

Lecture 04 exercises solutions

Becky Lin

Short exercises based on chapters from text

1. Explain the output of the following code chunk.

```
f <- function() {
  fe <- environment(f)
  ee <- environment()
  pe <- parent.env(ee)
  list(fe=fe,ee=ee,pe=pe)
}
f()
```

```
## $fe
## <environment: R_GlobalEnv>
##
## $ee
## <environment: 0x7fdc103dad80>
##
## $pe
## <environment: R_GlobalEnv>
```

Solution The variable `fe` holds the environment of the function `f`, which is the environment in which `f` was defined. This is the global environment in this example. The variable `pe` holds the parent environment of the environment inside `f`, which is the global environment. The variable `ee` holds the environment within the function.

2. Read the help files on the `exists()` and `get()` functions. Explain the output of the following code chunk.

```
f <- function(xx) {
  xx_parent <- if(exists("xx",envir=environment(f))) {
    get("xx",environment(f))
  } else {
    NULL
  }
  list(xx,xx_parent)
}
f(2)
```

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

```
xx <- 1
f(2)
```

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

```
rm(xx)
```

Solution On the first call to `f` `xx` does not exist in the global environment, so `xx_parent` is `NULL` in `f`. On the second call, `xx` does exist so `xx_parent` is its value in the global environment.

3. Write a function with argument `xx` that tests whether `xx` exists in the parent environment and, if so,
 - a. assigns the value of `xx` in the parent environment to the variable `xx_parent`, and
 - b. tests whether `xx` and `xx_parent` are equal. If the test is `FALSE`, throw a warning to alert the user to the fact that the two are not equal.

Solution

```
f <- function(xx) {
  if(exists("xx",envir=environment(f))) {
    xx_parent <- get("xx",envir=environment(f))
    if(xx_parent != xx) {
      warning("xx value in global env is ",xx_parent," and xx value in f is ",xx,"\n")
    }
  }
}
f(1)
xx <- 2
f(1)
```

```
## Warning in f(1): xx value in global env is 2 and xx value in f is 1
```

```
rm(xx)
```

4. Write an infix version of `c()` that concatenates two vectors.

```
`%c%` <- function(x,y) c(x,y)
1:3 %c% 4:6
```

```
## [1] 1 2 3 4 5 6
```

Map-Reduce

Refer to the Wikipedia page on algorithms for computing sample variances and covariances:

https://en.wikipedia.org/wiki/Algorithms_for_calculating_variance

(https://en.wikipedia.org/wiki/Algorithms_for_calculating_variance) We will implement the two-pass algorithm for computing the sample variance as a Map-Reduce.

1. Use the following code chunk from the lecture 4 notes to simulate data in 10 chunks and calculate the overall sample mean.

```
rfunc <- function(seed,n) { set.seed(seed); return(rnorm(n)) }
dat2 <- lapply(1:10,rfunc,n=1e5) # apply rfunc function with n=1e5 over 1:10
mfunc <- function(x) {
  return( data.frame(sum=sum(x),n=length(x)) ) }
sumdat <- lapply(dat2,mfunc) # apply mfunc to each element of list "dat2"
simple_reduce <- function(x, f) { # Text section 9.5
  out <- x[[1]]
  for (i in seq(2, length(x))) {
    out <- f(out, x[[i]])
  }
  out
}
allres <- simple_reduce(sumdat, rbind)
my_mean <- sum(allres[, "sum"])/sum(allres[, "n"]) # mean
```

2. Write a function that takes a vector and `my_mean` as input and returns
 - i. the sum of squared deviations between the vector's values and `my_mean` and
 - ii. the number of vector values (n).

Solution

```
mfunc3 <- function(x,my_mean) {
  return(data.frame(sum=sum((x-my_mean)^2),n=length(x)))
}
```

3. Use `lapply()` to call your function from (2) on each element of `dat2`.

Solution

```
sumdat <- lapply(dat2,mfunc3,my_mean=my_mean)
```

4. Use `simple_reduce()` to combine your results into a single data frame, and calculate the sample variance from this data frame. Compare your answer to `var(unlist(dat2))`.

Solution

```
allres <- simple_reduce(sumdat, rbind)
my_var <- sum(allres[, "sum"])/(sum(allres[, "n"])-1)
var(unlist(dat2)) # same
```

```
## [1] 1.003756
```

Recursive partitioning

- The following code chunk is the start of an implementation of recursive partitioning using a binary tree data structure to store the partition.
- Binary trees can be implemented as a linked list of nodes that contain
 1. data
 2. a pointer to the left child
 3. a pointer to the right child
- For our recursive partitioning example, the data will be a region of the original covariate space and the response/covariate data in that region.
- The following code chunk establishes node and region data structure.

```
# Constructor for the node data structure:
new_node <- function(data,childl=NULL,childr=NULL){
  nn <- list(data=data,childl=childl,childr=childr) # note that input data is a list
  class(nn) <- "node"
  return(nn)
}
# The data stored in the node are a partition, or region of the
# covariate space. Constructor for region data structure:
new_region <- function(coords=NULL,x,y){
  if(is.null(coords)) {
    # find range for each column in x and column binding result
    coords <- Reduce(cbind,lapply(x,range) )
  }
  out <- list(coords=coords,x=x,y=y)
  class(out) <- "region"
  return(out)
}
```

