

# Bios 6301: Assignment 6

Lydia Yao

*Due Tuesday, 26 October, 1:00 PM*

$5^{n=\text{day}}$  points taken off for each day late.

40 points total.

Submit a single knitr file (named `homework6.rmd`), along with a valid PDF output file. Inside the file, clearly indicate which parts of your responses go with which problems (you may use the original homework document as a template). Add your name as `author` to the file's metadata section. Raw R code/output or word processor files are not acceptable.

Failure to name file `homework6.rmd` or include author name may result in 5 points taken off.

## Question 1

### 16 points

Obtain a copy of the football-values lecture. Save the five 2021 CSV files in your working directory.

Modify the code to create a function. This function will create dollar values given information (as arguments) about a league setup. It will return a data.frame and write this data.frame to a CSV file. The final data.frame should contain the columns 'PlayerName', 'pos', 'points', 'value' and be ordered by value descendingly. Do not round dollar values.

Note that the returned data.frame should have `sum(posReq)*nTeams` rows.

Define the function as such (10 points):

```
# path: directory path to input files
# file: name of the output file; it should be written to path
# nTeams: number of teams in league
# cap: money available to each team
# posReq: number of starters for each position
# points: point allocation for each category
ffvalues <- function(path, file='outfile.csv', nTeams=12, cap=200, posReq=c(qb=1, rb=2, wr=3, te=1, k=1,
                                points=c(fg=4, xpt=1, pass_yds=1/25, pass_tds=4, pass_ints=-2,
                                rush_yds=1/10, rush_tds=6, fumbles=-2, rec_yds=1/20, rec_tds=6))) {

  ## read in CSV files
  positions <- c('k','qb','rb','te','wr')
  csvfile <- paste('proj_', positions, substr(2021, 3, 4), '.csv', sep='')
  files <- file.path(2021, csvfile)
  names(files) <- positions
  k <- read.csv(files['k'], header=TRUE, stringsAsFactors=FALSE)
  qb <- read.csv(files['qb'], stringsAsFactors=FALSE)
  rb <- read.csv(files['rb'])
  te <- read.csv(files['te'])
  wr <- read.csv(files['wr'])
  cols <- unique(c(names(k), names(qb), names(rb), names(te), names(wr)))
  k[, 'pos'] <- 'k'
```

```

qb[, 'pos'] <- 'qb'
rb[, 'pos'] <- 'rb'
te[, 'pos'] <- 'te'
wr[, 'pos'] <- 'wr'
cols <- c(cols, 'pos')
k[, setdiff(cols, names(k))] <- 0
qb[, setdiff(cols, names(qb))] <- 0
rb[, setdiff(cols, names(rb))] <- 0
te[, setdiff(cols, names(te))] <- 0
wr[, setdiff(cols, names(wr))] <- 0
x <- rbind(k[, cols], qb[, cols], rb[, cols], te[, cols], wr[, cols])
## calculate dollar values
for (i in names(points)){
  x[, paste('p_', i, sep='')] <- x[, i]*points[i]
}
x[, 'points'] <- rowSums(x[, grep("^p_", names(x))])
x2 <- x[order(x[, 'points'], decreasing=TRUE),]
k.ix <- which(x2[, 'pos']=='k')
qb.ix <- which(x2[, 'pos']=='qb')
rb.ix <- which(x2[, 'pos']=='rb')
te.ix <- which(x2[, 'pos']=='te')
wr.ix <- which(x2[, 'pos']=='wr')
if(posReq['k'] == 0) {x2[k.ix, 'marg'] <- x2[k.ix, 'points']}
} else{ x2[k.ix, 'marg'] <- x2[k.ix, 'points'] - x2[k.ix[nTeams* posReq['k']], 'points']}
x2[qb.ix, 'marg'] <- x2[qb.ix, 'points'] - x2[qb.ix[nTeams* posReq['qb']], 'points']
x2[rb.ix, 'marg'] <- x2[rb.ix, 'points'] - x2[rb.ix[nTeams* posReq['rb']], 'points']
x2[te.ix, 'marg'] <- x2[te.ix, 'points'] - x2[te.ix[nTeams* posReq['te']], 'points']
x2[wr.ix, 'marg'] <- x2[wr.ix, 'points'] - x2[wr.ix[nTeams* posReq['wr']], 'points']
x3 <- x2[x2[, 'marg'] >= 0,]
x3 <- x3[order(x3[, 'marg'], decreasing=TRUE),]
rownames(x3) <- NULL
x3[, 'value'] <- (nTeams*cap-nrow(x3)) * x3[, 'marg'] / sum(x3[, 'marg']) + 1
x4 <- x3[, c('PlayerName', 'pos', 'points', 'marg', 'value')]
## save dollar values as CSV file
write.csv(x, file = file, row.names = TRUE)
## return data.frame with dollar values
return (x4)
}

