

# TransUnif

Prof. Bruce E. Trumbo

## Transformations of Standard Uniform Distributions

We have seen that the **R** function `runif` uses a random number generator to simulate a sample from the standard uniform distribution  $UNIF(0,1)$ . All of our simulations use standard uniform random variables or are based on transforming such random variables to obtain other distributions of interest. Included in the **R** language are some functions that implement suitable transformations. For example, `rnorm`, `rexp`, `rbeta`, and `rbinom` simulate samples from normal, exponential, beta, and binomial distributions, respectively. Also, the function `sample` is based on simulated realizations of  $UNIF(0,1)$ .

A systematic study of the programming methods required to transform uniform distributions into other commonly used distributions involves technical details beyond the scope of this book. (For a more extensive treatment, see Chapter 3 of Fishman (1996).) However, if you are going to do simulations and trust the results, we feel you should have some idea how such transformations are accomplished—at least in a few familiar and elementary cases. The purpose of this section is to provide some of the basic theory and a few simple examples of transformations from uniform distributions to other familiar distributions. Also, this discussion provides the opportunity for a brief review of some distributions we will use later on.

### EXAMPLE 1

A real function (transformation) of a random variable is again a random variable. For example, if  $U \sim UNIF(0,1)$ , then the linear function  $X = g(U) = 4U + 2$  is a random variable uniformly distributed on the interval  $(2,6)$ . That is,  $X \sim UNIF(2,6)$ . The transformation  $g$  stretches the distribution of  $U$  by a factor of 4 and then shifts it two units to the right. Recalling that  $F_U(u) = P(U \leq u) = u$ , for  $0 < u < 1$ , we have the following formal demonstration. For  $2 < x < 6$ ,

$$F_X(x) = P(X \leq x) = P(g(U) \leq x) = P(4U + 2 \leq x) = P(g^{-1}(X) \leq g^{-1}(x))$$

$$= P(U \leq (x - 2)/4) = (x - 2)/4.$$

which is the density function of  $UNIF(2, 6)$ .

In **R**, the second and third parameters of the function `runif` specify the left and right endpoints, respectively, of  $UNIF(\theta_1, \theta_2)$ , the uniform distribution on the interval  $(\theta_1, \theta_2)$ . Thus each of the statements `4 * runif(10) + 2`, `4 * runif(10, 0, 1) + 2`, and `runif(10, 2, 6)` simulates 10 observations from  $UNIF(2, 6)$ . PROBLEM 1 asks you to consider a more general version of this example. ■

In EXAMPLE 1, we have found the CDF of the transformed random variable, and then used the CDF to find its density function. This method works in a large variety of situations. Next, we see that a particular nonlinear transformation of a standard uniform random distribution is a member of the beta family of distributions. We leave the formal demonstration to PROBLEM 2 and use a simulation and graphics to illustrate the effect of the transformation.

## EXAMPLE 2

Suppose  $U \sim UNIF(0, 1)$  and  $X = U$ . Then  $P(0 < X < 1) = 1$ . Also, because the square root of a number in  $(0, 1)$  is larger than the number itself, we know intuitively that the distribution of  $X$  must concentrate its probability toward the right end of  $(0, 1)$ . Specifically, the method of EXAMPLE 1 shows that  $X$  has the cumulative distribution function  $F_X(x) = x^2$ , and the density function  $f_X(x) = 2x$ , for  $0 < x < 1$ . Recall that if  $Y \sim BETA(\alpha, \beta)$  then its density function is

$$f_Y(y) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} y^{\alpha-1} (1 - y)^{\beta-1},$$

for  $0 < y < 1$ , and positive parameters  $\alpha$  and  $\beta$ . Here  $\Gamma$  denotes the gamma function, which has  $\Gamma(n + 1) = n!$  for positive integer  $n$ , and may be evaluated more generally in **R** using `gamma`. Thus  $X = \sqrt{U} \sim BETA(2, 1)$ .

