



## On Generating Random Intervals and Hyperrectangles

Luc Devroye; Peter Epstein; Jörg-Rüdiger Sack

*Journal of Computational and Graphical Statistics*, Vol. 2, No. 3. (Sep., 1993), pp. 291-307.

Stable URL:

<http://links.jstor.org/sici?sici=1061-8600%28199309%292%3A3%3C291%3AOGRIAH%3E2.0.CO%3B2-G>

*Journal of Computational and Graphical Statistics* is currently published by American Statistical Association.

---

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/astata.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 and preserving a digital archive of scholarly journals. For more information regarding JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).

# On Generating Random Intervals and Hyperrectangles

LUC DEVROYE\* , PETER EPSTEIN† , AND JÖRG-RÜDIGER SACK‡

We look at the problem of generating a random hyperrectangle in a unit hypercube such that each point of the hypercube has probability  $p$  of being covered. For random intervals of  $[0, 1]$ , we compare various methods based either on the distribution of the length or on the distribution of the midpoint. It is shown that no constant length solution exists. Nice versatility is achieved when a random interval is chosen from among the spacings defined by a Dirichlet process. A small simulation is included.

**Key Words:** Computational geometry; Dirichlet process; Monte Carlo; Random cover; Random variate generation; Uniform coverage probability.

## 1. INTRODUCTION

We consider the following problem: Given a set  $A$  that is either  $[0, 1]^d$  or  $\mathbb{R}^d$ , generate a random hyperrectangle  $\prod_{i=1}^d I_i$ , where  $I_1, \dots, I_d$  are random intervals with the property that each  $x \in A$  has equal probability of being covered. The coverage probability  $p$  is specified beforehand. Stated in this manner, there are some trivial solutions, such as:

the all-or-nothing method

```
generate a uniform  $[0, 1]$  random variable  $U$ 
if  $U \leq p$  then return  $A$ 
else return the empty set
```

This example shows that we should provide additional restrictions. Before we go further, however, some observations will help us to better focus. By continuous, strictly monotone bijections between  $[0, 1]$  and  $\mathbb{R}$ , it is easy to see that we need only consider  $[0, 1]$ . Applying transformations coordinatewise, a similar observation is valid for  $\mathbb{R}^d$ . Furthermore, assume that if we can solve a one-dimensional problem with coverage probability  $p$ , then we can solve the  $d$ -dimensional problem with coverage probability  $p^d$  just by forming the product of  $d$  independent intervals. For these reasons we will now

---

\*School of Computer Science, McGill University, Montreal, Canada H3A 2A7

†School of Computer Science, Carleton University, Ottawa, Canada K1S 5B6

‡School of Computer Science, Carleton University, Ottawa, Canada K1S 5B6

©1993 American Statistical Association, Institute of Mathematical Statistics,  
and Interface Foundation of North America

*Journal of Computational and Graphical Statistics*, Volume 2, Number 3, Pages 291–307

concentrate on  $[0, 1]$ . There is a caveat: If one specifies the shape of the sets in  $\mathbb{R}^d$ , such as circles or squares, then the  $d$ -dimensional problem becomes nontrivial and cannot be solved by marrying  $d$  one-dimensional solutions.

We will present several models. Let  $(L, R)$  be the random interval that we are studying. In some cases we produce closed intervals  $[L, R]$ , but the distinction is rather minor. We thus require of all the methods the following property: Given  $p$ ,

$$\Pr\{x \in (L, R)\} = p$$

for all  $x \in (0, 1)$ . There are many possible solutions to this problem. Additional issues that govern a user's choice include the following:

1. The distribution of the length  $D = R - L$ . It is impossible to have  $D$  equal a positive constant and still ensure uniform probability coverage. While  $\mathbf{E}D = p$  in all cases, the oscillatory nature of  $D$  is to some extent captured in the variance,  $\text{var } D$ .
2. The distribution of  $L$ ,  $R$ , and  $M = (L+R)/2$ . Somehow, one feels that  $M$  should be almost uniformly distributed in most cases, but again this is not necessarily true.
3. The probability that two independently generated intervals are nested could have a particular significance.
4. The presence or absence of atoms in the distributions of  $L$  and  $R$  may influence the attractiveness of certain approaches. It is disappointing to note that uniform probability coverage is impossible without introducing atoms in both the distributions of  $L$  and  $R$ .
5. The flexibility in molding and designing the distributions is of the utmost importance. The more things users can "try out," the more attractive a method becomes.
6. Intervals generated for the purpose of testing algorithms are to be stored in a data structure. For example,  $n$  intervals induce an interval graph on  $n$  nodes, in which two nodes are connected if their intervals intersect. Random interval graphs have been studied by Scheinerman (1988, 1990). The edge density is, of course, controlled by the distribution of  $(L, R)$ ; in all cases of interest to us, the expected number of edges is  $\Theta(pn^2)$ , so this can be controlled by the choice of  $p$ .
7. Finally, the computational complexity of a method has to be rigorously controlled.

Uniform coverage probabilities are required in several types of statistical simulation experiments. For example, consider a simulation of customers queuing for bank tellers. The work shift of a bank teller is a time interval. The collection of work shifts for all tellers is a set of (not necessarily equal) intervals. A random time shift for a teller is defined as one in which the teller works at any given time of the day with equal probability. Collections of random time shifts are advantageous as they assure uniformly good bank service throughout the day. In a second example, consider an application that produces city maps from satellite images. Cloud cover causes problems for such a system. To model the effects of cloud cover before launching the satellite, a random cloud generator may be used. For simplicity, such a generator could always produce rectangular clouds. The appropriate probability distribution would be uniform coverage probability, because that would imply that any given point in the map area has the same probability of

being covered by cloud. Finally, we should mention the experimental verification of the performance of algorithms for reporting all the rectangles in a collection of rectangles that cover a given point  $x$ . For expected time performance we require that the rectangles be iid and that each point  $x$  has an equal probability of being covered by a given rectangle.

We begin with a few theoretical properties of random intervals that achieve uniform coverage probability. Then six methods are proposed that culminate in a method where we generate one of  $n$  intervals defined by  $n - 1$  random points with a Dirichlet distribution on  $[0, 1]$ . We conclude with a small simulation study that highlights the flexibility of the Dirichlet process method afforded by its distributional parameters.