```

1. Call `x1 <- ffvalues('.')`

```
x1 <- ffvalues('.')
```

1. How many players are worth more than \$20? (1 point)

```
print(nrow(x1[x1$value > 20, ]))
```

```
## [1] 44
```

1. Who is 15th most valuable running back (rb)? (1 point)

```
print(x1[x1[, 'pos'] %in% c('rb'), ][15,]$PlayerName)
```

```
## [1] "Chris Carson"
```

1. Call `x2 <- ffvalues(getwd(), '16team.csv', nTeams=16, cap=150)`

```
x2 <- ffvalues(getwd(), '16team.csv', nTeams = 16, cap=150)
```

1. How many players are worth more than \$20? (1 point)

```
print(nrow(x2[x2$value > 20, ]))
```

```
## [1] 44
```

1. How many wide receivers (wr) are in the top 40? (1 point)

```
x3 <- x2[1:40,]
nrow(x3[x3[, 'pos'] %in% c('wr'),])
```

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

1. Call:

```
x3 <- ffvalues('.', 'qbheavy.csv', posReq=c(qb=2, rb=2, wr=3, te=1, k=0),
           points=c(fg=0, xpt=0, pass_yds=1/25, pass_tds=6, pass_ints=-2,
                    rush_yds=1/10, rush_tds=6, fumbles=-2, rec_yds=1/20, rec_tds=6))
```

1. How many players are worth more than \$20? (1 point)

```
print(nrow(x3[x3$value > 20, ]))
```

```
## [1] 46
```

1. How many quarterbacks (qb) are in the top 30? (1 point)

```
x4 <- x3[1:30,]
nrow(x4[x4[, 'pos'] %in% c('qb'),])
```

```
## [1] 14
```

## Question 2

### 24 points

Import the HAART dataset (haart.csv) from the GitHub repository into R, and perform the following manipulations: (4 points each) 1. Convert date columns into a usable (for analysis) format. Use the `table` command to display the counts of the year from `init.date`.

```
df <- read.csv('haart.csv', header=TRUE, stringsAsFactors=FALSE)
df$init.date <- as.POSIXct(df$init.date, format = "%m/%d/%y")
df$last.visit <- as.POSIXct(df$last.visit, format = "%m/%d/%y")
df$date.death <- as.POSIXct(df$date.death, format = "%m/%d/%y")
table(df$init.date)
```

```
##
## 1998-04-08 2000-05-01 2000-06-01 2000-06-15 2000-10-01 2000-11-23 2001-01-15
##          1          1          1          1          1          1          1
## 2001-02-14 2001-02-26 2001-03-09 2001-04-01 2001-04-22 2001-06-12 2001-08-01
##          1          1          1          1          1          1          1
## 2001-08-15 2001-09-01 2001-09-06 2001-11-01 2001-11-08 2001-11-21 2001-12-19
##          1          1          1          1          1          1          1
## 2001-12-20 2001-12-28 2002-01-01 2002-01-07 2002-01-08 2002-01-21 2002-02-02
##          1          1          1          1          1          1          1
## 2002-02-14 2002-02-18 2002-02-20 2002-02-26 2002-02-27 2002-02-28 2002-03-06
##          1          1          1          1          1          1          1
## 2002-03-07 2002-03-09 2002-03-15 2002-03-21 2002-03-24 2002-04-01 2002-04-05
```