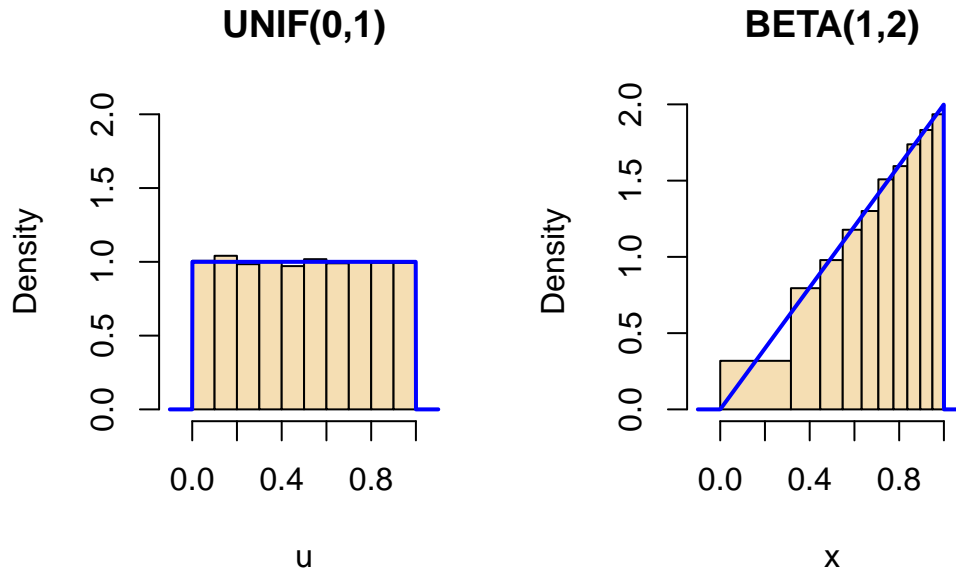
The following simulation shows what happens when one takes the square root of randomly chosen points `u` in the ten intervals  $(0, 0.1]$ ,  $(0.1, 0.2]$ , through  $(0.9, 1)$ . In **R**, the names of density functions of programmed distributions begin with the letter `d`: thus the functions `dunif` and `dbeta` in the code above.

```
set.seed(1212)
m <- 10000
u <- runif(m)
x <- sqrt(u)
```

```

par(mfrow=c(1,2))
hh = seq(-.1, 1.1, length=1000); cutp = seq(0, 1, by = .1)
hist(u, breaks=cutp, prob=T, col="wheat", ylim=c(0,2),
xlim=c(-.1, 1.1), main="UNIF(0,1)")
lines(hh, dunif(hh), col="blue", lwd=2)
hist(x, breaks=sqrt(cutp), prob=T, col="wheat",
xlim=c(-.1, 1.1), main="BETA(1,2)")
lines(hh, dbeta(hh, 2, 1), col="blue", lwd=2)

```



```

par(mfrow=c(1,1))

```

Graphical results are shown in FIGURE A. Each bar in each histogram represents about a thousand points, representing one tenth of the total probability. Density functions of  $UNIF(0, 1)$  and  $BETA(2, 1)$  are superimposed on their respective histograms.

By taking different powers of a standard uniform random variable one can obtain random variables with distributions  $BETA(, 1)$  (see PROBLEM 2). More intricate methods are required to sample from some other members of the distribution family  $BETA(, )$  (see PROBLEM 3). Optimal methods for all cases are available in **R** as the function *rbeta*. Thus either of the statements *sqrtrunif(10)* or *rbeta(10, 2, 1)* could be used to simulate 10 observations from  $BETA(2, 1)$ , but the latter code is more convenient because it can be used for any member of the beta family. ■

Now we summarize what we have seen so far.