- Some tests of the above constructors are given in the next code chunk.

```
set.seed(123); n <- 10
x <- data.frame(x1=rnorm(n),x2=rnorm(n))
y <- rnorm(n)
new_region(x=x,y=y)
```

```
## $coords
##           init
## [1,] -1.265061 -1.966617
## [2,]  1.715065  1.786913
##
## $x
##           x1           x2
## 1 -0.56047565  1.2240818
## 2 -0.23017749  0.3598138
## 3  1.55870831  0.4007715
## 4  0.07050839  0.1106827
## 5  0.12928774 -0.5558411
## 6  1.71506499  1.7869131
## 7  0.46091621  0.4978505
## 8 -1.26506123 -1.9666172
## 9 -0.68685285  0.7013559
## 10 -0.44566197 -0.4727914
##
## $y
## [1] -1.0678237 -0.2179749 -1.0260044 -0.7288912 -0.6250393 -1.6866933
## [7]  0.8377870  0.1533731 -1.1381369  1.2538149
##
## attr(,"class")
## [1] "region"
```

```
new_node(new_region(x=x,y=y))
```

```

## $data
## $coords
##          init
## [1,] -1.265061 -1.966617
## [2,]  1.715065  1.786913
##
## $x
##          x1          x2
## 1 -0.56047565  1.2240818
## 2 -0.23017749  0.3598138
## 3  1.55870831  0.4007715
## 4  0.07050839  0.1106827
## 5  0.12928774 -0.5558411
## 6  1.71506499  1.7869131
## 7  0.46091621  0.4978505
## 8 -1.26506123 -1.9666172
## 9 -0.68685285  0.7013559
## 10 -0.44566197 -0.4727914
##
## $y
## [1] -1.0678237 -0.2179749 -1.0260044 -0.7288912 -0.6250393 -1.6866933
## [7]  0.8377870  0.1533731 -1.1381369  1.2538149
##
## attr(,"class")
## [1] "region"
##
## $child1
## NULL
##
## $childr
## NULL
##
## attr(,"class")
## [1] "node"

```

```
new_node(new_region(x=x,y=y))$data
```

```
## $coords
##           init
## [1,] -1.265061 -1.966617
## [2,]  1.715065  1.786913
##
## $x
##           x1           x2
## 1  -0.56047565  1.2240818
## 2  -0.23017749  0.3598138
## 3   1.55870831  0.4007715
## 4   0.07050839  0.1106827
## 5   0.12928774 -0.5558411
## 6   1.71506499  1.7869131
## 7   0.46091621  0.4978505
## 8  -1.26506123 -1.9666172
## 9  -0.68685285  0.7013559
## 10 -0.44566197 -0.4727914
##
## $y
## [1] -1.0678237 -0.2179749 -1.0260044 -0.7288912 -0.6250393 -1.6866933
## [7]  0.8377870  0.1533731 -1.1381369  1.2538149
##
## attr(,"class")
## [1] "region"
```

- The recursive partitioning function is shown below. We'll discuss this in class.

```

#-----#
# Recursive partitioning function.
recpart <- function(x,y){
  init <- new_node(new_region(x=x,y=y))
  tree <- recpart_recursive(init)
  class(tree) <- c("tree",class(tree))
  return(tree)
}
recpart_recursive <- function(node) {
  R <- node$data
  # stop recursion if region has a single data point
  if(length(R$y) == 1) { return(node) } # NB: was return(NULL)
  # else find a split that minimizes a LOF criterion
  # Initialize
  lof_best <- Inf
  # Loop over variables and splits
  for(v in 1:ncol(R$x)){
    tt <- split_points(R$x[,v]) # Exercise: write split_points()
    for(t in tt) {
      gdat <- data.frame(y=R$y,x=as.numeric(R$x[,v] <= t))
      lof <- LOF(y~.,gdat) # Exercise: write LOF()
      if(lof < lof_best) {
        lof_best <- lof
        childRs <- split(R,xvar=v,spt=t) # Exercises: write split.region()
      }
    }
  }
  # Call self on best split
  node$childl <- recpart_recursive(new_node(childRs$Rl))
  node$childr <- recpart_recursive(new_node(childRs$Rr))
  return(node)
}

