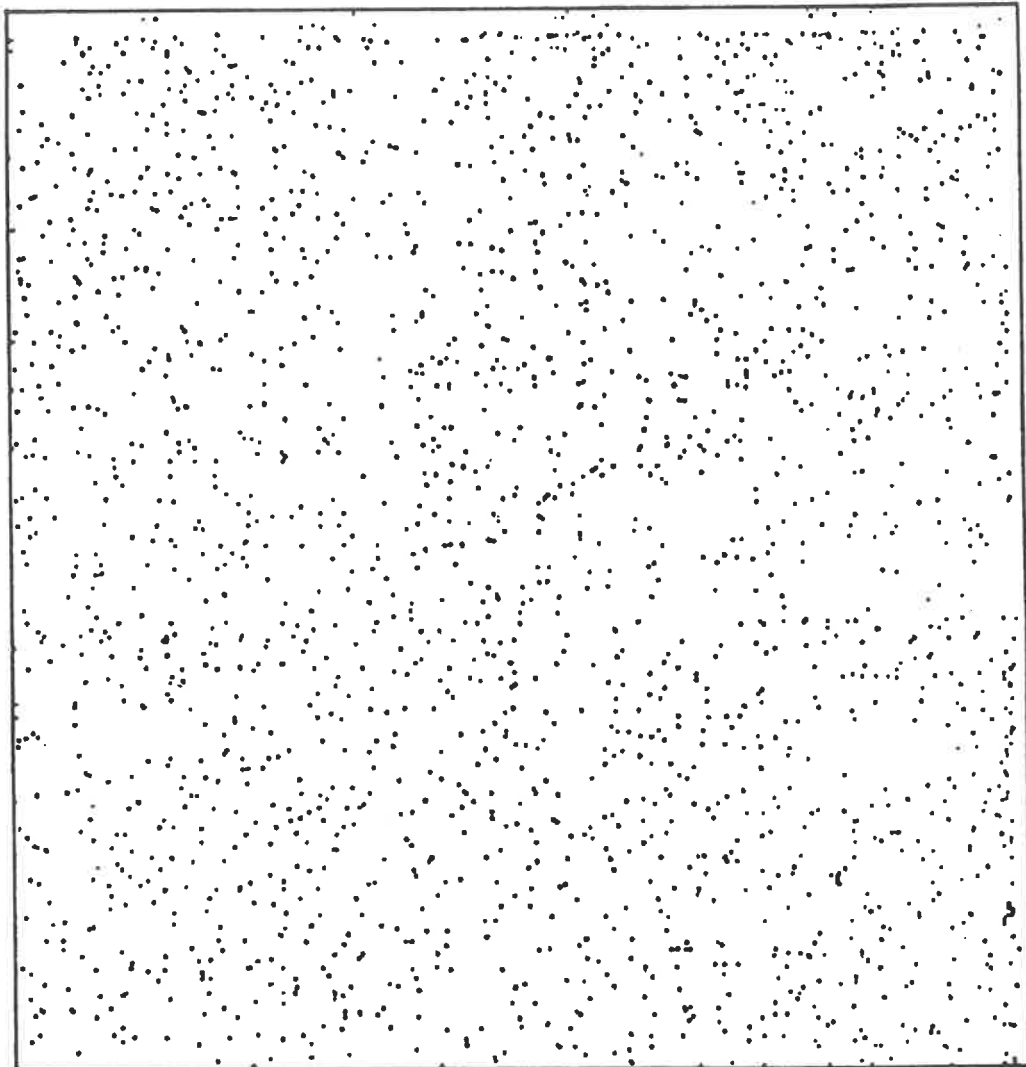


Figure 3.7 Locations of trees in Lansing Woods. Smaller dots represent oaks, larger dots represent hickories and maples. The data are analysed in Exercise 22. Reproduced with permission from Clayton (1984).

ments of such distances may be collected into a histogram or empirical distribution function. A goodness of fit test such as χ^2 (see Chapter 1) or Kolmogorov–Smirnov (see for example Hoel, 1971; or Afifi and Azen, 1979) can be carried out. Note that **edge effects** must be minimized since the density of R_1 was obtained on the basis of an infinite forest.

Assuming a Poisson forest the parameter λ may be estimated as follows. Let $\{X_i, i = 1, 2, \dots, n\}$ be a random sample for the random variable with the density (3.11). Then it is shown in Exercise 21 that an **unbiased estimator** (see

Exercise 6) of $1/\lambda$ is

$$\hat{\Lambda}^{-1} = \frac{\pi}{n} \sum_{i=1}^n X_i^2.$$

An estimate of λ is thus made and hence, if the total area A is known, the total number of trees may be estimated as λA . For further details see Diggle (1975, 1983), Ripley (1981) and Upton and Fingleton (1985).

Method 2—Counting

Another method of testing the hypothesis of a Poisson forest is to subdivide the area of interest into N equal smaller areas called cells. The numbers N_k of cells containing k plants can be compared using a χ^2 -test with the expected numbers under the Poisson assumption using (3.10), with \bar{n} = the mean number of plants per cell.

Extensions to three and four dimensions

Suppose objects are randomly distributed throughout a 3-dimensional region. The above concepts may be extended by defining a Poisson point process in \mathbb{R}^3 . Here, if A is a subset of \mathbb{R}^3 , the number of objects in A is a Poisson random variable with parameter $\lambda|A|$, where λ is the mean number of objects per unit volume and $|A|$ is the volume of A . Such a point process will be useful in describing distributions of organisms in the ocean or the earth's atmosphere, distributions of certain rocks in the earth's crust and of objects in space. Similarly, a Poisson point process may be defined on subsets of \mathbb{R}^4 with a view to describing random events in space-time.

3.7 COMPOUND POISSON RANDOM VARIABLES

Let $X_k, k = 1, 2, \dots$ be independent identically distributed random variables and let N be a non-negative integer-valued random variable, independent of the X_k . Then we may form the following sum:

$$S_N = X_1 + X_2 + \dots + X_N, \quad (3.12)$$

where the number of terms is determined by the value of N . Thus S_N is a **random sum of random variables**: we take S_N to be zero if $N = 0$. If N is a Poisson random variable, S_N is called a **compound Poisson random variable**. The mean and variance of S_N are then as follows.

Theorem 3.8 Let $E(X_1) = \mu$ and $\text{Var}(X_1) = \sigma^2$, $|\mu| < \infty, \sigma < \infty$. If N is