- In EXAMPLE 1, the CDF of  $X$  is  $F_X(x) = (x - 2)/4$ , for  $2 < x < 6$ . The inverse of the CDF is called the **quantile function**. Here it is  $F_X^{-1}(u) = 2 + 4u$ , obtained by solving  $F_X(x) = u$  for  $x$  in terms of  $u$ . This is the function  $g$  we used to transform  $U \sim UNIF(0, 1)$  to get the random variable  $X = g(U) \sim UNIF(2, 6)$ .
- In EXAMPLE 2, the CDF  $F_X(x) = x^2$ , is used to obtain  $f_X(x) = 2x$ , for  $0 < x < 1$ . Thus  $X$  has quantile function  $F_X^{-1}(u) = \sqrt{u}$ , which is the function  $g$  used to transform  $U \sim UNIF(0, 1)$  to get the the random variable  $X \sim BETA(2, 1)$ .

Suppose we want to simulate values from a distribution whose quantile function is known. A general principle is that this quantile function is the function  $g$  such that  $X = g(U)$  has the desired distribution, where  $U \sim UNIF(0, 1)$ . Specifically, in the next example, we want to simulate observations  $X \sim EXP(1)$ , the exponential distribution with rate 1. Accordingly, we find the quantile function of  $EXP(1)$  and use it to transform observations from  $UNIF(0, 1)$ .

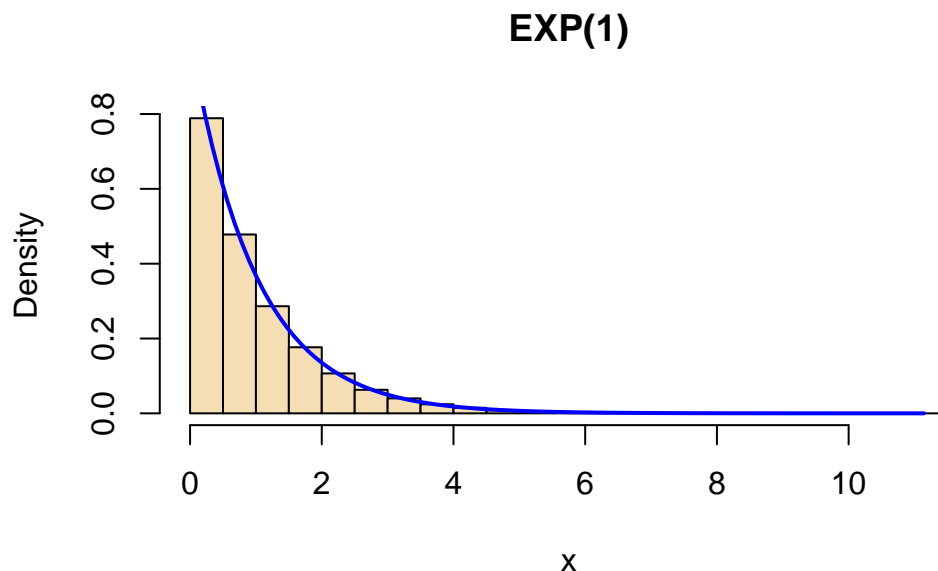
### EXAMPLE 3

Throughout this example let  $x > 0$  and  $0 < u < 1$ . We wish to simulate observations from the distribution  $EXP(1)$ , which has density function  $f(x) = e^{-x}$  and CDF  $F(x) = 1 - e^{-x}$ . Solving  $u = 1 - e^{-x}$  for  $x$  in terms of  $u$ , we have the quantile function  $F^{-1}(u) = -\ln(1 - u)$ . Thus  $X = -\ln(1 - U) \sim EXP(1)$ . Because  $1 - U \sim UNIF(0, 1)$  it is simpler to simulate observations from this exponential distribution as  $X = -\ln U$  (see PROBLEM 1).

