

Lab 2: Indexing and Iteration

Name: Rufus Petrie

This week's agenda: basic indexing, with a focus on matrices; some more basic plotting; vectorization; using `for()` loops.

Q1. Back to some R basics

- **1a.** Let's start easy by working through some R basics, to continue to brush up on them. Define a variable `x.vec` to contain the integers 1 through 100. Check that it has length 100. Report the data type being stored in `x.vec`. Add up the numbers in `x.vec`, by calling a built-in R function. How many arithmetic operations did this take? **Challenge:** show how Gauss would have done this same calculation as a 7 year old, using just 3 arithmetic operations.

```
x.vec <- 1:100
length(x.vec)
```

```
## [1] 100
```

```
typeof(x.vec)
```

```
## [1] "integer"
```

```
sum(x.vec)
```

```
## [1] 5050
```

```
100*101/2
```

```
## [1] 5050
```

To find the sum, R has to perform 99 operations, but you can use Gauss' formula $\sum_{i=1}^n i = n(n+1)/2$ to find it in 3.

- **1b.** Convert `x.vec` into a matrix with 20 rows and 5 columns, and store this as `x.mat`. Here `x.mat` should be filled out in the default order (column major order). Check the dimensions of `x.mat`, and the data type as well. Compute the sums of each of the 5 columns of `x.mat`, by calling a built-in R function. Check (using a comparison operator) that the sum of column sums of `x.mat` equals the sum of `x.vec`.

```
# YOUR CODE GOES HERE
```

- **1c.** Extract and display rows 1, 5, and 17 of `x.mat`, with a single line of code. Answer the following questions, each with a single line of code: how many elements in row 2 of `x.mat` are larger than 40? How many elements in column 3 are in between 45 and 50? How many elements in column 5 are odd? Hint: take advantage of the `sum()` function applied to Boolean vectors.

```
# YOUR CODE GOES HERE
```

- **1d.** Using Boolean indexing, modify `x.vec` so that every even number in this vector is incremented by 10, and every odd number is left alone. This should require just a single line of code. Print out the result to the console. **Challenge:** show that `ifelse()` can be used to do the same thing, again using just a single line of code.

```
# YOUR CODE GOES HERE
```

- **1e.** Consider the list `x.list` created below. Complete the following tasks, each with a single line of code: extract all but the second element of `x.list`—seeking here a list as the final answer. Extract the first and third elements of `x.list`, then extract the second element of the resulting list—seeking here a vector as the final answer. Extract the second element of `x.list` as a vector, and then extract the first 10 elements of this vector—seeking here a vector as the final answer. Note: pay close attention to what is asked and use either single brackets `[]` or double brackets `[[]]` as appropriate.

```
x.list = list(rnorm(6), letters, sample(c(TRUE,FALSE),size=4,replace=TRUE))  
# YOUR CODE GOES HERE
```

Prostate cancer data set

We're going to look at a data set on 97 men who have prostate cancer (from the book *The Elements of Statistical Learning*). There are 9 variables measured on these 97 men:

1. `lpsa`: log PSA score
2. `lcavol`: log cancer volume
3. `lweight`: log prostate weight
4. `age`: age of patient
5. `lbph`: log of the amount of benign prostatic hyperplasia
6. `svi`: seminal vesicle invasion
7. `lcp`: log of capsular penetration
8. `gleason`: Gleason score
9. `pgg45`: percent of Gleason scores 4 or 5

To load this prostate cancer data set into your R session, and store it as a matrix `pros.dat`:

```
pros.dat =  
  as.matrix(read.table("http://www.stat.cmu.edu/~ryantibs/statcomp/data/pros.dat"))
```

Q2. Basic indexing and calculations

- **2a.** What are the dimensions of `pros.dat` (i.e., how many rows and how many columns)? Using integer indexing, print the first 6 rows and all columns; again using integer indexing, print the last 6 rows and all columns.

```
# YOUR CODE GOES HERE
```

- **2b.** Using the built-in R functions `head()` and `tail()` (i.e., do *not* use integer indexing), print the first 6 rows and all columns, and also the last 6 rows and all columns.

```
# YOUR CODE GOES HERE
```

