



Lecture 1: Introduction to the Course Basics of Data

S. Morteza Najibi

Shiraz University

February 25, 2019

- ▶ Arrays
- ▶ Matrices
- ▶ Lists
- ▶ Dataframes
- ▶ Structures of structures

Many data structures in R are made by adding bells and whistles to vectors, so “vector structures”

Most useful: **arrays**

```
x <- c(7, 8, 10, 45)
x.arr <- array(x,dim=c(2,2))
x.arr
```

```
##      [,1] [,2]
## [1,]    7   10
## [2,]    8   45
```

dim says how many rows and columns; filled by columns

Can have 3, 4, ... n dimensional arrays; **dim** is a length- n vector

Some properties of the array:

```
dim(x.arr)
```

```
## [1] 2 2
```

```
is.vector(x.arr)
```

```
## [1] FALSE
```

```
is.array(x.arr)
```

```
## [1] TRUE
```

```
typeof(x.arr)
```

```
## [1] "double"
```

```
str(x.arr)
```

```
## num [1:2, 1:2] 7 8 10 45
```

```
attributes(x.arr)
```

```
## $dim
```

```
## [1] 2 2
```

`typeof()` returns the type of the *elements*

`str()` gives the **structure**: here, a numeric array, with two dimensions, both indexed 1–2, and then the actual numbers

Exercise: try all these with `x`

Can access a 2-D array either by pairs of indices or by the underlying vector:

```
x.arr[1,2]
```

```
## [1] 10
```

```
x.arr[3]
```

```
## [1] 10
```

Omitting an index means “all of it”:

```
x.arr[c(1:2),2]
```

```
## [1] 10 45
```

```
x.arr[,2]
```

```
## [1] 10 45
```

Using a vector-style function on a vector structure will go down to the underlying vector, *unless* the function is set up to handle arrays specially:

```
which(x.arr > 9)
```

```
## [1] 3 4
```


Many functions *do* preserve array structure:

```
y <- -x  
y.arr <- array(y,dim=c(2,2))  
y.arr + x.arr
```

```
##      [,1] [,2]  
## [1,]    0    0  
## [2,]    0    0
```

Others specifically act on each row or column of the array separately:

```
rowSums(x.arr)
```

```
## [1] 17 53
```

We will see a lot more of this idea

Running example: resource allocation (“mathematical programming”)

Factory makes cars and trucks, using labor and steel

- ▶ a car takes 40 hours of labor and 1 ton of steel
- ▶ a truck takes 60 hours and 3 tons of steel
- ▶ resources: 1600 hours of labor and 70 tons of steel each week

In R, a matrix is a specialization of a 2D array

```
factory <- matrix(c(40,1,60,3),nrow=2)
is.array(factory)
```

```
## [1] TRUE
```

```
is.matrix(factory)
```

```
## [1] TRUE
```

could also specify `ncol`, and/or `byrow=TRUE` to fill by rows.

Element-wise operations with the usual arithmetic and comparison operators (e.g., `factory/3`)

Compare whole matrices with `identical()` or `all.equal()`

Gets a special operator

```
six.sevens <- matrix(rep(7,6),ncol=3)
six.sevens
```

```
##      [,1] [,2] [,3]
## [1,]    7    7    7
## [2,]    7    7    7
```

```
factory %*% six.sevens # [2x2] * [2x3]
```

```
##      [,1] [,2] [,3]
## [1,]  700  700  700
## [2,]   28   28   28
```

What happens if you try `six.sevens %*% factory`?

Numeric vectors can act like proper vectors:

```
output <- c(10,20)
factory %*% output
```

```
##      [,1]
## [1,] 1600
## [2,]   70
```

```
output %*% factory
```

```
##      [,1] [,2]
## [1,]  420  660
```

R silently casts the vector as either a row or a column matrix

Transpose:

```
t(factory)
```

```
##      [,1] [,2]  
## [1,]   40   1  
## [2,]   60   3
```

Determinant:

```
det(factory)
```

```
## [1] 60
```

The `diag()` function can extract the diagonal entries of a matrix:

```
diag(factory)
```

```
## [1] 40  3
```

It can also *change* the diagonal:

```
diag(factory) <- c(35,4)  
factory
```

```
##      [,1] [,2]  
## [1,]   35  60  
## [2,]    1   4
```