## 2. A FEW THEORETICAL UNDERPINNINGS

The fundamental property is that

$$\mathbf{E}D = \mathbf{E}(R - L) = p . \tag{2.1}$$

This is best seen as

$$\mathbf{E}(R - L) = \mathbf{E} \int_0^1 I_{x \in (L, R)} dx = \int_0^1 \Pr\{x \in (L, R)\} dx = \int_0^1 p dx = p .$$

With intervals of the form  $(L, R) \subset [0, 1]$ , we cannot cover the endpoints 0 and 1. If one desires that  $\Pr\{x \in I\} = p$  holds for  $x \in [0, 1]$ , then  $I$  must be of the form  $[L, R]$ . As it turns out, in both instances,  $L$  must have an atom of weight  $p$  at 0 and  $R$  must have an atom of weight  $p$  at 1. So, as far as the distribution of  $(L, R)$  is concerned, the differences between open and closed intervals are only cosmetic.

The mean of  $D$  is fixed at  $p$ , but the variance of  $D$  is much less restricted. A trivial upper bound is obtained by considering that  $D \in [0, 1]$ , and thus  $\text{var } D \leq p(1 - p)$ . Equality is achieved by the all-or-nothing method. Random intervals with that much variability in the length look awkward; somehow, one's idea of a random interval includes the notion of a random length with a smooth distribution. If, as  $p \rightarrow 0$ ,  $\text{var } D = o(p^2)$ , then  $D/\mathbf{E}D \rightarrow 1$  in probability by Chebyshev's inequality. Visually, one will notice that all intervals will have about the same length, leading to uninteresting data for applications. The aesthetically most appealing cases occur when  $\text{var } D = \Theta(p^2)$  as  $p \rightarrow 0$ . We say that  $a_n = \Theta(b_n)$  if there are positive constants  $A$  and  $B$  such that  $Ab_n \leq a_n \leq Bb_n$  for all  $n$  large enough.

Next we observe that the distributions of  $L$  and  $R$  must have atoms:  $L$  has an atom of weight  $\geq p$  at 0 and  $R$  has an atom of weight  $\geq p$  at 1. To prove this, assume that  $L$  has an atom of size  $q$  at 0. Then let  $R'$  be the random variable equal to  $R$  conditional on  $L = 0$ .

$$p = \Pr\{x \in (L, R)\} \geq q\Pr\{R' > x\} \rightarrow q\Pr\{R' > 0\}$$

as  $x \downarrow 0$ . Furthermore,

$$p = \Pr\{x \in (L, R)\} \leq q\Pr\{R' > x\} + \Pr\{0 < L < x\} \rightarrow q\Pr\{R' > 0\}$$

as  $x \downarrow 0$ . Thus  $\Pr\{R' > 0\} = p/q$ . This implies that  $q \geq p$ . Also,

$$\Pr\{R = 0\} = \Pr\{L = R = 0\} = \Pr\{L = 0\}\Pr\{R' = 0\} = q - p .$$

It is also futile to consider constant length intervals. Indeed, if  $0 \leq L < R \leq 1$  is such that  $R - L = c \in (0, 1)$  for some constant  $c$ , then

$$\sup_{x \in (0,1)} \Pr\{x \in (L, R)\} > \inf_{x \in (0,1)} \Pr\{x \in (L, R)\} .$$

We argue by contradiction, assuming a uniform coverage probability of  $p$ . We have seen earlier that  $q \stackrel{\text{def}}{=} \Pr\{L = 0\} \geq p$ . Because  $\Pr\{c \in (L, R)\} = p$ , there exists a small  $\epsilon > 0$  such that

$$\Pr\{0 < L < c - \epsilon < c < R\} > p/2 .$$

But then

$$\Pr\{c - \epsilon \in (L, R)\} \geq \Pr\{0 < L < c - \epsilon < c < R\} + \Pr\{L = 0\} \geq \frac{3p}{2} ,$$

which is a contradiction.

It is noteworthy that it is impossible to generate an interval's midpoint independently of its length:  $R - L$  and  $(R + L)/2$  cannot be independent. The proof is omitted.

If intervals are generated as a first step in the generation of a random interval graph, then there are simple relations between  $n$ ,  $p$ , and the expected number of edges in the interval graph. We introduce the notion of a regular interval  $(L, R)$ : It has the property that  $\Pr\{L = R\} = 0$  and that no atom of  $L$  has weight more than  $p$ . All of the interval-generating schemes given are regular. The all-or-nothing method yields nonregular intervals. If  $N$  is the number of intersections in the interval graph induced by  $n$  random iid intervals with uniform coverage probability  $p$ , and if the intervals are regular, then  $\mathbf{E}N = \Theta(pn^2)$ . In fact,

$$p(1 - 2p) \binom{n}{2} \leq \mathbf{E}N \leq 4p \binom{n}{2} .$$

The proof is omitted.

### 3. METHODS

#### 3.1 METHOD NUMBER 1: UNWRAPPING THE CIRCLE

We begin with a simple method obtained by considering the interval as wrapped around the perimeter of a circle. The open-interval and closed-interval versions of the algorithm follow. Only the open intervals are analyzed.

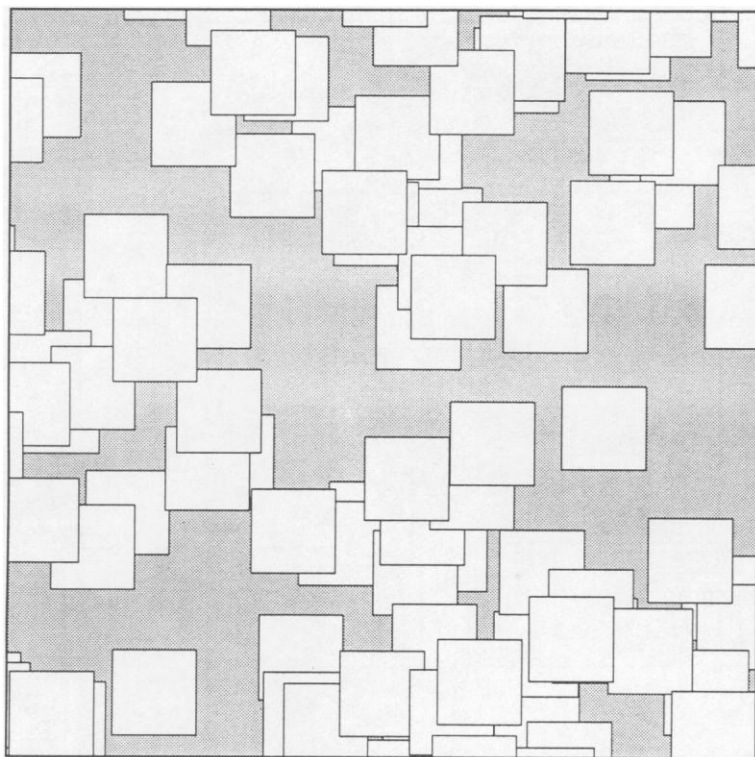


Figure 1. The Unwrap-the-Circle Method.