##	1	1	1	1	1	1	1
##	2002-04-18	2002-04-26	2002-04-30	2002-05-15	2002-06-11	2002-06-15	2002-06-21
##	1	1	1	1	1	1	1
##	2002-06-24	2002-06-28	2002-07-04	2002-07-05	2002-07-18	2002-07-23	2002-07-30
##	1	1	1	1	1	2	1
##	2002-08-02	2002-08-05	2002-08-06	2002-08-13	2002-08-14	2002-08-17	2002-08-27
##	2	2	1	1	1	2	2
##	2002-09-11	2002-10-02	2002-10-11	2002-10-15	2002-10-18	2002-10-31	2002-11-11
##	1	1	1	1	1	1	1
##	2002-11-26	2002-11-30	2002-12-05	2002-12-06	2002-12-09	2002-12-16	2002-12-26
##	1	1	1	1	1	1	2
##	2003-01-01	2003-01-02	2003-01-07	2003-01-09	2003-01-16	2003-01-17	2003-01-20
##	1	1	1	1	1	1	1
##	2003-01-28	2003-02-01	2003-02-04	2003-02-05	2003-02-06	2003-02-11	2003-02-13
##	1	1	1	2	1	2	1
##	2003-02-20	2003-02-25	2003-02-26	2003-03-01	2003-03-06	2003-03-12	2003-03-13
##	1	2	4	1	1	1	2
##	2003-03-15	2003-03-20	2003-03-24	2003-03-25	2003-03-26	2003-03-28	2003-04-01
##	1	2	1	3	1	1	2
##	2003-04-02	2003-04-04	2003-04-07	2003-04-08	2003-04-14	2003-04-15	2003-04-24
##	1	2	1	1	2	2	1
##	2003-04-29	2003-04-30	2003-05-01	2003-05-02	2003-05-05	2003-05-06	2003-05-07
##	1	3	1	1	1	1	3
##	2003-05-08	2003-05-09	2003-05-10	2003-05-14	2003-05-15	2003-05-20	2003-05-21
##	3	1	1	1	2	1	2
##	2003-05-22	2003-05-26	2003-05-27	2003-05-28	2003-05-29	2003-05-30	2003-06-01
##	4	2	3	4	2	3	1
##	2003-06-03	2003-06-04	2003-06-06	2003-06-09	2003-06-12	2003-06-13	2003-06-15
##	1	2	2	2	1	1	2
##	2003-06-16	2003-06-20	2003-06-23	2003-06-24	2003-06-30	2003-07-01	2003-07-02
##	1	1	1	2	2	5	5
##	2003-07-03	2003-07-04	2003-07-07	2003-07-08	2003-07-11	2003-07-14	2003-07-15
##	4	2	1	2	2	2	1
##	2003-07-16	2003-07-18	2003-07-19	2003-07-21	2003-07-22	2003-07-24	2003-07-28
##	4	4	1	3	1	1	2
##	2003-07-30	2003-07-31	2003-08-04	2003-08-05	2003-08-06	2003-08-08	2003-08-12
##	1	1	1	2	2	1	2
##	2003-08-16	2003-08-18	2003-08-20	2003-08-22	2003-08-25	2003-08-26	2003-08-27
##	1	2	3	1	1	1	1
##	2003-08-28	2003-08-29	2003-09-01	2003-09-02	2003-09-05	2003-09-08	2003-09-09
##	1	4	4	1	1	5	2
##	2003-09-10	2003-09-11	2003-09-15	2003-09-16	2003-09-17	2003-09-18	2003-09-19
##	2	3	6	3	1	2	1
##	2003-09-22	2003-09-23	2003-09-24	2003-09-29	2003-10-01	2003-10-02	2003-10-03
##	1	2	2	3	1	1	1
##	2003-10-04	2003-10-06	2003-10-07	2003-10-10	2003-10-12	2003-10-14	2003-10-15
##	1	1	1	1	1	1	3
##	2003-10-16	2003-10-17	2003-10-21	2003-10-22	2003-10-23	2003-10-24	2003-10-28
##	3	1	2	1	1	3	1
##	2003-10-29	2003-11-01	2003-11-03	2003-11-05	2003-11-07	2003-11-20	2003-11-21
##	1	2	2	1	1	2	1
##	2003-11-25	2003-11-26	2003-11-27	2003-11-28	2003-12-01	2003-12-02	2003-12-05
##	1	1	3	3	3	1	1
##	2003-12-08	2003-12-09	2003-12-11	2003-12-12	2003-12-15	2003-12-19	2003-12-22