The following **R** code demonstrates that a histogram of 100,000 observations generated in this way very nearly fits the density function of  $EXP(1)$ , as seen in FIGURE B. Furthermore, the mean and standard deviation of the simulated values are both nearly 1, which is the mean and standard deviation of the distribution  $EXP(1)$ .

```
set.seed(1234)
m <- 100000
u <- runif(m)
x <- -log(u)

hist(x, prob=T, col="wheat", main="EXP(1)")
xx = seq(0, max(x), length=100)
lines(xx, dexp(xx, 1), col="blue", lwd=2)
```



```
mean(x)
```

```
[1] 0.9988505
```

```
sd(x)
```

```
[1] 0.9984966
```

For most purposes, any of the following statements could be used to sample 10 observations from  $EXP(1)$  :  $-\log(\text{runif}(10))$ ,  $\text{qexp}(\text{runif}(1),1)$ , or  $\text{rexp}(10,1)$ . The second statement works because  $\text{qexp}$  (with second parameter 1) is the quantile function of  $EXP(1)$ . (PROBLEM 4 uses the quantile transformation to sample from  $EXP(1/2)$ .) However, the method using  $\text{rexp}$  is preferable because it uses an algorithm that is technically superior to our log-transform method, especially in its treatment of very large simulated values. ■

So far, all of our examples have dealt with continuous distributions. Now we turn to an example where we sample from a binomial distribution.

#### EXAMPLE 4

According to genetic theory the probability that any one offspring of a particular pair of guinea pigs will have straight hair is  $1/4$ . Suppose we want to simulate births of six offspring. That is, we want to simulate one realization of  $X \sim \text{BINOM}(6, 1/4)$ . One way to do this is to simulate

six observations from  $UNIF(0, 1)$ . The probability that any one of these uniform observations is less than  $1/4$  is  $1/4$ . So  $X$  can be simulated as the sum of six logical variables, where **FALSE** is interpreted as 0 and **TRUE** as 1:  $\text{sum}(\text{runif}(6) < 1/4)$ . The sample function is also programmed to use *runif*. So  $\text{sum}(\text{sample}(c(0, 1), 6, \text{repl} = T, \text{prob} = c(3/4, 1/4)))$  is an equivalent way to simulate  $X$  as a sum.

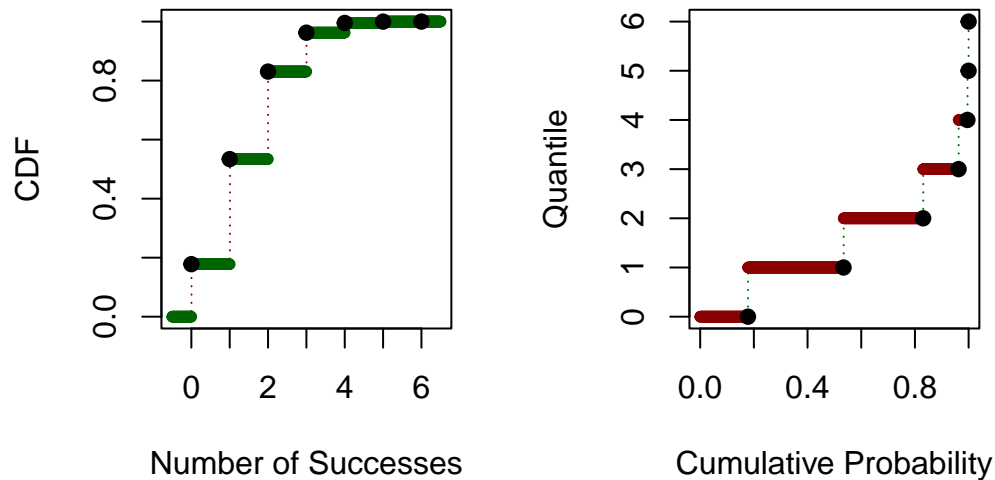
Because **R** defines the quantile function for a discrete random variable in just the right way, one can use the quantile function approach:  $q\text{binom}(\text{runif}(1), 6, 1/4)$ . The second method has the advantage of requiring only one random value from  $UNIF(0, 1)$ , while the first—somewhat wastefully—requires six. In this case, it turns out that the quantile transform method is exactly equivalent to  $r\text{binom}(1, 6, 1/4)$ .