```

Exercises

1. Write `split_points()`. The function should take a vector of covariate values as input and return the sorted unique values. You will need to trim off the maximum unique value, because this can't be used as a split point. (As yourself why not.) Write a snippet of R code that tests your function.
2. Write the function `LOF()` that returns the lack-of-fit criterion for a model. The function should take a model formula and data frame as input, pass these to `lm()` and return the residual sum of squares. Write a snippet of R code that tests your function.
3. Write `split.region()`. The function should take a region `R`, the variable to split on, `v`, and the split point, `t`, as arguments. Split the region into left and right partitions and return a list of two regions labelled `Rl` and `Rr`. Note: It is tempting to split the `x` and `y` data and calculate the coordinates matrix from the `x`'s, as the constructor does when not passed a coordinates matrix. However, this will leave gaps in the covariate space. (Ask yourself why.) Write a snippet of R code that tests your function.
4. Run `recpart()` with your versions of `split_points()`, `LOF()` and `split.region()`. Use the test data `x` and `y` defined in the testing code chunk. At this point you do not need to check that the output is correct; you will get a chance to do that in lab 3.


```

split_points <- function(x) {
  x <- sort(unique(x))
  x <- x[-length(x)] # remove the last element which is the max in the list of unique values
  return(x)
}

testx <- c(2,4,2,3,7,7,5,5,8) # 9 elements with 6 unique values
split_points(testx)

```

```
## [1] 2 3 4 5 7
```

```

LOF <- function(form,data) {
  ff <- lm(form,data)
  return(sum(residuals(ff)^2))
}

set.seed(360)
xvar <- rnorm(20,1,1)
yvar <- 1+2*xvar+rnorm(20,0,1)
dxy <- as.data.frame(cbind(xvar=xvar,yvar=yvar))
LOF(yvar~xvar,dxy)

```

```
## [1] 20.67728
```

```

split.region <- function(R,xvar,spt){
  r1_ind <- (R$x[,xvar] <= spt)
  c1 <- c2 <- R$coords
  c1[2,xvar] <- spt; c2[1,xvar] <- spt
  Rl <- new_region(c1,R$x[r1_ind,,drop=FALSE],R$y[r1_ind])
  Rr <- new_region(c2,R$x[!r1_ind,,drop=FALSE],R$y[!r1_ind])
  return(list(Rl=Rl,Rr=Rr))
}

init <- new_node(new_region(x=x,y=y))
init

```

```
## $data
## $coords
##           init
## [1,] -1.265061 -1.966617
## [2,]  1.715065  1.786913
##
## $x
##           x1           x2
## 1 -0.56047565  1.2240818
## 2 -0.23017749  0.3598138
## 3  1.55870831  0.4007715
## 4  0.07050839  0.1106827
## 5  0.12928774 -0.5558411
## 6  1.71506499  1.7869131
## 7  0.46091621  0.4978505
## 8 -1.26506123 -1.9666172
## 9 -0.68685285  0.7013559
## 10 -0.44566197 -0.4727914
##
## $y
## [1] -1.0678237 -0.2179749 -1.0260044 -0.7288912 -0.6250393 -1.6866933
## [7]  0.8377870  0.1533731 -1.1381369  1.2538149
##
## attr(,"class")
## [1] "region"
##
## $child1
## NULL
##
## $childr
## NULL
##
## attr(,"class")
## [1] "node"
```

```
test <- recpart(x,y)
test$child1$child1$childr
```

```
## $data
## $coords
##           init
## [1,] -1.265061 -1.9666172
## [2,] -0.445662  0.4978505
##
## $x
##           x1           x2
## 10 -0.445662 -0.4727914
##
## $y
## [1] 1.253815
##
## attr(,"class")
## [1] "region"
##
## $childl
## NULL
##
## $childr
## NULL
##
## attr(,"class")
## [1] "node"
```

```
test$childl$childl
```

```
## $data
## $coords
##          init
## [1,] -1.265061 -1.9666172
## [2,] -0.445662  0.4978505
##
## $x
##          x1          x2
## 8  -1.265061 -1.9666172
## 10 -0.445662 -0.4727914
##
## $y
## [1] 0.1533731 1.2538149
##
## attr(,"class")
## [1] "region"
##
## $child1
## $data
## $coords
##          init
## [1,] -1.265061 -1.9666172
## [2,] -1.265061  0.4978505
##
## $x
##          x1          x2
## 8  -1.265061 -1.966617
##
## $y
## [1] 0.1533731
##
## attr(,"class")
## [1] "region"
##
## $child1
## NULL
##
## $childr
## NULL
##
## attr(,"class")
## [1] "node"
##
## $childr
## $data
## $coords
##          init
## [1,] -1.265061 -1.9666172
## [2,] -0.445662  0.4978505
##
## $x
##          x1          x2
```

```
## 10 -0.445662 -0.4727914
##
## $y
## [1] 1.253815
##
## attr(,"class")
## [1] "region"
##
## $child1
## NULL
##
## $childr
## NULL
##
## attr(,"class")
## [1] "node"
##
## attr(,"class")
## [1] "node"
```

```
plot(test$data$x)
dc = test$data$coords
rect(dc[1,1],dc[1,2],dc[2,1],dc[2,2],border="red")
dc = test$child1$child1$data$coords
rect(dc[1,1],dc[1,2],dc[2,1],dc[2,2],border="blue")
dc = test$child1$childr$data$coords
rect(dc[1,1],dc[1,2],dc[2,1],dc[2,2],border="green")
```