##	1	1	1	1	2	1	1
##	2003-12-23	2003-12-26	2004-01-01	2004-01-05	2004-01-06	2004-01-07	2004-01-09
##	2	1	1	1	1	1	4
##	2004-01-11	2004-01-12	2004-01-13	2004-01-14	2004-01-15	2004-01-16	2004-01-21
##	1	1	1	1	1	2	1
##	2004-01-29	2004-01-30	2004-02-01	2004-02-02	2004-02-03	2004-02-04	2004-02-06
##	2	1	3	1	1	3	1
##	2004-02-09	2004-02-13	2004-02-15	2004-02-17	2004-02-19	2004-02-20	2004-02-24
##	2	2	2	1	1	1	2
##	2004-02-25	2004-03-01	2004-03-03	2004-03-05	2004-03-07	2004-03-08	2004-03-09
##	1	2	1	1	1	1	2
##	2004-03-10	2004-03-11	2004-03-12	2004-03-15	2004-03-16	2004-03-18	2004-03-21
##	1	2	1	2	2	1	1
##	2004-03-23	2004-03-26	2004-04-01	2004-04-02	2004-04-05	2004-04-06	2004-04-07
##	1	2	2	1	1	1	2
##	2004-04-13	2004-04-14	2004-04-15	2004-04-20	2004-04-21	2004-04-22	2004-04-23
##	1	2	3	1	1	1	1
##	2004-04-27	2004-04-29	2004-04-30	2004-05-04	2004-05-05	2004-05-06	2004-05-07
##	1	1	1	1	1	1	2
##	2004-05-11	2004-05-13	2004-05-14	2004-05-17	2004-05-18	2004-05-19	2004-05-21
##	1	1	2	1	1	1	1
##	2004-05-24	2004-05-26	2004-05-31	2004-06-01	2004-06-07	2004-06-08	2004-06-09
##	2	1	2	1	2	1	2
##	2004-06-11	2004-06-15	2004-06-22	2004-06-24	2004-06-25	2004-07-01	2004-07-05
##	1	4	1	3	1	2	1
##	2004-07-06	2004-07-12	2004-07-13	2004-07-14	2004-07-15	2004-07-16	2004-07-17
##	2	1	3	4	2	1	1
##	2004-07-19	2004-07-21	2004-07-22	2004-07-23	2004-07-24	2004-07-26	2004-07-27
##	1	1	1	1	1	1	3
##	2004-07-28	2004-07-29	2004-08-01	2004-08-02	2004-08-03	2004-08-04	2004-08-05
##	1	1	1	1	1	1	1
##	2004-08-06	2004-08-07	2004-08-10	2004-08-11	2004-08-12	2004-08-13	2004-08-14
##	1	1	1	1	1	2	3
##	2004-08-15	2004-08-16	2004-08-17	2004-08-18	2004-08-19	2004-08-20	2004-08-21
##	1	1	1	1	2	2	1
##	2004-08-26	2004-08-31	2004-09-01	2004-09-02	2004-09-05	2004-09-08	2004-09-13
##	2	1	4	1	1	2	3
##	2004-09-14	2004-09-15	2004-09-16	2004-09-17	2004-09-18	2004-09-20	2004-09-21
##	3	4	1	2	1	4	2
##	2004-09-23	2004-09-28	2004-10-03	2004-10-04	2004-10-05	2004-10-06	2004-10-07
##	1	1	1	3	3	1	2
##	2004-10-08	2004-10-12	2004-10-13	2004-10-14	2004-10-15	2004-10-17	2004-10-20
##	1	3	2	1	2	1	1
##	2004-10-21	2004-10-22	2004-10-23	2004-10-25	2004-10-27	2004-10-28	2004-10-29
##	2	2	2	1	2	1	2
##	2004-11-01	2004-11-08	2004-11-09	2004-11-10	2004-11-11	2004-11-12	2004-11-13
##	2	2	5	3	1	3	1
##	2004-11-15	2004-11-16	2004-11-18	2004-11-22	2004-11-23	2004-11-25	2004-11-26
##	4	3	1	1	4	1	1
##	2004-12-01	2004-12-03	2004-12-06	2004-12-07	2004-12-08	2004-12-09	2004-12-10
##	5	4	2	1	1	1	2
##	2004-12-12	2004-12-13	2004-12-15	2004-12-16	2004-12-23	2004-12-27	2004-12-28
##	1	1	1	1	3	1	1
##	2004-12-29	2004-12-30	2005-01-04	2005-01-05	2005-01-07	2005-01-10	2005-01-11

