

STAT 40001/STAT 50001 Statistical Computing

Lecture 7

Department of Mathematics and Statistics



PURDUE
UNIVERSITY
NORTHWEST

Probability-Simulating Experiments

```
> x=sample(c("H","T"), 10, replace=T)# Toss a coin 10 times
> x
[1] "H" "T" "T" "T" "H" "H" "T" "T" "T" "T"
> table(x)
x
H T
3 7
> table(x)/10 # calculates the probabilities
x
  H    T
0.3 0.7
```

You can check and observe that the estimated probabilities approach the true value as the number of tosses increase.

```
> x=sample(c("H","T"), 1000, replace=T)# Tossing a coin 1000
> table(x)/1000 # calculates the probabilities
x
```

Example-Coin tossing experiment

```
> space <- c(0,1)
> pi <- 0.5 # this is a fair coin
> N <- 50 # we want to flip a coin 50 times
> flips<-sample(space,size=N,replace=TRUE,prob=c(pi,1-pi))
> flips
```

cumsum function

```
> flips <-sample(c("heads","tails"),size=500,replace=TRUE)
> plot(cumsum(flips=="heads")/(1:length(flips)),type="l", ylim=c(0,1))
> abline(h = 0.5, col = "red")
```

Counting, Permutation and Combination

Factorial, permutation and combination are essentials to calculate the probability:

$$\textbf{Factorial} : n! = n \times (n-1) \times \cdots \times 3 \times 2 \times 1$$

$$\textbf{Permutation} : P(n, k) = \frac{n!}{(n-k)!} = n(n-1)(n-2) \cdots (n-(k-1))$$

$$\textbf{Combination} : C(n, k) = \frac{n!}{(n-k)!k!}$$

Factorial: `factorial(n)`

Permutation: `prod(n:n-k+1)`

Combination: `choose(n,k)`

Examples

```
> factorial(5)
```

```
[1] 120
```

```
> prod(10:(10-3+1))
```

```
[1] 720
```

```
> choose(10,3)
```

```
[1] 120
```

Birthday Problem

Suppose there are 25 randomly chosen people in a room. What is the probability two or more of them have the same birthday?

Model simplifications for a simple combinatorial solution:

- Ignore leap years
- Assume all 365 days equally likely

Solution: Let X be the number of matches.

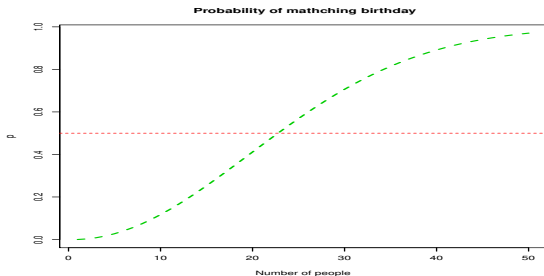
$$\begin{aligned}P(X \leq 1) &= 1 - P(X = 0) \\&= 1 - \frac{365 \times 364 \times \cdots \times (365 - 24)}{365^{25}} \\&= 1 - \prod_{i=0}^{24} \left(1 - \frac{i}{365}\right) \\&= 1 - 0.4313 = 0.5687\end{aligned}$$

In R

```
> 1-prod(1-(0:24)/365)
[1] 0.5686997
```

Birthday Problem

```
> p <- numeric(50)
> for (n in 1:50) {
+   q <- 1 - (0:(n - 1))/365
+   p[n] <- 1 - prod(q)}
> p
> plot(p, type="l", lty=2, xlab="Number of people",
main="Probability of matching birthday",lwd=3, col=3 )
```



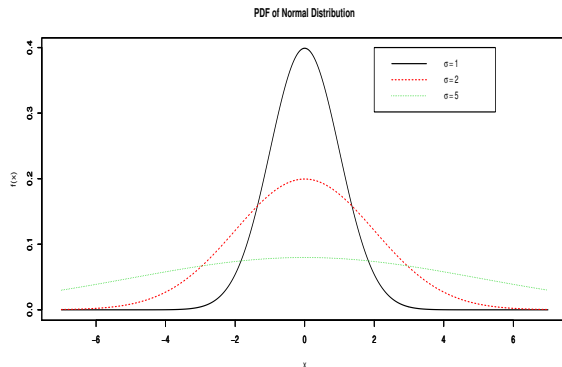
R supports a large number of distributions. Usually, four types of functions are provided for each distribution:

- d: density function
- p: cumulative distribution function, $P(X \leq x)$
- q: quantile function
- r: draw random numbers from the distribution

Normal Distribution

X has a normal distribution with parameter μ and σ if its pdf is given by

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}, \quad -\infty < x < \infty$$



Normal Distribution-R code

```
> x<-seq(-7,7,.01)
> pdf1<-dnorm(x,mean=0, sd=1)
> pdf2<-dnorm(x,mean=0, sd=2)
> pdf3<-dnorm(x,mean=0, sd=5)
> yrange<-range(pdf1,pdf2,pdf3)
> par(mfrow=c(1,1), oma=c(0,0,2,0))
> plot(x,pdf1,lty= 1,type="l",xlab="x",ylab="f(x)",xlim=c(-7,7),
      ylim=yrange,col=1, main="PDF of Normal Distribution")
> par(new=TRUE)
> plot(x,pdf2,lty= 2,type="l",xlab="",ylab="",xlim=c(-7,7),
      ylim=yrange,col=2)
> par(new=TRUE)
> plot(x,pdf3,lty= 3,type="l",xlab="",ylab="",xlim=c(-7,7),
      ylim=yrange,col=3)
> legend(2,0.4, legend=c(expression(sigma==1),expression(sigma==2),
      expression(sigma==5)),lty=1:3,col=1:3)
```

Probability Distributions

The distributions supported include continuous distributions:

- *unif*: Uniform distribution
- *norm*: Normal distribution
- *t*: Student's t- distribution
- *chisq*: Chi-square distribution
- *f*: Fisher's F distribution
- *gamma*: Gamma distribution
- *exp*: Exponential distribution
- *beta*: Beta distribution
- *lnorm*: Log-normal distribution
- *cauchy*: Cauchy distribution
- *logis*: Logistic distribution
- *weibull*: Weibull distribution

Probability Distributions

The distributions supported include discrete distributions:

- *binom*: Binomial distribution
- *geom* : Geometric distribution
- *hyper*: Hypergeometric distribution
- *nbinom*: Negative Binomial distribution
- *pois*: Poisson distribution

Using help: ?runif

Uniform {stats} R Documentation

The Uniform Distribution

```
dunif(x, min=0, max=1, log = FALSE)
```

```
punif(q, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
```

```
qunif(p, min=0, max=1, lower.tail = TRUE, log.p = FALSE)
```

```
runif(n, min=0, max=1)
```

Arguments

`x,q` vector of quantiles.

`p` vector of probabilities.

`n` number of observations. If `length(n) > 1`, the length is taken to be the number required.

`min,max` lower and upper limits of the distribution.

`log, log.p` logical;if TRUE, probabilities `p` are given as

`lower.tail` logical;if TRUE (default),probabilities are $P[X \leq x]$ otherwise, $P[X > x]$.

Computer-generated random numbers are not really random; they are “pseudo-random.” This means that the computer generates a fixed sequence of random-looking numbers. The sequence is very, very long, so the pseudo-random numbers are not repeated within any typical computation. Therefore when we generate a set of random numbers from any distribution every time we will get different set. But if we want to generate the same set of “random” numbers at later we need to use the **set.seed()** function. The **set.seed()** function puts the random number generator in a reproducible state. To avoid using the same random numbers we have to set the seed to different numbers each time.

```
> set.seed(20)
> rnorm(5)
[1] 1.1626853 -0.5859245 1.7854650 -1.3325937 -0.4465668
> set.seed(20)
> rnorm(5)
[1] 1.1626853 -0.5859245 1.7854650 -1.3325937 -0.4465668
```

Continuous Distributions

The cumulative probability function is a straightforward notion: it is an S-shaped curve showing, for any value of x , the probability of obtaining a sample value that is less than or equal to x . Here is what it looks like for the normal distribution:

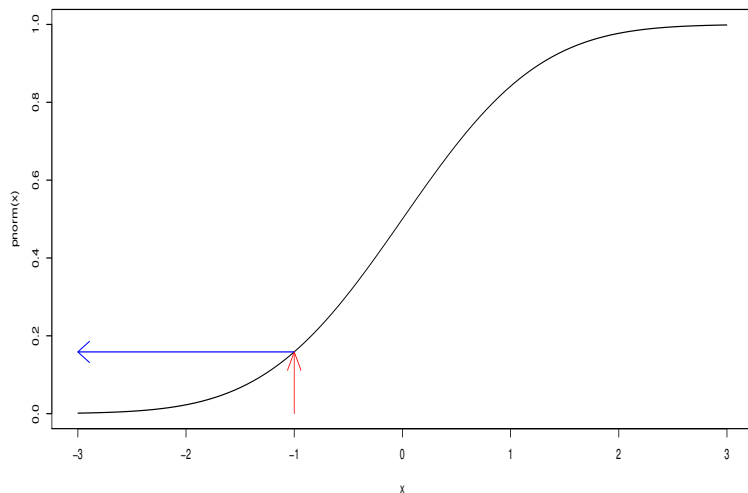
```
>curve(pnorm(x),-3,3, main="CDF of Standard Normal Distribut  
>arrows(-1,0,-1,pnorm(-1),col="red")  
>arrows(-1,pnorm(-1),-3,pnorm(-1),col="blue")
```

The value of $x(-1)$ leads up to the cumulative probability (red arrow) and the probability associated with obtaining a value of this size (-1) or smaller is on the y axis (blue arrow). The value on the y axis is 0.1586553:

```
>pnorm(-1)  
[1] 0.1586553
```

CDF- Normal Distribution

CDF of Standard Normal Distribution



Continuous Distributions

For a continuous random variable X , the function $f(x)$ is said to be a probability density function if it satisfies the following properties:

a) $f(x) \geq 0$ for all values of x .

b) $P(a < x < b) = \int_a^b f(x)dx$ for all a and b .

c) $\int_{-\infty}^{\infty} f(x)dx = 1$.

For a normal distribution the slope starts out very shallow up to about $x=-2$, increases up to a peak (at $x = 0$ in this example) then gets shallower, and becomes very small indeed above about $x=2$.

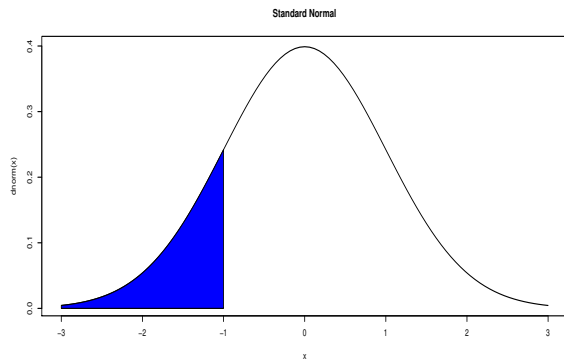
Here is what the density function of the normal (dnorm) looks like:

```
>curve(dnorm(x),-3,3,main="pdf of standard Normal Distribution")
```

```
>dnorm(-1)
```

```
[1] 0.2419707
```

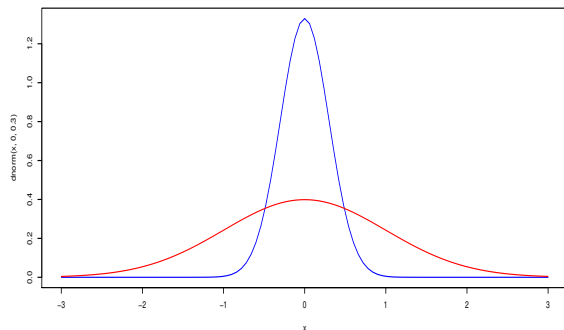

PDF- Normal Distribution



```
cord.x <- c(-3,seq(-3,-1,0.01),-1)
cord.y <- c(0,dnorm(seq(-3,-1,0.01))),0)
curve(dnorm(x),xlim=c(-3,3),main='Standard Normal')
polygon(cord.x,cord.y,col='blue')
```

PDF- Two Normal Distribution

Superimpose two PDFs:



```
curve(dnorm(x,0,0.3),from=-3,to=3,col="blue")  
curve(dnorm(x,0,1),from=-3, to=3, col="red", add=T)
```

The "from" and "to" can be omitted.