For a discrete random variable  $X$ , **R** defines  $F_X(u)$  as the minimum of the values  $x$  such that  $F_X(x)/geu$ . The left panel of FIGURE C shows the CDF of  $BINOM(6, 1/4)$ , where the vertical reference segments (dotted) represent individual binomial probabilities  $P(X = i)$ ,  $i = 0, 1, \dots, 6$ . The right panel shows the corresponding quantile function, where the horizontal segments of the function (heavy) represent these same probabilities. PROBLEM 5 shows **R** code for a simplified version of this figure. ■

```
par(mfrow=c(1,2))

xx <- seq(-.5, 6.5, length=1000)
plot(xx, pbinom(xx, 6, 1/4), type="s", lty="dotted", col="darkred", xlab="Number of Successes", ylab="Probability", pch=19)
points(xx, pbinom(xx, 6, 1/4), pch=20, col="darkgreen")
points(0:6, pbinom(0:6, 6, 1/4), pch=19)

qq <- seq(0, 1, length=1000)
plot(qq, qbinom(qq, 6, 1/4), type="s", lty="dotted", col="darkgreen", xlab="Cumulative Probability", ylab="Quantile", pch=19)
points(qq, qbinom(qq, 6, 1/4), pch=20, col="darkred")
q <- pbinom(0:6, 6, 1/4)
points(q, qbinom(q, 6, 1/4), pch=19)
```



```
par(mfrow=c(1,1))
```

In practice, when available, it is best to use random functions programmed into R (for example, *rbeta*, *rbeta*, *rbinom*) because they implement algorithms that are fast and accurate. However, some useful distributions are not programmed into the base package of **R**. It may be possible to use the quantile transformation of standard uniform to simulate observations from such a distribution.

### EXAMPLE 5.

The Pareto family of distributions is sometimes useful in economics, actuarial science, geology, and other sciences, but it is not included the base package of **R**. One member of this family has density function  $f(x) = 3/x^4$  and CDF  $F(x) = 1 - x^{-3}$ , for  $x > 1$ ; mean 1.5 and variance 0.75. The following **R** code simulates a sample of 5000 observations from this distribution.

```
set.seed(123)
m <- 5000
kap <- 3

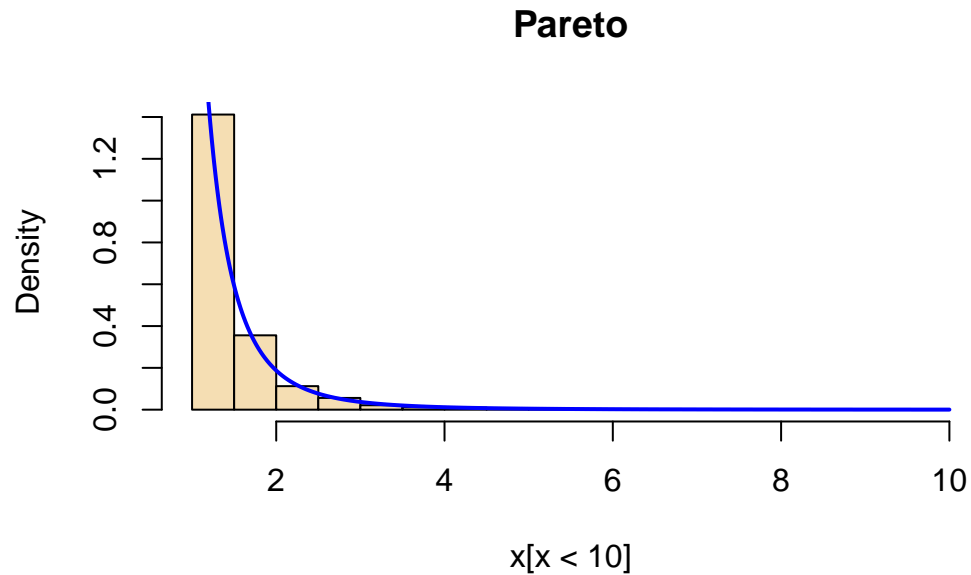
xx <- seq(1, 10, length=1000)
pdf <- kap/xx^(kap+1)
x <- (1 - runif(m))^(1/kap)
mean(x)
```

```
[1] 1.492558
```

```
var(x)
```

```
[1] 0.7048778
```

```
cutp <-seq(0, max(x)+.5, by=.5)  
hist(x[x<10], prob=T, col="wheat", main="Pareto")  
lines(xx, pdf, col="blue", lwd=2)
```



```
mean(x)
```

```
[1] 1.492558
```

```
var(x)
```

```
[1] 0.7048778
```