- **2c.** Does the matrix `pros.dat` have names assigned to its rows and columns, and if so, what are they? Use `rownames()` and `colnames()` to find out. Note: these would have been automatically created by the `read.table()` function that we used above to read the data file into our R session. To see where `read.table()` would have gotten these names from, open up the data file: <http://www.stat.cmu.edu/~ryantibs/statcomp/data/pros.dat> in your web browser. Only the column names here are actually informative.

```
# YOUR CODE GOES HERE
```

- **2d.** Using named indexing, pull out the two columns of `pros.dat` that measure the log cancer volume and the log cancer weight, and store the result as a matrix `pros.dat.sub`. (Recall the explanation of variables at the top of this lab.) Check that its dimensions make sense to you, and that its first 6 rows are what you'd expect. Did R automatically assign column names to `pros.dat.sub`?

```
# YOUR CODE GOES HERE
```

- **2e.** Using the log cancer weights and log cancer volumes, calculate the log cancer density for the 97 men in the data set (note: $\text{density} = \text{weight} / \text{volume}$). There are in fact two different ways to do this; the first uses three function calls and one arithmetic operation; the second just uses one arithmetic operation. Note: in either case, you should be able to perform this computation for all 97 men *with a single line of code*, taking advantage of R's ability to vectorize. Write code to do it both ways, and show that both ways lead to the same answer, using `all.equal()`.

```
# YOUR CODE GOES HERE
```

- **2f.** Append the log cancer density to the columns of `pros.dat`, using `cbind()`. The new `pros.dat` matrix should now have 10 columns. Set the last column name to be `ldens`. Print its first 6 rows, to check that you've done all this right.

```
# YOUR CODE GOES HERE
```

Q3. Exploratory data analysis with plots

- **3a.** Using `hist()`, produce a histogram of the log cancer volume measurements of the 97 men in the data set; also produce a histogram of the log cancer weight. In each case, use `breaks=20` as an argument to `hist()`. Comment just briefly on the distributions you see. Then, using `plot()`, produce a scatterplot of the log cancer volume (y-axis) versus the log cancer weight (x-axis). Do you see any kind of relationship? Would you expect to? **Challenge:** how would you measure the strength of this relationship formally? Note that there is certainly more than one way to do so. We'll talk about statistical modeling tools later in the course.

```
# YOUR CODE GOES HERE
```

- **3b.** Produce scatterplots of log cancer weight versus age, and log cancer volume versus age. Do you see relationships here between the age of a patient and the volume/weight of his cancer?

```
# YOUR CODE GOES HERE
```

- **3c.** Produce a histogram of the log cancer density, and a scatterplot of the log cancer density versus age. Comment on any similarities/differences you see between these plots, and the corresponding ones you produced above for log cancer volume/weight.

```
# YOUR CODE GOES HERE
```

- **3d.** Delete the last column, corresponding to the log cancer density, from the `pros.dat` matrix, using negative integer indexing.

```
# YOUR CODE GOES HERE
```

Q4. A bit of Boolean indexing never hurt anyone

- **4a.** The `svi` variable in the `pros.dat` matrix is binary: 1 if the patient had a condition called "seminal vesicle invasion" or SVI, and 0 otherwise. SVI (which means, roughly speaking, that the cancer invaded into the muscular wall of the seminal vesicle) is bad: if it occurs, then it is believed the prognosis for the patient is poorer, and even once/if recovered, the patient is more likely to have prostate cancer return in the future. Compute a Boolean vector called `has.svi`, of length 97, that has a `TRUE` element if a row (patient) in `pros.dat` has SVI, and `FALSE` otherwise. Then using `sum()`, figure out how many patients have SVI.

```
# YOUR CODE GOES HERE
```

- **4b.** Extract the rows of `pros.dat` that correspond to patients with SVI, and the rows that correspond to patients without it. Call the resulting matrices `pros.dat.svi` and `pros.dat.no.svi`, respectively.

