# Reading in the Active Duty Marital Data

Sean Conway

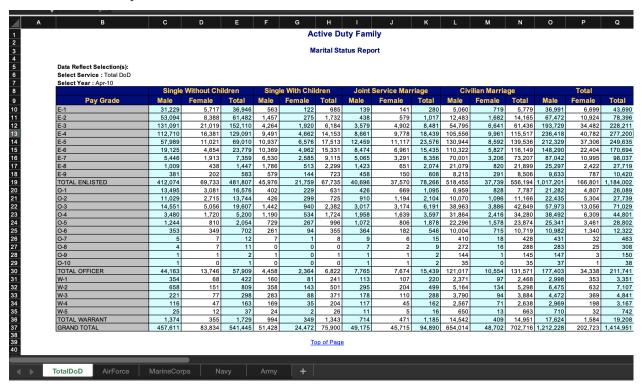
8/19/2021

### Marital DOD Data

The data come from data.gov. The file contains count data on the marital statuses of military service members, as well as their pay grade/family status. The data file is called ActiveDuty\_MaritalStatus.xls.

### BUT...

These data are not optimized for R.



So we will need to do some work to read in the data.

## Reading in the data - First Excel Sheet

While eventually we will want to read in all the sheets at once, we will start out by reading in the first Excel Sheet, TotalDoD (see above).

We will use the read\_excel() function from the readxl package.

First, we're going to manually specify the names of our columns. This involves doing a bit hard coding (reading in messy data is the only time hard coding is recommended!).

Note that we named these columns so they can be adequately separated later on.

Next we have to use read\_excel(), but we have to specify a number of argumnets. We first specify the path (note - file\_path is a variable that I created ahead of time to be specific to my computer. Yours will be different). Next, we specify sheet, the number of the sheet we wish to read in (we can also specify the sheet name). Next, we carefully chooose the range of cells in the file we read in, based on our visual inspection of the file. This is another case where we must hard-code it. Here, we want the range to go from B10 to Q37. We also need to manually specify col\_names from our col\_names\_dod vector that we created above. Once we read it in, we will remove any rows containing the word "total" from the column pay\_grade

```
file_path
```

## [1] "/Users/seanconway/Github/Datasets/R/data/ActiveDuty\_MaritalStatus.xls"

```
## # A tibble: 24 x 16
##
      pay_grade single_withoutc~ single_withoutc~ single_withoutc~ single_withchil~
##
      <chr>
                            <dbl>
                                              <dbl>
                                                               <dbl>
                                                                                 <dbl>
##
   1 E-1
                            31229
                                              5717
                                                               36946
                                                                                   563
  2 E-2
                                              8388
##
                            53094
                                                               61482
                                                                                  1457
## 3 E-3
                           131091
                                              21019
                                                              152110
                                                                                  4264
##
  4 E-4
                           112710
                                                              129091
                                                                                  9491
                                              16381
  5 E-5
                                                                                 10937
##
                            57989
                                              11021
                                                               69010
  6 E-6
##
                                               4654
                                                               23779
                                                                                 10369
                            19125
   7 E-7
##
                             5446
                                               1913
                                                                7359
                                                                                  6530
## 8 E-8
                             1009
                                                438
                                                                1447
                                                                                  1786
## 9 E-9
                              381
                                                202
                                                                                   579
                                                                 583
## 10 0-1
                            13495
                                               3081
                                                               16576
                                                                                   402
## # ... with 14 more rows, and 11 more variables:
## #
       single_withchildren_female <dbl>, single_withchildren_total <dbl>,
## #
       married_jointservice_male <dbl>, married_jointservice_female <dbl>,
       married_jointservice_total <dbl>, married_civilian_female <dbl>,
## #
```

```
## # married_civilian_male <dbl>, married_civilian_total <dbl>,
## # married_male_total <dbl>, married_female_total <dbl>,
## # married_total total <dbl>
```

We've read in the data!