FIGURE D shows a histogram of the results (except for the six observations that exceed 10) along with the density function. ■



## Transformations Involving Standard Normal Distributions

Normal distributions play an important role in probability and statistics, and so it is important to know how to simulate samples from normal distributions. The **R** function `rnorm` samples from the standard normal distribution. At the end of this section we indicate how to transform standard uniform observations into standard normal ones. In the first example below, we look at some relationships between standard normal and other distributions.

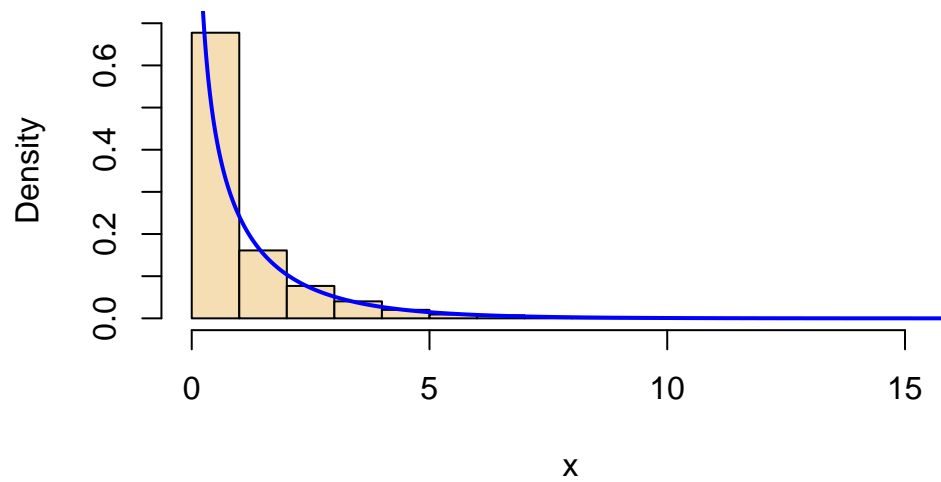
### EXAMPLE 6

If  $Z \sim NORM(0, 1)$ , then  $X = Z^2 \sim CHISQ(1)$ , that is, the chi-squared distribution with one degree of freedom. Also, if  $Z_1$  and  $Z_2$  are independently standard normal, then  $Q = Z_1^2 + Z_2^2 \sim CHISQ(2) = EXP(1/2)$ , where  $E(Q) = 2$  and  $V(Q) = 4$ . These are standard results from probability theory used in mathematical statistics. Formal proofs, not shown here, use transformation theory or moment generating functions. We illustrate these results via simulations.

```
set.seed(12)
m <- 10000
z1 <- rnorm(m)
z2 <- rnorm(m)
x <- z1^2
q <- z1^2 + z2^2

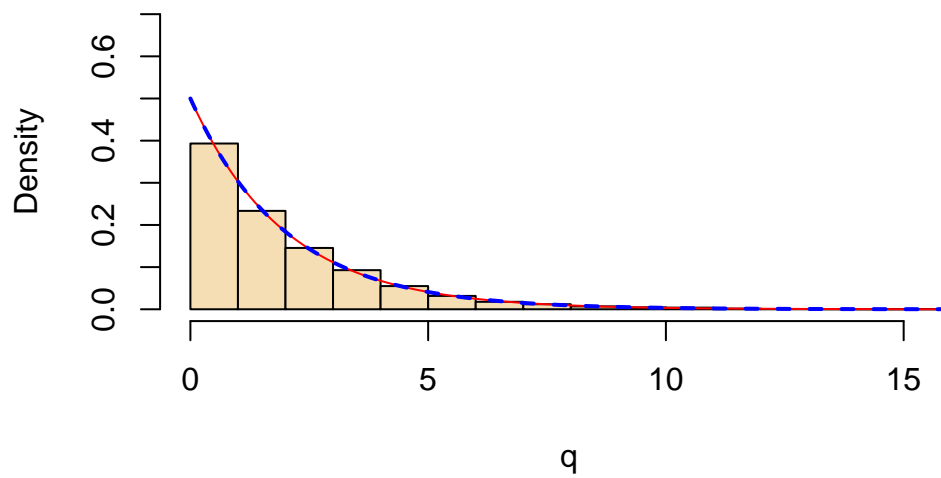
#par(mfrow=c(2,1))
mx <- max(x, q)
xx <- seq(0, mx, length=1000)
hist(x, prob=T, ylim=c(0,.7), xlim=c(0, mx), col="wheat", main="CHISQ(1)")
lines(xx, dchisq(xx, 1), col="blue", lwd=2)
```

