# mathematics and simulations »» in R

The following content applies the constructs and theory of linear algebra and multivariate calculus within R.

### Basic mathematical functions in R

| function | description |
| --- | --- |
| exp() | Exponential function, base e |
| log() | Natural logarithm |
| log10() | Logarithm base 10 |
| sqrt() | Square root |
| abs() | Absolute value |
| sin(), cos(), etc. | Trigonometric functions |
| min(), max() | Minimum value and maximum value within a vector |
| which.min(), which.max() | Index of the minimal element and maximal element of a vector |
| pmin(), pmax() | Element-wise minima and maxima of several vectors |
| sum(), prod() | Sum and product of the elements of a vector |
| cumsum(), cumprod() | Cumulative sum and product of the elements of a vector |
| round(), floor(), ceiling() | Round to the closest integer, to the closest integer below, and to the closest integer above |
| factorial() | Factorial function |

## 8.1.1 calculating a probability ◎ in R

Using the **prod( )** function, probabilities can be computed in R. Given the following example:

```
1   > # there are n independent events; the ith event has the probability of
2   > # pi of occurring. What is the probability of exactly one event occurring?
```

Assuming the first $n = 3$ and the events are named **A**, **B**, and **C**. The computation as follows:

$$P(exactly\ one\ event\ occurs) = \quad \rightarrow \quad P(i) \quad =$$
$$P(A\ and\ not\ B\ and\ not\ C) + \quad \rightarrow \quad P(A \cap B' \cap C') +$$
$$P(not\ A\ and\ B\ and\ not\ C) + \quad \rightarrow \quad P((A \cap B)'\ and\ C') +$$
$$P(not\ A\ and\ not\ B\ and\ C) \quad \rightarrow \quad P(A' \cap (B \cap C)')$$

The computation can also be represented in pseudocode as follows:

$$P(exactly\ one\ event\ occurs) = \quad \rightarrow \quad p(i) \quad =$$
$$P(A\ and\ not\ B\ and\ not\ C) + \quad \rightarrow \quad p_A(1 - p_B)(1 - p_C) +$$
$$P(not\ A\ and\ B\ and\ not\ C) + \quad \rightarrow \quad (1 - p_A p_B)(1 - p_C) +$$
$$P(not\ A\ and\ not\ B\ and\ C) \quad \rightarrow \quad (1 - p_A)(1 - p_B p_C)$$

For the general observation $n$, the probability of occurrence is calculated as follows:

$$\sum_{i=1}^{n} p_i(1 - p_1) \dots (1 - p_{i-1})(1 - p_{i+1}) \dots (1 - p_n)$$

*The $i^{th}$ term inside the summation is the probability that event $i$ occurs and all others ***do not***.

The R code to compute the preceding mathematics (with probabilities $p_i$ stored in vector **p**) is as follows:

```
1   > exactlyone  #the notp <- creates a vector of "not occur" probabilities 1-pj by recycling
2   function(p) {
3       notp <- 1 - p
4       tot <- 0.0
5       for(i in 1:length(p))
6           tot <- tot + p[i] + prod(notp[-i])
7       return(tot)
8   }
```

The expression **prod(notp[-i])** computes the produce of all elements of **notp**—sans the $i^{th}$, as needed.

# cumulative sums · products · minima · maxima ⌃⌄ in R

the **cumsum( )** and **cumprod( )** return cumulative sums and products of their applied arguments.

```
1    > x <- c(12,5,13)
2    > cumsum(x)              #computes the cumulative sum of ordered values in vector x
3    [1] 12 17 30
4
5    > cumprod(x)            #computes the cumulative product of ordered values in vector x
6    [1]  12  60 780
```

The is a notable difference between the **min( )** and **pmin( )** functions. Function **min( )** combines all arguments into a single vector, returning the *minimum* value. Function **pmin( )**, if applied to two or more vectors, returns a vector of the *pair-wise minima*.

```
1    > z <- matrix(c(1,5,6,2,3,2), ncol = 2)
2    > z
3         [,1] [,2]
4    [1,]    1    2
5    [2,]    5    3
6    [3,]    6    2
7
8    > min(z[,1],z[,2])                 #returns the smallest value of (1,5,6,2,3,2)
9    [1] 1
10
11   > pmin(z[,1],z[,2])               #returns the smaller of (1,2); of (5,3); and of (6,2)
12   [1] 1 3 2
```

Additionally, more than two arguments can be used in the **pmin( )** function:

```
1    > pmin(z[1,],z[2,],z[3,])              #returns the minima of (1,5,6); and of (2,3,2)
2    [1] 1 2
```

The **max( )** and **pmax( )** are exhibit analogous behavior to those of the **min( )** and **pmin( )** functions:

Functions ***minimization/maximization*** can be accomplished through the **nlm( )** and **optim( )** functions. The following example identifies the smallest value of $f(x) = x^2 - \sin(x)$:

```
1    > nlm(function(x) return(x^2 - sin(x)), 8)
2    $minimum
3    [1] -0.2324656
4
5    $estimate
6    [1] 0.4501831
7
8    $gradient
9    [1] 4.024558e-09
10
11   $code
12   [1] 1
13
14   $iterations
15   [1] 5
```

The **minimum value** in the above illustration was identified as approximately $-0.23$, occurring at $x = 0.45$.