You can do this in two ways: using the `has.svi` Boolean vector created above, or using on-the-fly Boolean indexing, it's up to you. Check that the dimensions of `pros.dat.svi` and `pros.dat.no.svi` make sense to you.

```
# YOUR CODE GOES HERE
```

- **4c.** Using the two matrices `pros.dat.svi` and `pros.dat.no.svi` that you created above, compute the means of each variable in our data set for patients with SVI, and for patients without it. Store the resulting means into vectors called `pros.dat.svi.avg` and `pros.dat.no.svi.avg`, respectively. Hint: for each matrix, you can compute the means with a single call to a built-in R function. What variables appear to have different means between the two groups?

```
# YOUR CODE GOES HERE
```

Q5. Computing standard deviations using iteration

- **5a.** Take a look at the starter code below. The first line defines an empty vector `pros.dat.svi.sd` of length `ncol(pros.dat)` (of length 9). The second line defines an index variable `i` and sets it equal to 1. Write a third line of code to compute the standard deviation of the `i`th column of `pros.dat.svi`, using a built-in R function, and store this value in the `i`th element of `pros.dat.svi.sd`.

```
pros.dat.svi.sd = vector(length=ncol(pros.dat))
i = 1
# YOUR CODE GOES HERE
```

- **5b.** Repeat the calculation as in the previous question, but for patients without SVI. That is, produce three lines of code: the first should define an empty vector `pros.dat.no.svi.sd` of length `ncol(pros.dat)` (of length 9), the second should define an index variable `i` and set it equal to 1, and the third should fill the `i`th element of `pros.dat.no.svi.sd` with the standard deviation of the `i`th column of `pros.dat.no.svi`.

```
pros.dat.no.svi.sd = vector(length=ncol(pros.dat))
i = 1
# YOUR CODE GOES HERE
```

- **5c.** Write a `for()` loop to compute the standard deviations of the columns of `pros.dat.svi` and `pros.dat.no.svi`, and store the results in the vectors `pros.dat.svi.sd` and `pros.dat.no.svi.sd`, respectively, that were created above. Note: you should have a single `for()` loop here, not two `for` loops. And if it helps, consider breaking this task down into two steps: as the first step, write a `for()` loop that iterates an index variable `i` over the integers between 1 and the number of columns of `pros.dat` (don't just manually write 9 here, pull out the number of columns programmatically), with an empty body. As the second step, paste relevant pieces of your solution code from Q5a and Q5b into the body of the `for()` loop. Print out the resulting vectors `pros.dat.svi.sd` and `pros.dat.no.svi.sd` to the console. Comment, just briefly (informally), by visually inspecting these standard deviations and the means you computed in Q4c: which variables exhibit large differences in means between the SVI and non-SVI patients, relative to their standard deviations?

```
# YOUR CODE GOES HERE
```

- **5d.** The code below computes the standard deviations of the columns of `pros.dat.svi` and `pros.dat.no.svi`, and stores them in `pros.dat.svi.sd.master` and `pros.dat.no.svi.sd.master`, respectively, using `apply()`. (We'll learn `apply()` and related functions a bit later in the course.) Remove `eval=FALSE` as an option to the Rmd code chunk, and check using `all.equal()` that the standard deviations you computed in the previous question equal these "master" copies. Note: use `check.names=FALSE` as a third argument to `all.equal()`, which instructs it to ignore the names of its first two arguments. (If `all.equal()` doesn't succeed in both cases, then you must have done something wrong in computing the standard deviations, so go back and fix them!)

```
pros.dat.svi.sd.master = apply(pros.dat.svi, 2, sd)
pros.dat.no.svi.sd.master = apply(pros.dat.no.svi, 2, sd)
all.equal(pros.dat.svi.sd, pros.dat.svi.sd.master, check.names=FALSE)
all.equal(pros.dat.no.svi.sd, pros.dat.no.svi.sd.master, check.names=FALSE)
```