## CHISQ(1)



```
hist(q, prob=T, ylim=c(0,.7), xlim=c(0, mx),col="wheat", main="CHISQ(2)")  
lines(xx, dexp(xx, 1/2), col="red")  
lines(xx, dchisq(xx, 2), col="blue", lwd=2,lty="dashed")
```

## CHISQ(2)



```
#ar(mfrow=c(1,1))  
  
mean(x)
```

```
[1] 1.003402
```

```
var(x)
```

```
[1] 1.939824
```

```
mean(q)
```

```
[1] 2.001475
```

```
var(q)
```

```
[1] 3.917715
```

Graphical results are shown in FIGURE E. In the lower panel, the double plotting with two line styles shows that the density functions of  $EXP(1/2)$  and  $CHISQ(2)$  are the same. ■

The following example illustrates the idea behind the most common method of generating standard normal random variables from standard uniform random variables.

## EXAMPLE 7

Suppose an archer shoots arrows at a distant target. She is aiming at the bull's eye, which we take to be the origin of a plot, but the hits are subject to random error. We model the vertical and horizontal displacements from the origin as independent standard normal random variables  $Z_1$  and  $Z_2$ . We know from EXAMPLE 6 that each arrow hits at a random distance  $D = \sqrt{Z_1^2 + Z_2^2}$  from the origin, where  $D^2 = Q \sim EXP(1/2)$ .

Now consider a line through the arrow's position to the origin, and the angle  $\Theta$  it makes with the positive  $Z_1$ -axis measured in degrees counter-clockwise. Intuitively, it seems that  $\Theta \sim UNIF(0, 360)$ , which is illustrated by the following simulation. In the code below, the arctangent takes values between  $-90$  and  $90$  degrees. Adding  $180$  degrees precisely when  $Z_1$  is negative completes the circle from  $-90$  to  $270$  degrees, and taking the resulting value modulo  $360$  (code `%%`) adjusts the values to lie in the interval  $(0, 360)$ . The resulting graph is shown in Figure F.

```
set.seed(1212)
m <- 10000
```

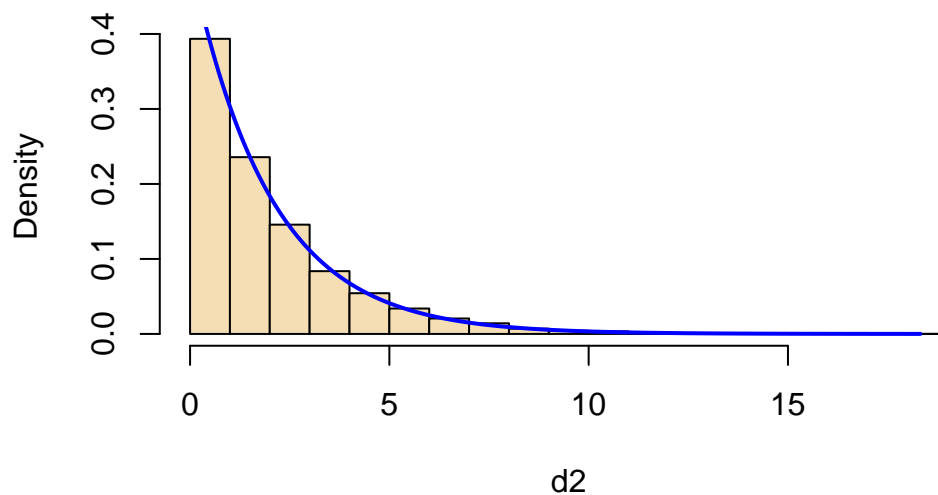
```

z1 <- rnorm(m)
z2 = rnorm(m)

#par(mfrow=c(2,1))
# squared distance from origin
d2 <- z1^2 + z2^2
hist(d2, prob=T, col="wheat")
dd <- seq(0, max(d2), length=1000)
lines(dd, dchisq(dd, 2), col="blue", lwd=2)