**unwrap-the-circle method: closed cover**

input parameter:  $z \in [0, 1]$   
 generate  $U$  uniformly on  $[-z, 1]$   
 return  $[0, 1] \cap [U, U + z]$

**unwrap-the-circle method: open cover**

input parameter:  $z \in [0, 1]$   
 generate  $U$  uniformly on  $[-z, 1]$   
 return  $(0, 1) \cap (U, U + z)$

The coverage probability  $p$  is trivially related to the parameter  $z$ : For fixed  $z > 0$ ,

$$\Pr\{x \in (L, R)\} = \frac{z}{z + 1}$$

for all  $x \in (0, 1)$ . To prove this we note that every point  $x \in (0, 1)$  has probability  $z/(z + 1)$  of being covered by the open random intervals (this is certainly true before the intersection with  $[0, 1]$  and remains true after the intersection).

This method has the advantage that all of the intervals, except those near the endpoints, have equal length. By choosing  $z$ , the coverage probability can be any number

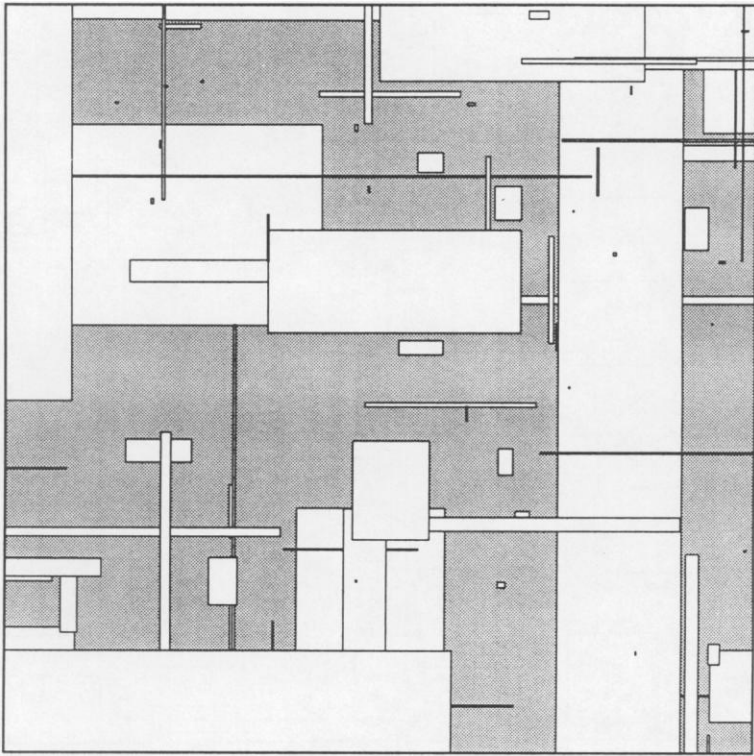


Figure 2. The Random Version of Unwrap-the-Circle Method.  $a = .3$ ,  $b = a(1 - p)/p$ .

between 0 and 1. It is easy to see that  $D/ED \rightarrow 1$  in probability as  $p \rightarrow 0$ . Sometimes more control is desired over the lengths. Nevertheless, the collection of generated intervals does not appear very natural. For example, nested intervals occur only when one of the endpoints is zero or one. Note also that  $M$  is not uniformly distributed on  $[0, 1]$ .

**3.2 METHOD NUMBER 2: A RANDOMIZED UNWRAP-THE-CIRCLE METHOD**

We replace  $z$  in the previous section by a random variable  $Z \geq 0$ . Given  $Z$ , the expected length of a generated interval is  $Z/(Z + 1)$ . Thus we need to choose the distribution of  $Z$  such that

$$E \left\{ \frac{Z}{Z + 1} \right\} = p .$$

A brief example might illustrate this: Let  $Z$  be beta of the second kind with parameters  $a$  and  $b$ ; that is, with density

$$f(z) = \frac{z^{a-1}}{B(a, b)(1 + z)^{a+b}} , z > 0 .$$

Then  $Z/(1 + Z)$  is beta  $(a, b)$ , with mean  $a/(a + b) = p$ . The parameters  $a$  and  $b$  can be picked to achieve equality here. For example, for any  $a$ , pick  $b = a(1 - p)/p$ . A beta

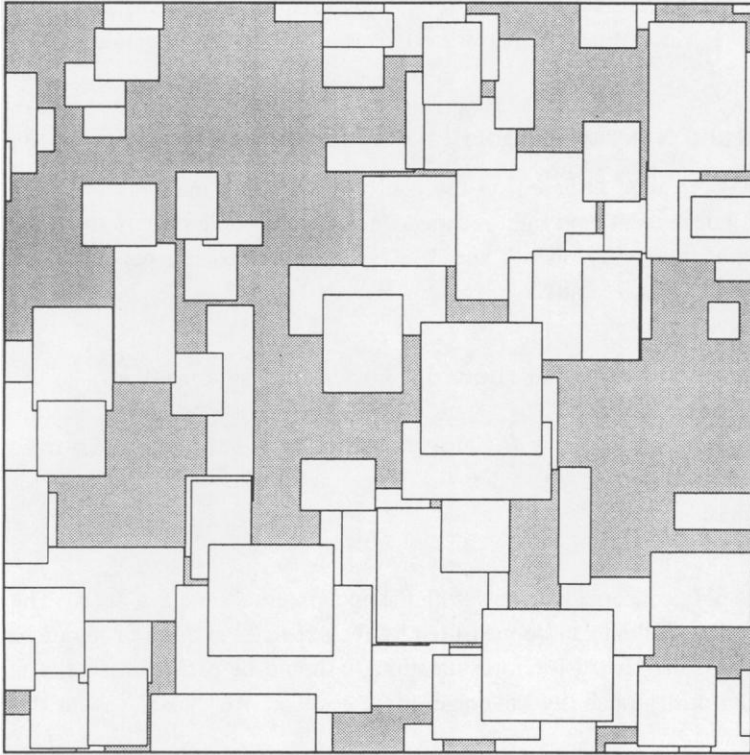


Figure 3. The Random Version of Unwrap-the-Circle Method.  $a = 8$ ,  $b = a(1 - p)/p$ .

variate  $B$  can be generated by the methods described in Devroye (1986), Cheng (1978), or Schmeiser and Babu (1980), and a beta of the second kind is obtained as  $B/(1 - B)$ . In general, if  $Z$  has distribution function  $F$ , then  $D$  has distribution function given by

$$\Pr\{D \leq x\} = F(x) + 2xE\{(Z + 1)^{-1}I_{Z > x}\} = F(x) + 2x \int_x^\infty \frac{1}{y + 1} F(dy).$$

In general, if a user specifies a desirable distribution for  $D$ , then one has to solve this equation for  $F$ , provided that a solution exists.