This tibble looks okay, but there's still much work to be done. First, we need to remove any of the columns that contain the word "total". We don't need these aggregated totals, as they will only muddle the data (plus we can calculate them ourselves if needed).

```
marital_dod_2 <- marital_dod_1 %>%
  select(c(pay_grade,!contains("total")))
marital_dod_2
```

```
## # A tibble: 24 x 9
##
      pay_grade single_withoutc~ single_withoutc~ single_withchil~ single_withchil~
##
      <chr>
                             <dbl>
                                               <dbl>
                                                                 <dbl>
                                                                                   <dbl>
##
   1 E-1
                            31229
                                                5717
                                                                   563
                                                                                     122
##
    2 E-2
                            53094
                                                8388
                                                                  1457
                                                                                     275
##
   3 E-3
                                               21019
                                                                  4264
                                                                                    1920
                           131091
##
   4 E-4
                                                                                    4662
                           112710
                                               16381
                                                                  9491
##
   5 E-5
                            57989
                                               11021
                                                                                    6576
                                                                 10937
    6 E-6
##
                            19125
                                                4654
                                                                 10369
                                                                                    4962
   7 E-7
##
                                                                                    2585
                             5446
                                                1913
                                                                  6530
##
   8 E-8
                              1009
                                                 438
                                                                  1786
                                                                                     513
## 9 E-9
                                                                   579
                               381
                                                 202
                                                                                     144
## 10 0-1
                            13495
                                                3081
                                                                   402
                                                                                     229
## # ... with 14 more rows, and 4 more variables: married_jointservice_male <dbl>,
       married_jointservice_female <dbl>, married_civilian_female <dbl>,
## #
       married_civilian_male <dbl>
```

Next, we'll use pivot\_longer() to combine the column names (except for pay\_grade) into a single column, status. We do this because the variable status is currently spread across columns (it is wide). We want our data to be tidy - where each row is a single observation.

We specify cols as every column except pay\_grade with !contains(pay\_grade). We also specify that the column names will be moved to status and the column values will be moved to count.

```
## # A tibble: 192 x 3
##
      pay_grade status
                                                count
##
                <chr>>
      <chr>>
                                                <dbl>
##
   1 E-1
                single_withoutchildren_male
                                                31229
    2 E-1
##
                single_withoutchildren_female
                                                 5717
##
   3 E-1
                single_withchildren_male
                                                  563
##
  4 E-1
                single withchildren female
                                                  122
  5 E-1
##
                married_jointservice_male
                                                  139
## 6 E-1
                married jointservice female
                                                  141
## 7 E-1
                married_civilian_female
                                                 5060
```

```
## 8 E-1 married_civilian_male 719
## 9 E-2 single_withoutchildren_male 53094
## 10 E-2 single_withoutchildren_female 8388
## # ... with 182 more rows
```

Our data is now tidy! Now we have just a bit more to do. We should use separate() to separate pay\_grade into two columns (enlisted and pay\_grade), as well as separate status into three columns (relationship, family\_status, and gender).

```
## # A tibble: 192 x 6
##
      enlisted pay grade relationship family status
                                                        gender count
##
      <chr>
               <chr>
                          <chr>
                                        <chr>
                                                        <chr>
                                                               <dbl>
##
   1 E
               1
                          single
                                       withoutchildren male
                                                                31229
    2 E
##
                          single
                                       withoutchildren female
                                                               5717
               1
##
   3 E
                          single
                                       withchildren
                                                        male
                                                                  563
               1
  4 E
                                       withchildren
##
               1
                          single
                                                        female
                                                                  122
##
  5 E
                          married
                                        jointservice
                                                        male
                                                                  139
## 6 E
               1
                                        jointservice
                                                        female
                                                                  141
                          married
   7 E
##
               1
                          married
                                        civilian
                                                        female
                                                                 5060
## 8 E
               1
                                        civilian
                          married
                                                        male
                                                                  719
## 9 E
                          single
                                       withoutchildren male
                                                                53094
## 10 E
               2
                          single
                                       withoutchildren female 8388
## # ... with 182 more rows
```

Whala! We successfully read in and cleaned a very messy Excel spreadsheet. Next, we will use the purrr package to read multiple sheets at once.

# Reading in the Data - All Sheets

As we saw, there are multiple sheets in ActiveDuty\_MaritalStatus.xls. We will need to do a bit of work to read them all in, while maintaining best practices (i.e., NOT copying and pasting code).