Sampling from Uniform Distribution

Use R to simulate 1000 times the sampling distribution of the mean, \bar{X} of 1, 50, and 100 observations from the uniform distribution. Create a histogram and determine the mean and standard deviation of these simulations.

```
> xbar1=numeric(1000)
> for (i in 1:1000){x=runif(1);xbar1[i]=mean(x)}
> hist(xbar1,col="green",main="CLT for uniform n=1")

> windows()
> xbar2=numeric(1000)
> for (i in 1:1000){x=runif(50);xbar2[i]=mean(x)}
> hist(xbar2,col="blue",main="CLT for uniform n=50")

> windows()
> xbar3=numeric(1000)
> for (i in 1:1000){x=runif(100);xbar3[i]=mean(x)}
> hist(xbar3,col="orange",main="CLT for uniform n=100")
```

Few Examples- Marking and shading

1) Plot the standard normal PDF and mark the 90th percentile:

```
curve(dnorm,-3,3)
lines(qnorm(0.9),dnorm(qnorm(0.9)),type="h", col="red")
```

2) Shade the area under the $N(0,1)$ pdf to the right of the 90th percentile:

```
x1=seq(qnorm(0.9),3,0.01);
y1=dnorm(x1)
curve(dnorm,-3,3); lines(x1,y1,type="h",col="red")
```

3) More options:

```
curve(dnorm,-3,3)
polygon(c(rep(1,201),rev(seq(1,3,.01))),c(dnorm(seq(1,3,.01)),
dnorm(rev(seq(1,3,.01))))),col="orange", lty=2, lwd=2,
border="red")
```

Example

Generate 5 random numbers from $N(2, 2^2)$

```
> rnorm(5, mean=2, sd=2)
[1] 3.1985264 0.4473159 -0.8669186 -1.1561986 3.0654766
```

```
> rnorm(5) # This generates with mean 0 and sd=1
[1] 0.6998394 1.0070251 -1.0634292 -0.1712023 0.9539368
```

Obtain 95% quantile for the standard normal distribution

```
> qnorm(0.95)
[1] 1.644854
```

Binomial Distribution

A Binomial experiment consists of n independent Bernoulli trials - with the probability of success p remaining constant throughout the trials. Let X be the number of successes recorded. Then X is said to have the binomial distribution with parameters n and p .

Binomial Properties

- a) The experiment consists of a fixed number, n , of Bernoulli trials.
- b) The trials are identical and independent with probability of success, p .
- c) The random variable X denotes the the number of successes obtained in the n trials.

The probability mass function of the binomial random variable X is:

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}, \quad x = 0, 1, \dots, n;$$

When $n = 1$ the binomial distribution reduces to the bernoulli distribution.

Binomial Distribution

- `dbinom(x, size, prob)`
- `pbinom(q, size, prob)`
- `qbinom(p, size, prob)`
- `rbinom(n, size, prob)`

In Binomial distribution where size is the sample size and prob is the probability of success. Example:

What is the probability of 0 to 5 heads when a fair coin is tossed 10 times

`dbinom(0:5, 10, .5)`

Probability of 5 or less heads of fair coin out of 10 flips

`pbinom(5, 10, .5)`

Suppose you have a biased coin that has a probability of 0.8 of coming up heads. The probability of getting 5 heads in 16 tosses of this coin is

`dbinom(5,16, .8)`

The 0.25 quantile is

`qbinom(.25,16,.8)`

Poisson Distribution

Given a continuous interval (in time, length, etc), assume discrete events occur randomly throughout the interval. If the interval can be partitioned into subintervals of small enough length such that (i) the probability of more than one occurrence in a subinterval is zero;

(ii) the probability of one occurrence in a subinterval is the same for all subintervals and proportional to the length of the subinterval;

(iii) the occurrence of an event in one subinterval has no effect on the occurrence or non-occurrence in another non-overlapping subinterval,

If the mean number of occurrences in the interval is λ , the random variable X that equals the number of occurrences in the interval is said to have the Poisson distribution with parameter λ .