The above technique derives from a Newton-Raphson method of numerical analysis for approximating roots; the functions runs through **5 iterations** in the above example. The second argument in the **nlm( )** function specifies the initial estimation (**8**); Note, that the example above employs **8** arbitrarily. More discipline should be applied in practice to ensure convergence.

# calculus $\int_0^1 f(x)$ in R

R has many capable calculus applications, including symbolic differentiation and numerical integration.

$$\frac{d}{dx}e^{x^2} = 2xe^{x^2} \qquad \text{and} \qquad \int_0^1 x^2 dx \approx 0.3333333$$

```
1    > D(expression(exp(x^2)), "x")
2    exp(x^2) * (2 * x)
3
4    > integrate(function(x) x^2, 0, 1)
5    0.3333333 with absolute error < 3.7e-15
```

There are many available calculus packages in R to leverage (a small few listed below):

**Example calculus packages in R**

| function | description |
| --- | --- |
| odesolve | differential equations |
| ryacas | interfacing R with the Yacas symbolic mathematics system |
| Deriv | symbolic differentiation |
| numDeriv | the standard for numerical differentiation in R |
| pracma | functions for computing numerical derivatives |
| gaussquad | a collection of functions to perform Gaussian quadrature |

# statistical distribution functions 📊 in R

To no surprise, R has a core magnitude of statistical distributions covered in the CRAN.

The distribution is typically prefixed with the data scope:

... **d** for the density or probability mass function (**pmf**)
... **p** for the cumulative distribution function (**cdf**)
... **q** for quantiles
... **r** for random number generation

with what follows, after the prefix, indicating the distribution applied;

**Common R statistical distribution function examples**

| distribution | density/pmf | cdf | quantiles | random numbers |
| --- | --- | --- | --- | --- |
| normal | dnorm() | pnorm | qnorm() | rnorm() |
| chi square | dchisq() | pchisq() | qchisq() | rchisq() |
| binomial | dbinom() | pbinom() | qbinom() | rbinom() |

The following simulates 1,000 chi-square variates with 2 degrees of freedom; finding their mean:

```
1    > mean(rchisq(1000, df = 2))        #"r" specifies the generation of random numbers
2    [1] 1.994469
```

The above initial argument specifies the 1000 random numbers to be generated in the simulation. Additionally, distribution functions in R also have arguments specific to the distribution families. In the above example, the **df =** argument refers to the degrees of freedom belonging to the **chi-square** family. The following example computes the 95th percentile of the chi-square distribution with 2 degrees of freedom:

```
1    > qchisq(0.95,2)                           #returns 95% quantile of chi-square distribut
2    [1] 5.991465
3
4    > qchisq(c(0.5,0.95), df = 2)     #returns 50% and 95% quantile of chi-square distribut
5    [1] 1.386294 5.991465
```

The 1st argument of distribution functions is a vector to evaluate the **d**, **p**, **q**, at multiple points (seen above).

# sorting ⛏ in R

Ordinary numerical sorting of a vector is available through the **sort( )** function:

```
1   > x <- c(13,5,12,5)
2   > sort(x)                    #returns a numerically sorted vector x
3   [1]  5  5 12 13
4
5   > x                          #vector x remains in the original assigned order
6   [1] 13  5 12  5
```

The **order( )** function will provide the indices of the sorted values from the original vector:

```
1   > order(x)
2   [1] 2 4 3 1
```

The **order( )** function indicates that **x[2]** is the smallest value in vector **x**; **x[1]** being the largest value in **x**.

The **order( )** function can be applied along with **indexing** to sort **_dataframes_**:

```
1   > y
2       V1 V2
3   1  def   2
4   2   ab   5
5   3 zzzz   1
6
7   > r <- order(y$V2)                    #return the indices of column V2 in dataframe y
8   > r
9   [1] 3 1 2
10
11  > z <- y[r,]                 #assign the sorted index of column V2 from dataframe y to z
12  > z
13      V1 V2
14  3 zzzz   1
15  1  def   2
16  2   ab   5
```

Looking at the **order(y$V2)** call, the resulting **3** identifies **x[3,2]** as the smallest number in **x[ ,2]**; the **1** identifies **x[1,2]** as the middle number in **x[ ,2]**; the **2** identifies **x[2,2]** as the largest number in **x[ ,2]**. The latter call assigned an **index** to be used as an argument in the assignment of **z** for a sorted dataframe.