Q6. Computing t-tests using vectorization

- **6a.** Recall that the **two-sample (unpaired) t-statistic** between data sets $X = (X_1, \dots, X_n)$ and $Y = (Y_1, \dots, Y_m)$ is:

$$T = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{s_X^2}{n} + \frac{s_Y^2}{m}}},$$

where $\bar{X} = \sum_{i=1}^n X_i/n$ is the sample mean of X , $s_X^2 = \sum_{i=1}^n (X_i - \bar{X})^2/(n-1)$ is the sample variance of X , and similarly for \bar{Y} and s_Y^2 . We will compute these t-statistics for all 9 variables in our data set, where X will play the role of one of the variables for SVI patients, and Y will play the role of this variable for non-SVI patients. Start by computing a vector of the denominators of the t-statistics, called `pros.dat.denom`, according to the formula above. Take advantage of vectorization; this calculation should require just a single line of code. Make sure not to include any hard constants (e.g., don't just manually write 21 here for n); as always, programmatically define all the relevant quantities. Then compute a vector of t-statistics for the 9 variables in our data set, called `pros.dat.t.stat`, according to the formula above, and using `pros.dat.denom`. Again, take advantage of vectorization; this calculation should require just a single line of code. Print out the t-statistics to the console.

YOUR CODE GOES HERE

- **6b.** Given data X and Y and the t-statistic T as defined the last question, the **degrees of freedom** associated with T is:

$$\nu = \frac{\left(\frac{s_X^2}{n} + \frac{s_Y^2}{m}\right)^2}{\frac{\left(\frac{s_X^2}{n}\right)^2}{n-1} + \frac{\left(\frac{s_Y^2}{m}\right)^2}{m-1}}.$$

Compute the degrees of freedom associated with each of our 9 t-statistics (from our 9 variables), storing the result in a vector called `pros.dat.df`. This might look like a complicated/ugly calculation, but really, it's not too bad: it only involves arithmetic operators, and taking advantage of vectorization, the calculation should only require a single line of code. Hint: to simplify this line of code, it will help to first set short variable names for variables/quantities you will be using, as in `sx = pros.dat.svi.sd`, `n = nrow(pros.dat.svi)`, and so on. Print out these degrees of freedom values to the console.

YOUR CODE GOES HERE

- **6c.** The function `pt()` evaluates the distribution function of the t-distribution. E.g.,

```
pt(x, df=v, lower.tail=FALSE)
```

returns the probability that a t-distributed random variable, with v degrees of freedom, exceeds the value x . Importantly, `pt()` is vectorized: if x is a vector, and so is v , then the above returns, in vector format: the probability that a t-distributed variate with $v[1]$ degrees of freedom exceeds $x[1]$, the probability that a t-distributed variate with $v[2]$ degrees of freedom exceeds $x[2]$, and so on.

Call `pt()` as in the above line, but replace x by the absolute values of the t-statistics you computed for the 9 variables in our data set, and v by the degrees of freedom values associated with these t-statistics. Multiply the output by 2, and store it as a vector `pros.dat.p.val`. These are called **p-values** for the t-tests of mean difference between SVI and non-SVI patients, over the 9 variables in our data set. Print out the p-values to the console. Identify the variables for which the p-value is smaller than 0.05 (hence deemed to have a significant difference between SVI and non-SVI patients). Identify the variable with the smallest p-value (the most significant difference between SVI and non-SVI patients).

```
# YOUR CODE GOES HERE
```

Q7. My plot is at your command (optional)

- **Challenge.** Use the last code block from this week’s lecture as starter code to complete the following task. In the body of a **repeat** loop, prompt the user for a variable name to plot, using **readline()**. Check if the string that you collect from the user is one of the column names in **pros.dat**. Hint: use the **%in%** operator. If the string is indeed one of the column names, produce a histogram of the corresponding variable, with a title and x-axis label set appropriately. If the string is “quit”, then break out of the repeat loop. Otherwise, print to the console: “Oops! That’s not a variable in my data set.” In the Rmd code chunk for your solution code, *make sure to set **eval=FALSE**; otherwise your lab file will never finish knitting*. Try out your solution code by running it in your console.

```
# YOUR CODE GOES HERE
```

- **Challenge.** Extend your prompting code in the last question to allow for scatterplots as well as histograms, and any other options you deem interesting, that the user might want to specify.

```
# YOUR CODE GOES HERE
```