##	1	1	2	2	2	1	1
##	2005-01-12	2005-01-13	2005-01-15	2005-01-17	2005-01-26	2005-01-27	2005-01-28
##	2	3	1	2	3	2	2
##	2005-01-31	2005-02-01	2005-02-02	2005-02-04	2005-02-10	2005-02-11	2005-02-12
##	1	4	2	1	1	2	1
##	2005-02-14	2005-02-15	2005-02-16	2005-02-17	2005-02-18	2005-02-22	2005-02-24
##	4	2	2	3	2	2	1
##	2005-02-25	2005-02-28	2005-03-01	2005-03-02	2005-03-07	2005-03-10	2005-03-14
##	2	1	1	1	1	3	1
##	2005-03-18	2005-03-19	2005-03-22	2005-03-23	2005-03-24	2005-03-28	2005-03-29
##	1	1	1	4	1	1	1
##	2005-03-30	2005-04-01	2005-04-06	2005-04-07	2005-04-08	2005-04-09	2005-04-15
##	2	1	1	1	1	1	2
##	2005-04-22	2005-04-27	2005-04-28	2005-04-29	2005-04-30	2005-05-01	2005-05-03
##	1	1	1	4	5	1	3
##	2005-05-04	2005-05-05	2005-05-07	2005-05-08	2005-05-10	2005-05-11	2005-05-13
##	3	2	1	1	1	1	1
##	2005-05-14	2005-05-17	2005-05-18	2005-05-20	2005-05-24	2005-05-26	2005-05-27
##	1	1	1	2	1	1	3
##	2005-05-31	2005-06-01	2005-06-03	2005-06-07	2005-06-09	2005-06-10	2005-06-14
##	1	2	1	2	1	2	1
##	2005-06-15	2005-06-16	2005-06-17	2005-06-21	2005-06-22	2005-07-01	2005-07-02
##	1	1	2	2	1	2	1
##	2005-07-05	2005-07-06	2005-07-08	2005-07-11	2005-07-12	2005-07-14	2005-07-15
##	1	1	1	2	1	1	1
##	2005-07-16	2005-07-21	2005-07-29	2005-08-01	2005-08-02	2005-08-04	2005-08-05
##	1	1	1	3	1	1	3
##	2005-08-07	2005-08-08	2005-08-12	2005-08-16	2005-08-17	2005-08-26	2005-08-30
##	1	2	1	1	1	1	1
##	2005-09-01	2005-09-09	2005-09-20	2005-09-22	2005-09-23	2005-09-27	2005-10-01
##	1	1	1	1	2	1	1
##	2005-10-02	2005-10-05	2005-10-10	2005-10-11	2005-10-13	2005-10-14	2005-10-19
##	1	1	1	2	2	1	1
##	2005-10-26	2005-10-27	2005-10-28	2005-11-01	2005-11-03	2005-11-04	2005-11-06
##	1	2	1	1	1	1	1
##	2005-11-10	2005-11-13	2005-11-17	2005-11-22	2005-11-29	2005-12-03	2005-12-05
##	2	1	1	1	1	1	1
##	2005-12-06	2005-12-07	2005-12-08	2005-12-10	2005-12-12	2005-12-15	2005-12-28
##	2	2	1	1	1	1	1
##	2005-12-29	2005-12-30	2006-01-04	2006-01-12	2006-01-13	2006-01-19	2006-01-22
##	1	1	1	1	1	1	1
##	2006-01-23	2006-01-25	2006-01-26	2006-01-27	2006-01-31	2006-02-01	2006-02-07
##	1	1	1	2	1	1	1
##	2006-02-09	2006-02-17	2006-02-19	2006-02-21	2006-02-22	2006-02-23	2006-02-24
##	1	1	1	2	1	1	1
##	2006-02-27	2006-02-28	2006-03-01	2006-03-02	2006-03-04	2006-03-09	2006-03-10
##	1	1	2	1	1	1	1
##	2006-03-21	2006-03-25	2006-03-27	2006-03-28	2006-03-29	2006-03-30	2006-04-01
##	2	2	1	1	1	1	3
##	2006-04-06	2006-04-19	2006-04-21	2006-04-24	2006-05-04	2006-05-05	2006-05-11
##	1	1	1	1	1	1	1
##	2006-05-12	2006-05-13	2006-05-16	2006-05-19	2006-05-30	2006-05-31	2006-06-06
##	1	1	1	1	1	1	1
##	2006-06-08	2006-06-16	2006-07-06	2006-07-10	2006-07-14	2006-07-15	2006-07-20

```
##      1      2      1      2      2      1      1
## 2006-07-23 2006-07-27 2006-07-28 2006-08-01 2006-08-03 2006-08-15 2006-08-23
##      1      1      1      1      1      1      1
## 2006-08-26 2006-09-06 2006-09-08 2006-09-15 2006-09-18 2006-09-19 2006-09-25
##      1      1      2      1      1      2      1
## 2006-10-09 2006-10-10 2006-10-12 2006-10-14 2006-10-17 2006-10-18 2006-10-21
##      1      1      1      1      1      1      2
## 2006-10-26 2006-10-27 2006-10-29 2006-10-31 2006-11-07 2006-11-11 2006-11-14
##      1      1      1      1      1      1      1
## 2006-11-15 2006-11-21 2006-12-07 2006-12-11 2006-12-16 2006-12-22 2006-12-26
##      2      1      1      2      1      1      1
## 2007-01-03 2007-01-12 2007-01-13 2007-01-17 2007-01-23 2007-01-24 2007-01-31
##      1      2      3      1      1      1      1
## 2007-02-05 2007-02-13 2007-02-23 2007-03-02 2007-03-04 2007-03-05 2007-03-06
##      1      2      1      1      1      1      1
## 2007-03-09 2007-03-13 2007-03-16 2007-03-20 2007-03-23 2007-03-30 2007-04-17
##      1      1      1      1      1      1      1
## 2007-04-21 2007-04-24 2007-04-30 2007-05-02 2007-06-05 2007-06-29 2007-07-17
##      1      1      1      2      1      1      1
## 2007-07-30 2007-08-01 2007-08-13 2007-09-05 2007-09-18 2007-09-24 2007-09-26
##      1      2      1      1      1      1      1
## 2007-10-04 2007-10-29
##      2      1
```