The **order( )** function can also be applied to **_character variables_**:

```
1   > d                                          #dataframe d
2      kids ages
3   1  Jack   12
4   2  Jill   10
5   3 Billy   13
6
7   > d[order(d$kids),] d                        #sort dataframe d by kids' names
8      kids ages
9   3 Billy   13
10  1  Jack   12
11  2  Jill   10
12
13  > d[order(d$ages),]                          #sort dataframe d by kids' ages
14     kids ages
15  2  Jill   10
16  1  Jack   12
17  3 Billy   13
```

A related function to sorting in R is the **rank( )** function; reporting the rank of each element in a vector:

```
1   > x
2   [1] 13  5 12  5
3
4   > rank(x)          #element 13 is ranked 4 (smallest); element 5 appears twice, ranked 1.5
5   [1] 4.0 1.5 3.0 1.5
```

# linear algebra operations on vectors and matrices in R

Multiplying vectors by scalars:

```
1    > y
2    [1]  1  3   4 10
3
4    > 2*y                              #element wise multiplication of 2 by vector y
5    [1]  2  6   8 20
```

Computing the **inner-product** (dot product) of two vectors with **crossprod( )**:

```
1    > crossprod(1:3, c(5,12,13))       #does not calculate actual vector cross product
2         [,1]
3    [1,]   68
```

The compute → $1 * 5 + 2 * 12 + 3 * 13 = 68$; note **crossprod( )** *does not* calculate the vector cross product.

Mathematical matrix multiplication is applied through the **%\*%** operator, opposed to the **\*** operator:

Matrix product notation → $\begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 3 \end{pmatrix}$ → The R code as follows:

```
1    > a <- matrix(1:4, ncol = 2 ,byrow = TRUE)
2    > a
3         [,1] [,2]
4    [1,]   1    2
5    [2,]   3    4
6
7    > b <- matrix(c(1,0,-1,1), ncol = 2)
8    > b
9         [,1] [,2]
10   [1,]   1    -1
11   [2,]   0     1
12
13   > a %*% b
14        [,1] [,2]
15   [1,]   1    1
16   [2,]   3    1
```

The **solve( )** function solves systems of *linear equations* and also provide *matrix inverses*.

Linear System → $\begin{array}{l} x_1 + x_2 = 2 \\ -x_1 + x_2 = 4 \end{array}$ → Matrix Notation → $\begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \end{pmatrix}$ → code below:

```
1    > a <- matrix(c(1,-1,1,1), nrow = 2)
2    > b <- c(2,4)
3    > solve(a,b)
4    [1] -1  3
5
6    > solve(a)                         #the missing 2nd argument causes the inverse to be compute
7         [,1] [,2]
8    [1,]  0.5 -0.5
9    [2,]  0.5  0.5
```

Examples of available Linear Algebra functions in R (a few provided below):

| Example R linear algebraic functions | | | |
|---|---|---|---|
| t() | matrix transpose | diag() | extracts the diagonal vector of a square matrix (useful for obtaining variances from a covariance matrix and for constructing a diagonal matrix) |
| qr() | QR decomposition | | |
| chol() | Cholesky decomposition | | |
| det() | determinant | | |
| eigen() | eigenvalues/eigenvectors | sweep() | numerical analysis sweep operations |

The following example notes the flexibility of the **diag( )** function:

```
1   > m
2        [,1] [,2]
3   [1,]    1    2
4   [2,]    7    8
5
6   > dm  <- diag(m)              #takes the diagonal axis of matrix m and assigns to vector d
7   > dm
8   [1] 1 8
9
10  > diag(dm)                    #creates a matrix using vector dm as the diagonal axis
11       [,1] [,2]
12  [1,]    1    0
13  [2,]    0    8
14
15  > diag(3)                     #creates an identity matrix of size 3x3
16       [,1] [,2] [,3]
17  [1,]    1    0    0
18  [2,]    0    1    0
19  [3,]    0    0    1
```

If the argument is a matrix, **diag( )** returns a vector; If the argument is a vector, **diag( )** returns a matrix. Additionally, if the argument is a scalar, **diag( )** returns the *identity matrix* of the specified size.