Re-set it for later:

```
diag(factory) <- c(40,3)
```

```
diag(c(3,4))
```

```
##      [,1] [,2]  
## [1,]    3    0  
## [2,]    0    4
```

```
diag(2)
```

```
##      [,1] [,2]  
## [1,]    1    0  
## [2,]    0    1
```



```
solve(factory)
```

```
##           [,1]      [,2]  
## [1,]  0.05000000 -1.0000000  
## [2,] -0.01666667  0.6666667
```

```
factory %*% solve(factory)
```

```
##      [,1] [,2]  
## [1,]    1    0  
## [2,]    0    1
```

Solving the linear system $\mathbf{A}\vec{x} = \vec{b}$ for \vec{x} :

```
available <- c(1600,70)
solve(factory,available)
```

```
## [1] 10 20
```

```
factory %*% solve(factory,available)
```

```
##      [,1]
## [1,] 1600
## [2,]   70
```

We can name either rows or columns or both, with `rownames()` and `colnames()`

These are just character vectors, and we use the same function to get and to set their values

Names help us understand what we're working with

Names can be used to coordinate different objects

```
rownames(factory) <- c("labor","steel")  
colnames(factory) <- c("cars","trucks")  
factory
```

```
##      cars trucks  
## labor   40     60  
## steel    1      3
```

```
available <- c(1600,70)  
names(available) <- c("labor","steel")
```

```
output <- c(20,10)
names(output) <- c("trucks","cars")
factory %*% output # But we've got cars and trucks mixed up!
```

```
##      [,1]
## labor 1400
## steel  50
```

```
factory %*% output[colnames(factory)]
```

```
##      [,1]
## labor 1600
## steel  70
```

```
all(factory %*% output[colnames(factory)] <= available[rownames(factory)])
```

```
## [1] TRUE
```

Notice: Last lines don't have to change if we add motorcycles as output or rubber and glass as inputs (abstraction again)

Take the mean: `rowMeans()`, `colMeans()`: input is matrix, output is vector. Also `rowSums()`, etc.

`summary()`: vector-style summary of column

```
colMeans(factory)
```

```
##   cars trucks  
##  20.5   31.5
```

```
summary(factory)
```

```
##           cars           trucks  
## Min.      : 1.00   Min.      : 3.00  
## 1st Qu.:10.75   1st Qu.:17.25  
## Median :20.50   Median :31.50  
## Mean    :20.50   Mean    :31.50  
## 3rd Qu.:30.25   3rd Qu.:45.75  
## Max.    :40.00   Max.    :60.00
```

`apply()`, takes 3 arguments: the array or matrix, then 1 for rows and 2 for columns, then name of the function to apply to each

```
rowMeans(factory)
```

```
## labor steel  
##      50      2
```

```
apply(factory,1,mean)
```

```
## labor steel  
##      50      2
```

What would `apply(factory,1,sd)` do?

Sequence of values, *not* necessarily all of the same type

```
my.distribution <- list("exponential",7,FALSE)
my.distribution
```

```
## [[1]]
## [1] "exponential"
##
## [[2]]
## [1] 7
##
## [[3]]
## [1] FALSE
```

Most of what you can do with vectors you can also do with lists

Can use `[]` as with vectors
or use `[[]]`, but only with a single index
`[[]]` drops names and structures, `[]` does not

```
is.character(my.distribution)
```

```
## [1] FALSE
```

```
is.character(my.distribution[[1]])
```

```
## [1] TRUE
```

```
my.distribution[[2]]^2
```

```
## [1] 49
```

What happens if you try `my.distribution[2]^2`? What happens if you try `[[]]` on a vector?

Add to lists with `c()` (also works with vectors):

```
my.distribution <- c(my.distribution,7)
my.distribution
```

```
## [[1]]
## [1] "exponential"
##
## [[2]]
## [1] 7
##
## [[3]]
## [1] FALSE
##
## [[4]]
## [1] 7
```

Chop off the end of a list by setting the length to something smaller (also works with vectors):

```
length(my.distribution)
```

```
## [1] 4
```

```
length(my.distribution) <- 3  
my.distribution
```

```
## [[1]]  
## [1] "exponential"  
##  
## [[2]]  
## [1] 7  
##  
## [[3]]  
## [1] FALSE
```

We can name some or all of the elements of a list

```
names(my.distribution) <- c("family","mean","is.symmetric")
my.distribution
```

```
## $family
## [1] "exponential"
##
## $mean
## [1] 7
##
## $is.symmetric
## [1] FALSE
```

```
my.distribution[["family"]]
```

```
## [1] "exponential"
```

```
my.distribution["family"]
```

```
## $family
## [1] "exponential"
```