The freedom obtained by looking at a family of distributions for  $Z$ , such as the beta family suggested earlier, can be used to specify a desired variance for  $D$  if it is feasible (i.e.,  $\text{var } D \leq p(1 - p)$ ). The parameters in the family could be selected to achieve the given variance. For the beta family given previously, take the situation that  $p \rightarrow 0$ , while  $a/(a + b) = p$ . Thus  $b \sim a/p$ . Straightforward variance calculations show that if we take  $a \sim p^\gamma$  for a fixed constant  $\gamma$ , then several asymptotic behaviors may arise:

- $\gamma \leq -1$ :  $\text{var } D = \Theta(p^3)$ ; the intervals are of nearly constant size.
- $-1 \leq \gamma < 0$ :  $\text{var } D = \Theta(p^{2-\gamma})$ ; the intervals are of nearly constant size.
- $\gamma = 0$ :  $\text{var } D = \Theta(p^2)$ ; the intervals are of intermediate variability: the oscillations are of the same order of magnitude as the interval sizes.



- $\gamma \in (0, 1)$ :  $\text{var } D/p^2 \rightarrow \infty$ ; the interval sizes are highly variable.
- $\gamma = 1$ :  $\text{var } D = \Theta(p)$ ; the interval sizes are extremely variable.

### 3.3 METHOD NUMBER 3: THE REJECTION METHOD

The next method is based on the goal of generating intervals in which  $L$  and  $R$  are nearly independent—outright independence is achievable only in the trivial case that  $L = R$  with probability one.  $L$  and  $1 - R$  are generated independently with a given distribution function  $F$  until  $L < R$ .

rejection method

```

given is a distribution function  $F$  on  $[0, 1]$ 
repeat
  generate  $L, 1 - R$  independently with distribution
  function  $F$ 
until  $L < R$ 
return  $(L, R)$ 
    
```

We are only concerned for now with the open intervals and  $x \in (0, 1)$ . The efficiency of the rejection method can be measured by the expected number of iterations ( $N$ ) until we obtain  $L < R$ . The distribution function  $F$  should be picked in such a manner that we have uniform probability coverage at the level  $p$ . We chose  $F$  with the following properties:

- $F$  is continuous and nondecreasing on  $(0, 1)$ .
- $F(0) = q \stackrel{\text{def}}{=} e^{1-1/p}$ .
- For  $x \in (0, 1/2]$ ,  $F(x)F(1 - x) = q$ . (This implies that  $F(1) = 1$  and that  $F(1/2) = \sqrt{q}$ .)

With this choice we claim that the rejection method is valid (i.e.,  $\Pr\{x \in (L, R)\} = p$  for all  $x \in (0, 1)$ ), and that  $\mathbf{E}N = p/q$ . If  $L$  and  $1 - R$  are independent with distribution function  $F$ , then

$$\Pr\{x \in (L, R)\} = F(x)F(1 - x) = q$$

for all  $x \in (0, 1)$ . Also, the probability of a successful pair in a particular loop of the iteration is given by

$$\begin{aligned}
 \Pr\{L < R\} &= \int_0^1 F(1 - x) dF(x) \\
 &= \int_{0 < x \leq 1} \frac{q}{F(x)} dF(x) + F(0) \\
 &= \int_q^1 \frac{q}{y} dy + q \\
 &= q \log(1/q) + q \\
 &= \frac{q}{p}.
 \end{aligned}$$

Thus

$$\Pr\{x \in (L, R) | L < R\} = \frac{\Pr\{x \in (L, R), L < R\}}{\Pr\{L < R\}} = \frac{\Pr\{x \in (L, R)\}}{\Pr\{L < R\}} = p.$$

Finally, by results on rejection algorithms, the expected number of iterations is  $\mathbf{E}N = 1/\Pr\{L < R\} = p/q$ .

The choice of  $F$  is trivial: Take  $F(0) = q$ , and let  $F$  increase in a continuous manner until it reaches the value  $\sqrt{q}$  at  $x = 1/2$ . Then define  $F$  on  $(1/2, 1)$  by continuation based on the equation  $F(1-x) = q/F(x)$ .

### 3.4 METHOD NUMBER 4: THE CONDITIONAL METHOD

Assume that  $F$  is a distribution function as in the previous section. We will show that the following loopless algorithm is valid.

**conditional method**

```

given is a distribution function  $F$  on  $[0, 1]$ 
generate  $U$  and  $V$  independently and uniformly on  $[0, 1]$ 
generate  $L$ :
  if  $U \leq p$  then  $L \leftarrow 0$ 
    else  $L \leftarrow F^{\text{inv}}(qe^{(U-p)/p}) = F^{\text{inv}}(e^{(U-1)/p})$ 
define  $R \leftarrow 1 - F^{\text{inv}}(VF(1-L))$ 
return  $(L, R)$ 

```

We call this the conditional method because we first generate  $L$  and then generate  $R$  conditional on  $L$ . In multivariate random variate generation, a survey of applications of this methodology can be found in Johnson (1987) or in chapter XI of Devroye (1986). We first show that the algorithm is correct. Return to the method of the previous section, and let  $L$  and  $R$  be as given there.

$$\begin{aligned} \Pr\{L \leq x | L < R\} &= \frac{\Pr\{L = 0\} + \int_{0 < y \leq x} \Pr\{R \geq y | L = y\} dF(y)}{\Pr\{L < R\}} \\ &= \frac{q + \int_{0 < y \leq x} F(1-y) dF(y)}{q/p} \\ &= \frac{q + \int_{0 < y \leq x} q/F(y) dF(y)}{q/p} \\ &= p + p \log(F(x)/F(0)) \\ &= p + p \log\left(\frac{F(x)}{q}\right), \quad x \in [0, 1]. \end{aligned}$$

Call the latter distribution function  $G$ . The recipe given for  $L$  in the conditional method algorithm corresponds to using the inversion method for  $G$ . Given  $L$ , the distribution