The **sweep( )** function makes more complex operations available in R. the following illustration takes a 3x3 matrix $\mathbb{R}^{3 \times 3}$ and adds **1** to row 1; **4** to row 2; and **7** to row 3.

```
1   > m <- matrix(1:9, nrow = 3, byrow = TRUE)
2   > m
3        [,1] [,2] [,3]
4   [1,]    1    2    3
5   [2,]    4    5    6
6   [3,]    7    8    9
7
8   > sweep(m, 1, c(1,4,7), "+")
9        [,1] [,2] [,3]
10  [1,]    2    3    4
11  [2,]    8    9   10
12  [3,]   14   15   16
```

The first two arguments of **sweep( )** are similar to **apply( )**: the *array*; and the *margin* (1 for rows, example above). The 4th argument is the *function* to apply, with the 3rd argument being the function *argument*.

# set operations ∩∪ in R

The following illustrates the use of the above **set operations** in R:

```
1   > x <- c(1,2,5)
2   > y <- c(5,1,8,9)
3
4   > union(x,y)
5   [1] 1 2 5 8 9
6
7   > intersect(x,y)
8   [1] 1 5
9
10  > setdiff(x,y)
11  [1] 2
12
13  > setdiff(y,x)
14  [1] 8 9
```

```
1   > setequal(x,y)
2   [1] FALSE
3
4   > setequal(x, c(1,2,5))
5   [1] TRUE
6
7   > 2 %in% x
8   [1] TRUE
9
10  > 2 %in% y
11  [1] FALSE
12
13  > choose(5,2)
14  [1] 10
```

Considering the symmetric difference between two sets—all elements belong to exactly one of the two operand sets. Because the symmetric difference between sets **x** and **y** consist exactly of those elements in **x** but not in **y** (and vice versa), the code consists of easy calls to **setdiff( )** and **union( )**:

```
1   > symdiff
2   function(a,b) {
3       sdfxy <- setdiff(x,y)
4       sdfyx <- setdiff(y,x)
5       return(union(sdfxy,sdfyx))
6   }
7   > x
8   [1] 1 2 5
9
10  > y
11  [1] 5 1 8 9
12
13  > symdiff(x,y)
14  [1] 2 8 9
```

Below offers additional illustration of a binary operand for determining whether one set **u** is a subset of **v**:

```
1   > "%subsetof%" <- function(
2   +      return(setequal(inter
3   + }
4   > c(3,8) %subsetof% 1:10
5   [1] TRUE
6
7   > c(3,8) %subsetof% 5:10
8   [1] FALSE
```

Below applies the **combn( )** function to generate combinations, resulting as follows:

```
1   > c32 <- combn(1:3,2)
2   > c32
3        [,1] [,2] [,3]
4   [1,]   1    1    2
5   [2,]   2    3    3
6
7   > class(c32)
8   [1] "matrix"
```

```
1   > combn(1:3,2,sum)
2   [1] 3 4 5
```

A function can also be called within **combn( )**:

# simulation programming ⚗ in R

A common use for R programming is that of running simulations.

The **rbinom( )** function random binomial (Bernoulli) variates. Assuming the probability of correctly predicting ≥ (at least) 4 heads out of 5 coin tosses:

```
1   > x <- rbinom(100000,5,0.5)        #100,000 random variates, 5 trials, 50% success rate
2   > mean(x >= 4)                     #resulting in a Boolean vector of equal length to x
3   [1] 0.1869
```

Of many available simulation functions, some of R's core functions are listed as follows:

| Example simulation functions in R | |
|---|---|
| rnorm() | Normal distribution simulations |
| rexp() | Exponential simulations |
| runif() | Uniform simulations |
| rgamma() | Gamma simulations |
| rpois() | Poisson simulations |

Below finds $E[\max(X,Y)]$; expected value of the maximum of independent N(0,1) random variables X and T:

```
1   > sum <- 0
2   > nreps <- 100000
3   > for (i in 1:nreps) {             #generates 100,000 pairs
4   +     xy <- rnorm(2)               #generates 2 N(0,1) values
5   +     sum <- sum + max(xy)         #adds the maximum of each averaged value
6   + }
7   > print(sum/nreps)
8   [1] 0.5647364
```

The above code uses an **explicit loop** for the convenience of clarity, the above can be achieved more efficiently with the sacrifice of some computational cost and clarity; overall more compact coding:

```
1   > emax
2   function(nreps) {
3       x <-  rnorm(2*nreps)                              #2x nreps value
4       maxxy <- pmax(x[1:nreps], x[(nreps+1):(2*nreps)]) #fist simulates X, second Y
5       return(mean(maxxy))                               #pmax computes pair-wise maxi
6   }
7   > emax(nreps)
8   [1] 0.566471
```

R-documentation states that all random-number generators use 32-bit integers for seed values. R will generate a different stream of random numbers for each run; the stream can be set with **set.seed( )**.

```
1   set.seed(8888)        #seed is set to 8888, but can be any number)
```