2. Create an indicator variable (one which takes the values 0 or 1 only) to represent death within 1 year of the initial visit. How many observations died in year 1?

```
hold <- difftime(df$date.death, df$init.date, units="days")
hold.tf <- hold < 365
hold.tf[is.na(hold.tf)] <- 0
df[, 'deathly'] <- hold.tf
table(hold.tf)[2]
```

```
## 1
## 92
```

3. Use the `init.date`, `last.visit` and `death.date` columns to calculate a followup time (in days), which is the difference between the first and either the last visit or a death event (whichever comes first). If these times are longer than 1 year, censor them (this means if the value is above 365, set followup to 365). Print the quantile for this new variable.

```
for(i in 1:nrow(df)) {
  df[i, 'followup'] = min(365, difftime(min(c(df[i, 'last.visit'], df[i, 'date.death']), na.rm=T), df[i, 'init.date']))
}
quantile(df$followup, na.rm=TRUE)
```

```
##      0%      25%      50%      75%     100%
## 0.0000 320.7188 365.0000 365.0000 365.0000
```

4. Create another indicator variable representing loss to followup; this means the observation is not known to be dead but does not have any followup visits after the first year. How many records are lost-to-followup?

```
for(i in 1:nrow(df)) {
  if (df[i, 'followup'] < 365 & is.null(df[i, 'death.date'])){
    df[i, 'lost-to-followup'] = 1
  } else {
```

```

    df[i, 'lost-to-followup'] = 0
  }
}

```

5. Recall our work in class, which separated the `init.reg` field into a set of indicator variables, one for each unique drug. Create these fields and append them to the database as new columns. Which drug regimen are found over 100 times? 18 new columns are created, one for each drug. From the different regimens, "3TC, AZT, EFV" and "3TC, AZT, NVP" both had over 100 times.

```

over100 <- table(df$init.reg)
over100[over100 > 100]

```

```

##
## 3TC,AZT,EFV 3TC,AZT,NVP
##          421      284

```

```

for(i in 1:nrow(df)) {
  hold = strsplit(df[i, 'init.reg'], ',')
  for (j in hold){
    df[i, j] = 1
  }
}