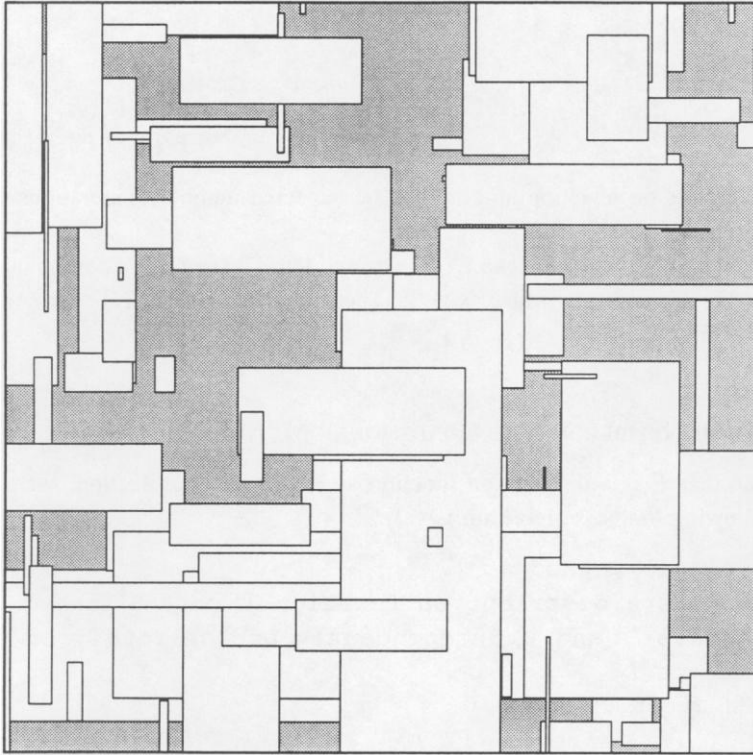


Figure 4. The Conditional Method.

function of  $1 - R$  is  $F$ , restricted to  $[0, 1 - L]$ :

$$\Pr\{1 - R \leq y | L\} = \frac{F(y)}{F(1 - L)}, \quad 0 \leq y \leq 1 - L .$$

The recipe for  $R$  in the algorithm given previously corresponds once again to the inversion method. The uniform probability coverage follows from the correctness of the algorithm of the previous section.

One particular example stands out because of its simplicity. Consider the distribution function  $F(x) = q^{1-x}$  in the rejection method algorithm. It satisfies all the necessary conditions as  $F(0) = q$ ,  $F$  is nondecreasing, and  $F(x)F(1 - x) = q$  for all  $x \in [0, 1]$ . In the conditional method algorithm, we obtained the distribution function for  $L$  as

$$G(x) = p + p \log(F(x)/q) = p + x(1 - p) .$$

This is easily seen to be a mixture of an atom of size  $p$  at the origin and a uniform density on  $[0, 1]$  carrying weight  $1 - p$ . The distribution function of  $(1 - R)/(1 - L)$  given  $L$  is

$$\Pr\left\{\frac{1 - R}{1 - L} \leq x | L\right\} = q^{(1-L)(1-x)}, \quad 0 \leq x \leq 1 .$$

This is a mixture of an atom of size  $q^{1-L}$  at the origin and an exponentially increasing density on  $(0, 1)$ . Using the fact that minus the logarithm of a uniform  $[0, 1]$  random variable is exponentially distributed, we see that, given  $L$ ,  $R$  is distributed as

$$L + \min\left(1 - L, \frac{E}{1/p - 1}\right),$$

where  $E$  is a standard exponential random variate. The conditional method is thus simplified as follows:

```
conditional method (worked out example)
  generate  $U$  and  $V$  independently and uniformly on  $[0, 1]$ 
  if  $U \leq p$  then  $L \leftarrow 0$  else  $L \leftarrow V$ 
  generate an exponential random variate  $E$ 
   $R \leftarrow L + \min(1 - L, E/(1/p - 1))$ 
  return  $(L, R)$ 
```

### 3.5 METHOD NUMBER 5: THE ORDER STATISTICS METHOD

Epstein and Sack (1992) proposed the following method: If  $p = 1/n$  for some integer  $n$ , then partition  $[0, 1]$  into  $n$  intervals induced by a random iid sample of size  $n-1$  drawn from the uniform distribution. Pick one of these intervals uniformly at random. Then, for any fixed  $x \in (0, 1)$ , the coverage probability is exactly  $p$ . If an ordered sample is generated directly in order (see Devroye 1986, chapter V), then this method takes  $O(n)$  time. This is a special case of a more general paradigm given in the following.

the partitioning method

```
generate an ordered sample  $0 = X_0 < X_1 < \dots < X_{n-1} < X_n = 1$ 
  (from any distribution, regardless of dependence)
generate  $Z$  uniformly in  $\{0, 1, \dots, n-1\}$ 
return  $(L, R) \leftarrow (X_Z, X_{Z+1})$ 
```

Among all possible distributions for the  $X_i$ 's, the uniform distribution occupies a special place, because the properties of its spacings are well understood. The length  $R - L$  is beta  $(1, n-1)$  distributed, with mean  $1/n$ . Because the  $k$ th smallest uniform order statistic in a sample of size  $n-1$  is beta  $(k, n-k)$  distributed, we can obtain an  $O(1)$  expected-time version of the algorithm, provided that all of the beta variates are generated by a uniformly fast algorithm (Cheng 1978; Devroye 1986; Schmeiser and Babu 1980).

the order statistics method

```
generate  $Z$  uniformly in  $\{0, 1, \dots, n-1\}$ 
if  $Z = 0$  then  $L \leftarrow 0$ 
  else  $L \leftarrow \text{beta}(Z, n-Z)$ 
if  $Z = n-1$  then  $W \leftarrow 1$ 
  else  $W \leftarrow \text{beta}(1, n-Z-1)$ 
define  $R \leftarrow L + (1-L)W$ 
return  $(L, R)$ 
```