The probability mass function of the Poisson random variable X is:

$$P(X = x) = \frac{\lambda^x \exp(-\lambda)}{x!}, \quad x = 0, 1, 2, \dots;$$

Poisson Distributions

- `dpois(x, lamda)`
- `ppois(q, lamda)`
- `qpois(p, lamda)`
- `rpois(n, lamda)`

In Poisson distribution with mean= λ probability of 0,1, or 2 events with $\lambda = 4$

`dpois(0:2, 4)`

probability of at least 3 events with $\lambda = 4$

1- `ppois(2,4)`

Suppose a certain region of California experiences about 5 earthquakes a year. Assume occurrences follow a Poisson distribution. What is the probability of 3 earthquakes in a given year?

Here $\lambda = 5$.

`dpois(3,5)`

Student's -t distribution

If X_1, X_2, \dots, X_n is a random sample from a normal distribution with mean μ and variance σ^2 then

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

This is an important result, but the major difficulty arise on application in which cases σ is unknown. In this case we replace σ with its estimate s and we study the distribution of $\frac{\bar{X} - \mu}{s/\sqrt{n}}$. The distribution of this expression will have the student's t-distribution. A random variable X is said to have t-distribution with n degrees of freedoms if its pdf is given by

$$f(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \left\{ 1 + \frac{x^2}{n} \right\}^{-\frac{n+1}{2}} \text{ for } -\infty < x < \infty.$$

R codes for t- distribution

- dt: density function of t-distribution
- pt: cumulative distribution function
- qt: quantile function of t- distribution
- rt: draw random numbers from the t-distribution

Need to choose the parameter n .

$t(df, ncp)$?noncentral t distribution with noncentrality parameter ncp

```
> pt(-1,df=10)
```

```
[1] 0.1704466
```

```
> pt(0, df=10)
```

```
[1] 0.5
```

```
> pt(1, df=10)
```

```
[1] 0.8295534
```

```
# Calculating percentiles
```

```
> # Find the 25th percentile with a degree of freedom=4
```

```
> qt(.25, df=4)
```

```
[1] -0.7406971
```

Quantile-Quantile Plots for Normal Distributions

One of the most useful graphical procedure for assessing distributions is the quantile-quantile plot. A quantile-quantile (Q-Q) plot plots the quantiles of one distribution against the quantiles of another distribution as (x,y) points. When two distributions have similar shapes, the points will fall along a straight line. The R function to draw a quantile-quantile plot is `qqplot(x,y)`. Histograms can be used to compare two distributions. However, it is rather challenging to put both histograms on the same graph.

R offers two statements: `qqnorm()`, to test the goodness of fit of a Gaussian distribution, or `qqplot()` for any kind of distribution.

Example

```
x.norm<-rnorm(n=200,m=10,sd=2)
hist(x.norm,main="Histogram of observed data")
plot(density(x.norm),main="Density estimate of data")
plot(ecdf(x.norm),main="Empirical CDF")
z.norm<-(x.norm-mean(x.norm))/sd(x.norm) ## standardized data
qqnorm(z.norm) ## drawing the QQplot
abline(0,1) ## drawing a 45-degree reference line
```

Example

```
unif50=runif(50)
unif100=runif(100)
norm50=rnorm(50)
norm100=rnorm(100)
lognorm50=exp(rnorm(50))
lognorm100=exp(rnorm(100))
par(mfrow=c(2,3))
plot(ecdf(unif50),pch="16")
plot(ecdf(unif100),pch="17")
plot(ecdf(norm50),pch="18")
plot(ecdf(norm100),pch="19")
plot(ecdf(lognorm50),pch="20")
plot(ecdf(lognorm100),pch="21")
```

Example

```
unif50=runif(50)
unif100=runif(100)
norm50=rnorm(50)
norm100=rnorm(100)
lognorm50=exp(rnorm(50))
lognorm100=exp(rnorm(100))
par(mfrow=c(2,3))
qqnorm(unif50,main="Normal QQ-plot\n unif50")
qqnorm(unif100,main="Normal QQ-plot\n unif100")
qqnorm(norm50,main="Normal QQ-plot\n norm50")
qqnorm(norm100,main="Normal QQ-plot\n norm100")
qqnorm(lognorm50,main="Normal QQ-plot\n lognorm50")
qqnorm(lognorm100,main="Normal QQ-plot\n lognorm100")
```