```

6. The dataset `haart2.csv` contains a few additional observations for the same study. Import these and append them to your master dataset (if you were smart about how you coded the previous steps, cleaning the additional observations should be easy!). Show the first five records and the last five records of the complete (and clean) data set.

```

library(plyr)
df2 <- read.csv('haart2.csv', header=TRUE, stringsAsFactors=FALSE)
df2$init.date <- as.POSIXct(df2$init.date, format = "%m/%d/%y")
df2$last.visit <- as.POSIXct(df2$last.visit, format = "%m/%d/%y")
df2$date.death <- as.POSIXct(df2$date.death, format = "%m/%d/%y")
hold <- difftime(df2$date.death, df2$init.date, units="days")
hold.tf <- hold < 365
hold.tf[is.na(hold.tf)] <- 0
df2[, 'death1y'] <- hold.tf
for(i in 1:nrow(df2)) {
  df2[i, 'followup'] = min(365, difftime(min(c(df2[i, 'last.visit'], df2[i, 'date.death']), na.rm=T), df2$init.date, units="days"))
}
for(i in 1:nrow(df2)) {
  if (df2[i, 'followup'] < 365 & is.null(df2[i, 'death.date'])){
    df2[i, 'lost-to-followup'] = 1
  } else {
    df2[i, 'lost-to-followup'] = 0
  }
}
for(i in 1:nrow(df2)) {
  hold = strsplit(df2[i, 'init.reg'], ',')
  for (j in hold){
    df2[i, j] = 1
  }
}
df3 <- rbind.fill(df, df2)
head(df3, 5)

```



```
##      male age aids cd4baseline logvl weight hemoglobin init.reg init.date
## 1      1  25   0         NA      NA      NA      NA 3TC,AZT,EFV 2003-07-01
## 2      1  49   0        143      NA 58.0608      11 3TC,AZT,EFV 2004-11-23
## 3      1  42   1        102      NA 48.0816       1 3TC,AZT,EFV 2003-04-30
## 4      0  33   0        107      NA 46.0000      NA 3TC,AZT,NVP 2006-03-25
## 5      1  27   0         52       4      NA      NA 3TC,D4T,EFV 2004-09-01
##      last.visit death date.death deathly followup lost-to-followup 3TC AZT EFV
## 1 2007-02-26      0      <NA>      0 365.00000      0 1 1 1
## 2 2008-02-22      0      <NA>      0 365.00000      0 1 1 1
## 3 2005-11-21      1 2006-01-11      0 365.00000      0 1 1 1
## 4 2006-05-05      1 2006-05-07      1 40.95833      1 1 1 NA
## 5 2007-11-13      0      <NA>      0 365.00000      0 1 NA 1
##      NVP D4T ABC DDI IDV LPV RTV SQV FTC TDF DDC NFV T20 ATV FPV
## 1      NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 2      NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 3      NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 4      1  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 5      NA  1  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
```

```
tail(df3, 5)
```

```
##      male      age aids cd4baseline      logvl weight hemoglobin init.reg
## 1000      0 40.00000      1      131      NA 46.2672       8 3TC,D4T,NVP
## 1001      0 27.00000      0      232      NA      NA      NA 3TC,AZT,NVP
## 1002      1 38.72142      0      170      NA 84.0000      NA 3TC,AZT,NVP
## 1003      1 23.00000      NA      154 3.995635 65.5000      14 3TC,DDI,EFV
## 1004      0 31.00000      0      236      NA 45.8136      NA 3TC,D4T,NVP
##      init.date last.visit death date.death deathly followup lost-to-followup
## 1000 2003-07-03 2008-02-29      0      <NA>      0 365.00000      0
## 1001 2003-12-01 2004-01-05      0      <NA>      0 35.00000      1
## 1002 2002-09-26 2004-03-29      0      <NA>      0 365.00000      0
## 1003 2007-01-31 2007-04-16      0      <NA>      0 74.95833      1
## 1004 2003-12-03 2007-10-11      0      <NA>      0 365.00000      0
##      3TC AZT EFV NVP D4T ABC DDI IDV LPV RTV SQV FTC TDF DDC NFV T20 ATV FPV
## 1000      1  NA  NA      1  1  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 1001      1   1  NA      1  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 1002      1   1  NA      1  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 1003      1  NA   1  NA  NA  NA      1  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
## 1004      1  NA  NA      1  1  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA  NA
```