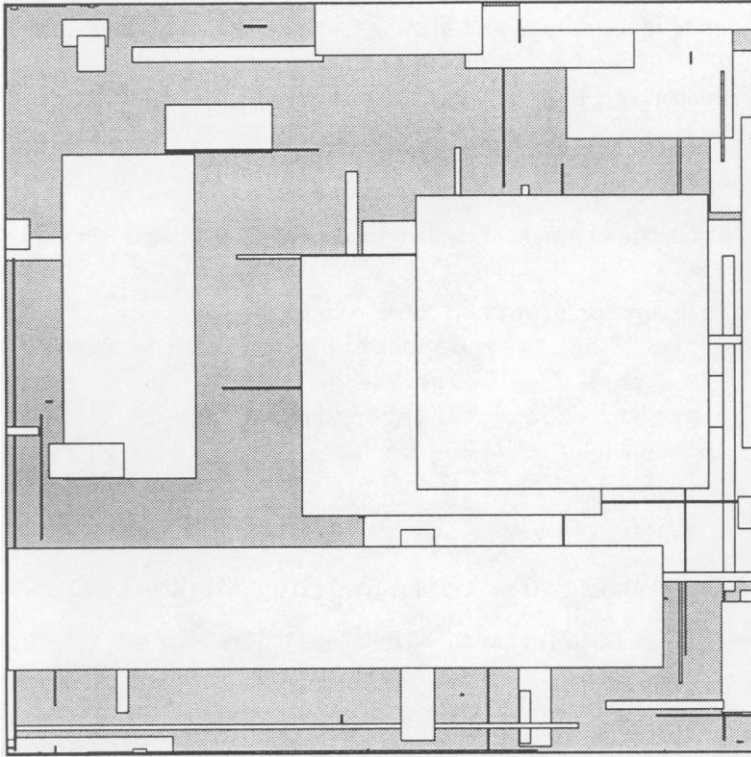


Figure 5. *The Dirichlet Process Method.  $\nu = 2$ .*

Thus far we can only achieve coverage probabilities that are  $1/n$  for some integer  $n$ . To get around this we replace  $n - 1$  by a Poisson random variate  $N$  with parameter  $\lambda > 0$ ; that is,

$$\Pr\{N = i\} = \frac{\lambda^i}{i!} e^{-\lambda}, \quad i \geq 0.$$

This can be done in constant expected time by algorithms such as those given in Ahrens and Dieter (1980, 1982, and 1987), Ahrens, Kohrt, and Dieter (1983), Devroye (1981, 1987), Pokhodzei (1984), Schmeiser and Kachitvichyanukul (1981), or Stadlober (1988). For fixed  $x \in (0, 1)$ , we have

$$\Pr\{x \in (L, R)\} = \mathbf{E} \left\{ \frac{1}{N+1} \right\} = \frac{1 - e^{-\lambda}}{\lambda}.$$

By varying  $\lambda$  from 0 to  $\infty$ , any coverage probability between 1 and 0 can be obtained in this manner. If the coverage probability  $p$  is given, we have to solve the equation  $p = (1 - e^{-\lambda})/\lambda$ . For large  $\lambda$ ,  $N$  is very concentrated around  $\lambda$ , so that the resulting intervals behave roughly the same as those of the raw-order statistics method. Unfortunately, this method does not have any flexibility with regard to the distribution of either  $L$  or the interval length  $D$ .

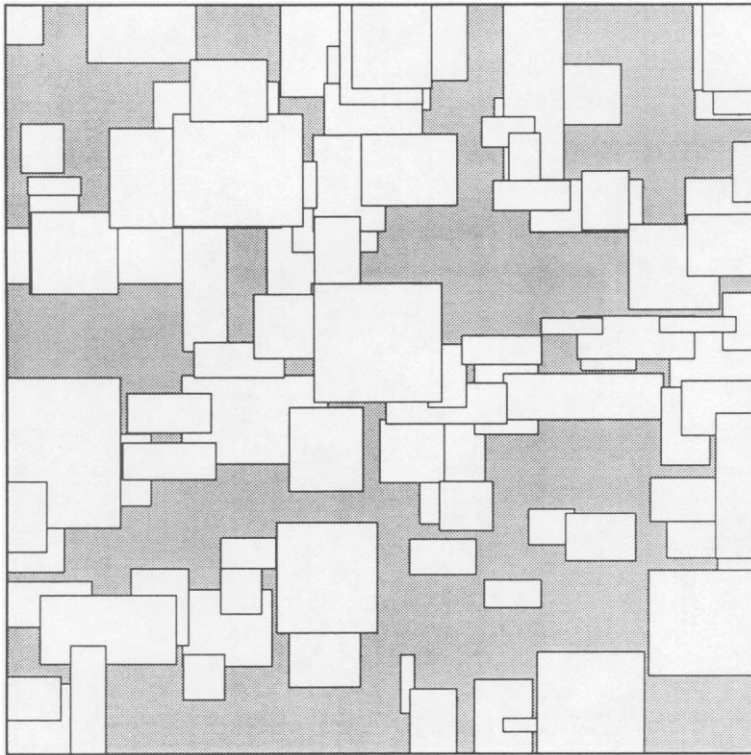


Figure 6. The Dirichlet Process Method.  $v = 4$ .

Assume that  $n - 1$  is replaced by a negative binomial random variable  $N$  with parameter 2 defined by

$$\Pr\{N = i\} = (i + 1)(1 - q)^2 q^i \quad (i \geq 0) .$$

This choice is convenient because the coverage probability ( $p$ ) is easily seen to be

$$\mathbf{E} \left\{ \frac{1}{N + 1} \right\} = \sum_{i=0}^{\infty} (1 - q)^2 q^i = 1 - q \stackrel{\text{def}}{=} p .$$

Also,  $N$  is easily generated as the sum of two independent geometric random variables, each of which is obtainable as  $\lfloor -E/\log(q) \rfloor$  (Devroye 1986).

### 3.6 METHOD NUMBER 6: THE DIRICHLET PROCESS METHOD

One may control the variability of the intervals much better by using properties of Dirichlet processes (see Wilks 1962; Aitchison 1963; Basu and Tiwari 1982; Devroye 1986, chapter XIV.4; Narayanan 1990). The partition of the previous section follows a Dirichlet distribution, in which each spacing is beta ( $v, v(n - 1)$ ) distributed, and  $v > 0$  is a variability parameter.