First, we define the column names we want. Note that this is the same as above, but we need to include am empty string, since we will be forced to read in the "first" column (which is empty).

```
"married_civilian_male",
"married_civilian_total",
"married_male_total",
"married_female_total",
"married_total_total")
```

Next, we use the excel\_sheets() function to extract a character vector of all sheet names from our file

```
sheets <- excel_sheets(file_path)
sheets</pre>
```

"MarineCorps" "Navy"

"Army"

Next, is the bulk of our workload. We write a custom function that will read in each sheet, as well as clean and tidy it. The function is called read\_mar\_sheets(), and it takes a single argument: sheet\_name.

## [1] "TotalDoD"

"AirForce"

The function first reads in the data, using read\_excel(). Note that we specify the pathfrom the file\_path variable I created above. The sheet is specified by the user as sheet\_name. We also ask read\_excel() to trim whitespace (trim\_ws=TRUE). The col\_names are defined with col\_names\_all from above, and skip=9 tells read\_excel() to skip the first 9 rows (this number comes from a visual inspection of the Excel file).

After the sheet is read in, we use the pipe (%>%) to continue perfoming operations on it. We create a column branch, equal to the value of sheet\_name. Then, we use select(2:last\_col()) to remove the first column from the data (recall that the first column is blank). We then use select(c(pay\_grade, !contains("total"))) to remove any columns with the word "total" in their name.

filter(str\_detect(pay\_grade, regex("total",ignore\_case = TRUE),negate = T)) removes the
word "total" from the column pay\_grade.

Just like above, we use pivot\_longer() to make the data tidy, and then use separate() to separate out pay\_grade and status into different columns.

This whole operation assigns a tibble to data (in the function scope), and a return statement returns this tibble to the user.

```
read_mar_sheets <- function(sheet_name){</pre>
  data <- read_excel(path = `file_path`,</pre>
             sheet = sheet_name,
             trim_ws = TRUE,
             col_names = col_names_all,
             skip=9) %>%
    mutate("branch"=sheet name) %>%
    select(2:last_col()) %>% # immediately remove out blank column
    select(c(pay grade,!contains("total"))) %>% # remove columns with the word total
    filter(str_detect(pay_grade, regex("total",ignore_case = TRUE),negate = T)) %>%
    pivot_longer(cols = !contains(c("pay_grade", "branch")), # these columns can remain as is
               names_to = "status", values_to = "count") %>%
    separate(col=pay_grade, into=c("enlisted","pay_grade"),
             sep="-") %>%
    separate(col=status, into = c("relationship", "family_status", "gender"),
           sep = "_")
  return(data)
}
```

Our function is obviously customized for this specific operation. It also contains some hard coding, which is necessary given the data file. However, note that the sheets are actually fairly consistent (e.g., the same

amount of white space, very similar row values, etc.). In the future, you may encounter even trickier Excel files, which will require creative programming to read into R.

We can efficiently use this function with purrr::map() to read in all our sheets simulataneously. purrr is package that's part of the tidyverse, and it contains functions (like map() that allow us to vectorize operations.

Here is an example, using a simple custom function called add\_1, which adds 1 to its input.

```
x \leftarrow c(1,2,3,4,5)
add_1 <- function(x) x+1
purrr::map(x, add_1)
## [[1]]
## [1] 2
##
## [[2]]
## [1] 3
##
## [[3]]
## [1] 4
## [[4]]
## [1] 5
##
## [[5]]
## [1] 6
```