Lists have a special short-cut way of using names, `$` (which removes names and structures):

```
my.distribution[["family"]]
```

```
## [1] "exponential"
```

```
my.distribution$family
```

```
## [1] "exponential"
```

Creating a list with names:

```
another.distribution <- list(family="gaussian",mean=7,sd=1,is.symmetric=TRUE)
```

Adding named elements:

```
my.distribution$was.estimated <- FALSE  
my.distribution[["last.updated"]] <- "2011-08-30"
```

Removing a named list element, by assigning it the value NULL:

```
my.distribution$was.estimated <- NULL
```

Lists give us a way to store and look up data by *name*, rather than by *position*

A really useful programming concept with many names: **key-value pairs**, **dictionaries**, **associative arrays**, **hashes**

If all our distributions have components named **family**, we can look that up by name, without caring where it is in the list

Dataframe = the classic data table, n rows for cases, p columns for variables

Lots of the really-statistical parts of R presume data frames **penn** from last time was really a dataframe

Not just a matrix because *columns can have different types*

Many matrix functions also work for dataframes (`rowSums()`, `summary()`, `apply()`)

but no matrix multiplying dataframes, even if all columns are numeric


```
a.matrix <- matrix(c(35,8,10,4),nrow=2)
colnames(a.matrix) <- c("v1","v2")
a.matrix
```

```
##      v1 v2
## [1,] 35 10
## [2,]  8  4
```

```
a.matrix[, "v1"] # Try a.matrix$v1 and see what happens
```

```
## [1] 35  8
```

```
a.data.frame <- data.frame(a.matrix,logicals=c(TRUE,FALSE))
a.data.frame
```

```
##   v1 v2 logicals
## 1 35 10      TRUE
## 2  8  4     FALSE
```

```
a.data.frame$v1
```

```
## [1] 35  8
```

```
a.data.frame[, "v1"]
```

```
## [1] 35  8
```

```
a.data.frame[1,]
```

```
##   v1 v2 logicals
## 1 35 10      TRUE
```

```
colMeans(a.data.frame)
```

```
##      v1      v2 logicals
## 21.5  7.0    0.5
```

We can add rows or columns to an array or data-frame with `rbind()` and `cbind()`, but be careful about forced type conversions

```
rbind(a.data.frame,list(v1=-3,v2=-5,logicals=TRUE))
```

```
##   v1 v2 logicals
## 1 35 10      TRUE
## 2  8  4     FALSE
## 3 -3 -5      TRUE
```

```
rbind(a.data.frame,c(3,4,6))
```

```
##   v1 v2 logicals
## 1 35 10         1
## 2  8  4         0
## 3  3  4         6
```

So far, every list element has been a single data value

List elements can be other data structures, e.g., vectors and matrices:

```
plan <- list(factory=factory, available=available, output=output)
plan$output
```

```
## trucks    cars
##      20     10
```

Internally, a dataframe is basically a list of vectors

List elements can even be other lists
which may contain other data structures
including other lists
which may contain other data structures...

This **recursion** lets us build arbitrarily complicated data structures from the basic ones

Most complicated objects are (usually) lists of data structures

`eigen()` finds eigenvalues and eigenvectors of a matrix
Returns a list of a vector (the eigenvalues) and a matrix (the eigenvectors)

```
eigen(factory)
```

```
## eigen() decomposition
## $values
## [1] 41.556171  1.443829
##
## $vectors
##           [,1]      [,2]
## [1,] 0.99966383 -0.8412758
## [2,] 0.02592747  0.5406062
```

```
class(eigen(factory))
```

```
## [1] "eigen"
```

With complicated objects, you can access parts of parts (of parts...)

```
factory %*% eigen(factory)$vectors[,2]
```

```
##           [,1]  
## labor -1.2146583  
## steel  0.7805429
```

```
eigen(factory)$values[2] * eigen(factory)$vectors[,2]
```

```
## [1] -1.2146583  0.7805429
```

```
eigen(factory)$values[2]
```

```
## [1] 1.443829
```

```
eigen(factory)[[1]][[2]] # NOT [[1,2]]
```

```
## [1] 1.443829
```

- ▶ Arrays add multi-dimensional structure to vectors
- ▶ Matrices act like you'd hope they would
- ▶ Lists let us combine different types of data
- ▶ Dataframes are hybrids of matrices and lists, for classic tabular data
- ▶ Recursion lets us build complicated data structures out of the simpler ones