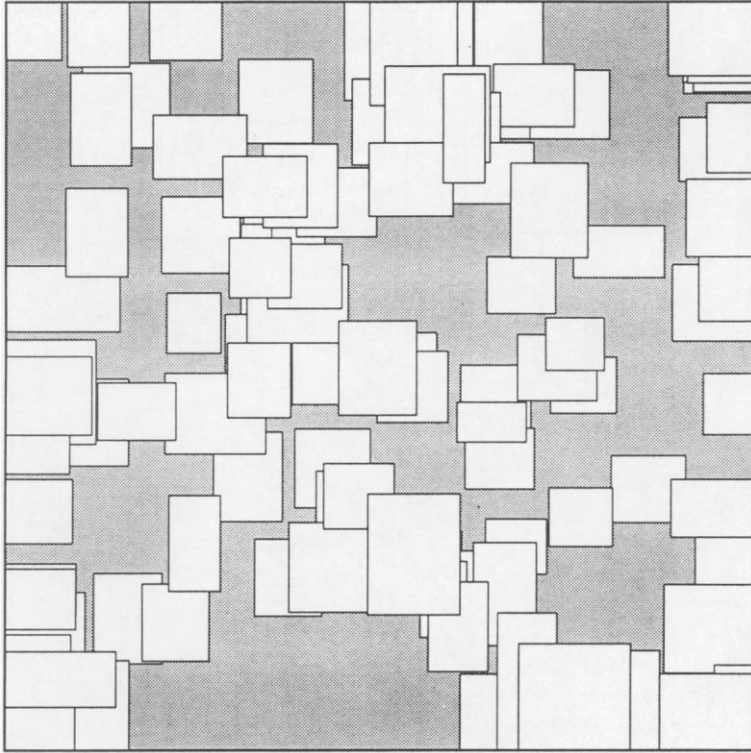


Figure 7. *The Dirichlet Process Method.  $v = 20$ .*

the Dirichlet process method

```

generate  $Z$  uniformly in  $\{0, 1, \dots, n-1\}$ 
if  $Z = 0$  then  $L \leftarrow 0$ 
    else  $L \leftarrow \text{beta}(vZ, v(n-Z))$ 
if  $Z = n-1$  then  $W \leftarrow 1$ 
    else  $W \leftarrow \text{beta}(v, v(n-Z-1))$ 
define  $R \leftarrow L + (1-L)W$ 
return  $(L, R)$ 

```

We see that the length  $D$  is beta  $(v, v(n-1))$ , so that in view of  $ED = p$ , we must have  $p = 1/n$ . Interestingly, we have the following simple expression for the variance:

$$\text{var}\{D\} = \frac{n-1}{n^2(vn+1)} = \frac{p(1-p)}{1+v/p}.$$

This restricts the design slightly as  $1/p$  must be an integer. By making  $v$  as a function of  $p$ , however, we obtain virtually unlimited freedom in the choice of  $\text{var}\{D\}$ : Any variance in the allowable range  $(0, p(1-p))$  is achievable by letting  $v$  vary from  $\infty$  down to 0. When  $p \rightarrow 0$  and  $v \rightarrow \infty$ , then  $D/ED \rightarrow 1$  in probability, as  $\text{var}\{D\} \sim p^2/v$ , leading to low variability. The case  $v = 1$  corresponds to the order-statistics method. For  $v$  constant,  $\text{var}\{D\} = \Theta(p^2)$ , which is the most interesting case, as  $(n-1)D$  and

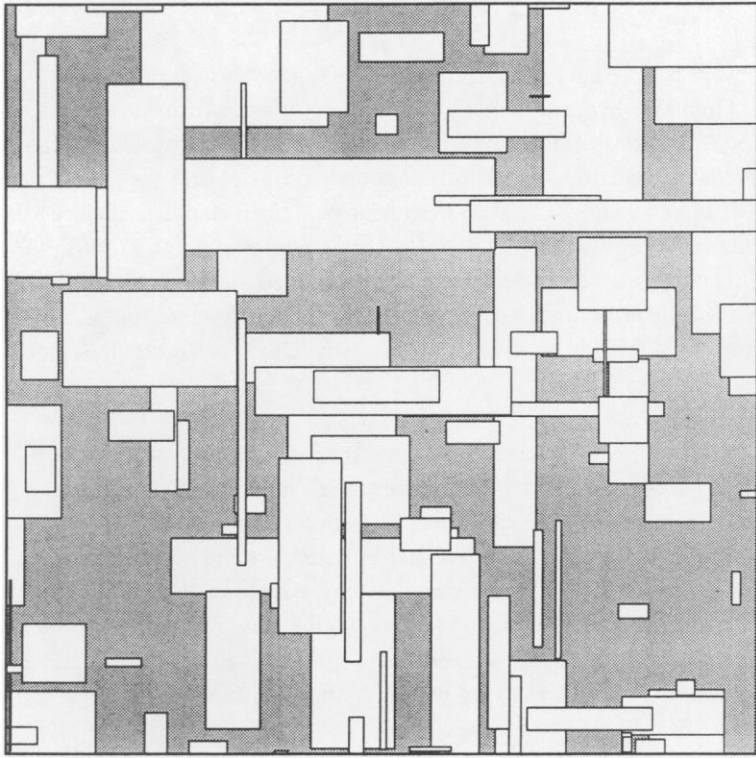


Figure 8. The Dirichlet Process Method.  $\nu = 1$ .

$D/ED$  both tend in distribution to a gamma  $(\nu, 1/\nu)$  random variable. For small  $\nu$ , the variance is gigantic: If  $\nu \rightarrow 0$  as  $p \rightarrow 0$ , then we have  $\text{var}\{D\}/E^2D \rightarrow \infty$ .

For general  $p$ , we may define  $n = \lfloor 1/p \rfloor$ , and use the order-statistics method or the Dirichlet process method in which  $n$  is used with probability  $q = n(n+1)p - n$ , and  $n+1$  is used otherwise. Observe that the coverage probability is

$$\Pr\{N = i\} = \frac{n(n+1)p - n}{n} + \frac{1 + n - n(n+1)p}{n+1} = p.$$

This manner of mixing keeps the variability of the interval sizes to a minimum.

#### 4. RANDOM SQUARES

Random squares in the plane pose a serious challenge. Suppose that we were to generate a uniform coverage random interval of  $(0, 1)$  with coverage probability  $\sqrt{p}$ . Assume that its length is  $D$ . This would fix the size of the square at  $D \times D$ . To fix the position of the square with respect to the other coordinate, we would need yet another random interval, with uniform coverage probability  $\sqrt{p}$  and length  $D$ . But we know that fixed-length intervals do not lead to uniform coverage probabilities. Also, we cannot play very much with the length of  $D$  because  $ED = \sqrt{p}$ . This leaves us nowhere, and an entirely new approach is necessary for random squares.



## 5. EXPERIMENTS

As we were mostly interested in computational geometric applications, we generated a number of random hyperrectangles with various methods. In all cases we wanted a uniform coverage probability of  $p = 1/100$ . This is achieved by demanding a uniform coverage probability of  $p' = 1/10$  for  $x$  intervals on  $[0, 1]$  and for  $y$  intervals on  $[0, 1]$  separately. If we choose  $n = 100$  independent random rectangles, then, as the number of rectangles covering any point  $x \in (0, 1)^2$  is binomial  $(100, 1/100)$ , the expected number of rectangles covering any point is just one. Also, the expected uncovered area is precisely the probability that a given  $x \in (0, 1)^2$  is not covered; that is,  $(1 - 1/100)^{100}$ , which is close to  $1/e$ . This provides a quick visual check of correctness, but it is also a nice calibrator in comparisons.