```

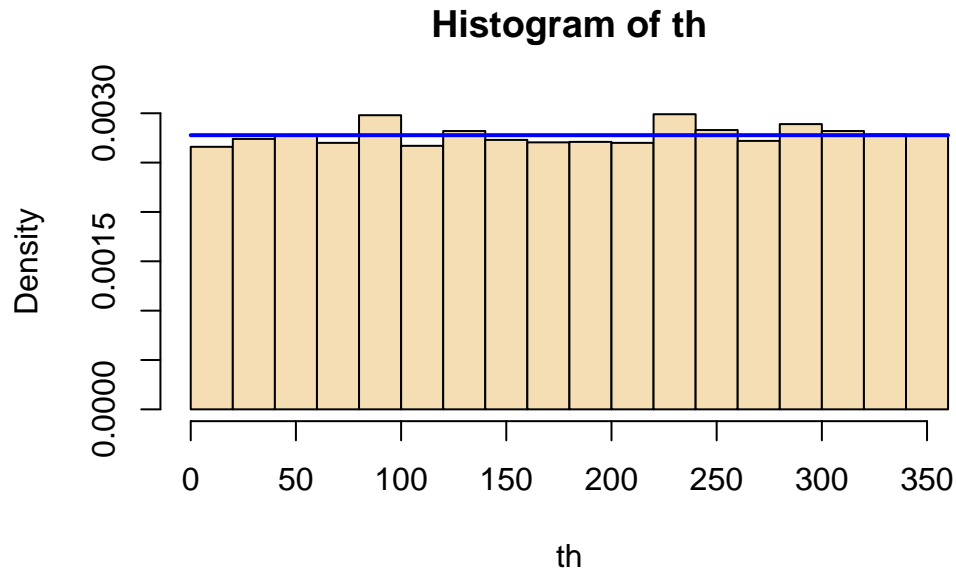
**Histogram of d2**



```

# angle in degrees (counterclockwise from right)
th <- (((180/pi)*atan(z1/z2) + 180*(z1<0)) %% 360)
hist(th, prob=T, col="wheat")
tt <- seq(0, 360, length = 1000)
lines(tt, dunif(tt, 0, 360), col="blue", lwd=2)

```



Thus the position of the hit can be modeled in polar coordinates by using two standard uniform random variables:

- The angle can be simulated as a linear transformation of a simulated observation from a standard uniform distribution (see EXAMPLE 1 and PROBLEM 1).
- The distance from the origin is the square root of an exponential random variable, and that exponential random variable can be obtained as a log transform of a standard uniform (see EXAMPLE 3 and PROBLEM 4).

Conversion from polar to rectangular coordinates reverts to the two independent standard normal random variables with which we started. This procedure of simulating two independent standard normal observations from two simulated independent standard uniform ones is known as the Box-Muller transformation. It is explored further in PROBLEM 9. ■

## PROBLEMS

- 1) *General linear transformation.* Let  $U \sim UNIF(0, 1)$  and  $X = aU + b$ , where  $a = 0$ . Use the method of EXAMPLE 1 to find the distribution of  $X$ . In particular, what is the distribution of  $Y = 1 - U$ ?
- 2) Two