Four methods were tested:

1. The unwrap-the-circle method. This shows that, except near the borders, all the rectangles are in fact equisized squares. See Figure 1 (p. 295).
2. The randomized version of unwrap-the-circle, in which  $Z$  is a beta random variable of the second kind with parameters  $a$  and  $b = a(1-p)/p$ . The first parameter was varied; for  $a = .3$ , the rectangles vary widely in size. This is reduced when  $a = 1$ , while for  $a = 8$  and up the rectangles have comparable sizes and tend to be more and more square shaped. See Figures 2 and 3 (pp. 296 and 297).
3. The conditional method based on the distribution function  $F(x) = q^{1-x}$  given in the text. See Figure 4 (p. 300).
4. The Dirichlet process method. The size  $n$  in the order-statistics method is picked automatically as described in Subsection 3.6. The variability parameter  $v$  is changed from  $v = .2$  (high variability in size and shape of the rectangles), to  $v = .5$ ,  $v = 4$ , and  $v = 20$ . For  $v = 1$  we obtain the standard order-statistics method. For the higher values of  $v$ , the rectangles become more and more like squares, and their sizes become increasingly similar. See Figures 5–8 (pp. 302–305).

The methods were coded directly in the Postscript language and sent to the printer. This means that the Postscript laser printer itself does all the work, from generating the necessary random variates to graphics computations. The host computer does no compilation or interpretation. This approach leads to short, fresh code, and can be carried out without any compilers or graphics packages. Also, one has the full power of the Postscript language at one's fingertips while keeping things simple. This often leads to striking displays not otherwise possible through standard software. For a beautiful introduction to Postscript, start with McGilton and Campione (1992).

*[Received August 1992. Revised May 1993]*

## ACKNOWLEDGMENTS

The authors' research was sponsored by NSERC Grants A3456 and OGP322, and FCAR Grant 90-ER-0291.

## REFERENCES

- Ahrens, J. H., and Dieter, U. (1980), "Sampling From Binomial and Poisson Distributions: A Method With Bounded Computation Times," *Computing*, 25, 193–208.
- (1982), "Computer Generation of Poisson Deviates From Modified Normal Distributions," *ACM Transactions on Mathematical Software*, 8, 163–179.
- (1987), "A Convenient Sampling Method With Bounded Computation Times for Poisson Distributions," in *First International Conference on Statistical Computation, Izmir, Turkey*, 4–17.
- Ahrens, J. H., Kohrt, K. D., and Dieter, U. (1983), "Algorithm 599. Sampling From Gamma and Poisson Distributions," *ACM Transactions on Mathematical Software*, 9, 255–257.
- Aitchison, J. (1963), "Inverse Distributions and Independent Gamma-distributed Products of Random Variables," *Biometrika*, 50, 505–508.
- Basu, D. and Tiwari, R. C. (1982), "A Note on the Dirichlet Process," in *Statistics and Probability: Essays in Honor of C. R. Rao*, eds. G. Kallianpur, P. R. Krishnaiah, and J. K. Ghosh, Amsterdam: North-Holland, pp. 89–103.
- Cheng, R. C. H. (1978), "Generating Beta Variates With Nonintegral Shape Parameters," *Communications of the ACM*, 21, 317–322.
- Devroye, L. (1981), "The Computer Generation of Poisson Random Variables," *Computing*, 26, 197–207.
- (1986), *Non-Uniform Random Variate Generation*, New York: Springer-Verlag.
- (1987), "A Simple Generator for Discrete Log-concave Distributions," *Computing*, 39, 87–91.
- Epstein, P., and Sack, J.-R. (1992), "Generating Triangulations and Intervals at Random," Technical Report, Carleton University, School of Computer Science.
- Johnson, M. E. (1987), *Multivariate Statistical Simulation*, New York: John Wiley.
- McGilton, H., and Campione, M. (1992), *PostScript by Example*, Reading, MA: Addison-Wesley.
- Narayanan, A. (1990), "Computer Generation of Dirichlet Random Vectors," *Journal of Statistical Computation and Simulation*, 36, 19–30.
- Pokhodzei, B. B. (1984), "Beta- and Gamma-methods of Modelling Binomial and Poisson Distributions," *USSR Computational Mathematics and Mathematical Physics*, 24, 114–118.
- Scheinerman, E. R. (1988), "Random Interval Graphs," *Combinatorica*, 8, 357–371.
- (1990), "An Evolution of Interval Graphs," *Discrete Mathematics*, 82, 287–302.
- Schmeiser, B. W., and Babu, A. J. G. (1980), "Beta Variate Generation via Exponential Majorizing Functions," *Operations Research*, 28, 917–926.
- Schmeiser, B. W., and Kachitvichyanukul, V. (1981), "Poisson Random Variate Generation," Research Memorandum 81-4, Purdue University, School of Industrial Engineering.
- Stadlober, E. (1988), "Sampling from Poisson, Binomial and Hypergeometric Distributions: Ratio of Uniforms as a Simple Fast Alternative," Habilitationsschrift, Institute of Statistics, Technical University of Graz, Austria.
- Wilks, S. S. (1962), *Mathematical Statistics*, New York: John Wiley.

## LINKED CITATIONS

- Page 1 of 1 -



You have printed the following article:

### **On Generating Random Intervals and Hyperrectangles**

Luc Devroye; Peter Epstein; Jörg-Rüdiger Sack

*Journal of Computational and Graphical Statistics*, Vol. 2, No. 3. (Sep., 1993), pp. 291-307.

Stable URL:

<http://links.jstor.org/sici?sici=1061-8600%28199309%292%3A3%3C291%3AOGRIAH%3E2.0.CO%3B2-G>

---

*This article references the following linked citations. If you are trying to access articles from an off-campus location, you may be required to first logon via your library web site to access JSTOR. Please visit your library's website or contact a librarian to learn about options for remote access to JSTOR.*

## **References**

### **Inverse Distributions and Independent Gamma-Distributed Products of Random Variables**

J. Aitchison

*Biometrika*, Vol. 50, No. 3/4. (Dec., 1963), pp. 505-508.

Stable URL:

<http://links.jstor.org/sici?sici=0006-3444%28196312%2950%3A3%2F4%3C505%3AIDAIGP%3E2.0.CO%3B2-O>

### **Beta Variate Generation via Exponential Majorizing Functions**

Bruce W. Schmeiser; A. J. G. Babu

*Operations Research*, Vol. 28, No. 4. (Jul. - Aug., 1980), pp. 917-926.

Stable URL:

<http://links.jstor.org/sici?sici=0030-364X%28198007%2F08%2928%3A4%3C917%3ABVGVEM%3E2.0.CO%3B2-K>