Going back to vector sheets, we can use purrr:map() to vectorize our read-in operation. This wil create a list of tibbles, which can then be easily bound together into one large tibble using bind\_rows().

```
data_list <- purrr::map(sheets, read_mar_sheets)
data_list</pre>
```

```
## [[1]]
## # A tibble: 192 x 7
##
      enlisted pay_grade branch
                                   relationship family_status
                                                                 gender count
               <chr>
##
      <chr>
                          <chr>
                                   <chr>
                                                 <chr>
                                                                  <chr>>
                                                                         <dbl>
   1 E
                                                                         31229
##
               1
                          TotalDoD single
                                                 withoutchildren male
##
    2 E
               1
                         TotalDoD single
                                                withoutchildren female 5717
##
   3 E
               1
                         TotalDoD single
                                                withchildren
                                                               \mathtt{male}
                                                                           563
##
   4 E
                         TotalDoD single
                                                withchildren
                                                                 female
                                                                           122
               1
    5 E
##
               1
                         TotalDoD married
                                                 jointservice
                                                                 male
                                                                           139
##
    6 E
               1
                         TotalDoD married
                                                                 female
                                                                           141
                                                 jointservice
##
   7 E
               1
                          TotalDoD married
                                                 civilian
                                                                 female
                                                                          5060
##
   8 E
                                                                           719
               1
                          TotalDoD married
                                                 civilian
                                                                 male
##
   9 E
                         TotalDoD single
                                                 withoutchildren male
                                                                         53094
## 10 E
               2
                          TotalDoD single
                                                 withoutchildren female 8388
## # ... with 182 more rows
##
## [[2]]
## # A tibble: 152 x 7
##
      enlisted pay_grade branch
                                   relationship family_status
                                                                 gender count
##
      <chr>
                          <chr>
                                   <chr>>
                                                 <chr>
                                                                 <chr> <dbl>
               <chr>
```

```
1 E
##
                          AirForce single
                                                  withoutchildren male
                                                                           7721
                1
##
    2 E
                1
                          AirForce single
                                                  withoutchildren female
                                                                           1550
   3 E
                          AirForce single
                                                  withchildren
##
                                                                  male
                                                                             27
   4 E
                                                  withchildren
                                                                   female
                                                                              5
##
                          AirForce single
                1
##
    5 E
                          AirForce married
                                                  jointservice
                                                                  male
                                                                             49
##
    6 E
                1
                          AirForce married
                                                                   female
                                                                             27
                                                  jointservice
##
   7 E
                1
                          AirForce married
                                                  civilian
                                                                   female
                                                                           1064
##
   8 F.
                          AirForce married
                                                  civilian
                                                                  male
                                                                            178
                1
## 9 E
                2
                          AirForce single
                                                 withoutchildren male
                                                                           4380
## 10 E
                2
                                                 withoutchildren female 1010
                          AirForce single
   # ... with 142 more rows
##
## [[3]]
## # A tibble: 192 x 7
      enlisted pay_grade branch
                                       relationship family_status
                                                                      gender count
##
      <chr>
                <chr>>
                          <chr>
                                       <chr>
                                                     <chr>
                                                                      <chr>
                                                                             <dbl>
##
    1 E
                          MarineCorps single
                                                     withoutchildren male
                                                                              6232
                1
   2 E
##
                1
                          MarineCorps single
                                                     withoutchildren female
                                                                               583
##
   3 E
                          MarineCorps single
                                                     withchildren
                                                                      male
                                                                                54
                1
  4 E
##
                          MarineCorps single
                                                     withchildren
                                                                      female
                                                                                 3
##
    5 E
                1
                          MarineCorps married
                                                     jointservice
                                                                      male
                                                                                20
##
   6 E
                          MarineCorps married
                                                     jointservice
                                                                      female
                                                                                19
   7 F.
##
                          MarineCorps married
                                                     civilian
                                                                      female
                                                                               611
                1
##
   8 E
                1
                          MarineCorps married
                                                     civilian
                                                                      male
                                                                                21
## 9 E
                2
                          MarineCorps single
                                                     withoutchildren male
                                                                             15916
## 10 E
                2
                          MarineCorps single
                                                     withoutchildren female 1336
## # ... with 182 more rows
## [[4]]
## # A tibble: 184 x 7
##
      enlisted pay_grade branch relationship family_status
                                                                gender count
##
      <chr>
                <chr>
                          <chr>
                                 <chr>>
                                                <chr>
                                                                 <chr>
                                                                        <dbl>
##
   1 E
                                                                         7820
                1
                          Navy
                                  single
                                               withoutchildren male
##
    2 E
                          Navy
                                               withoutchildren female
                                                                         2275
                1
                                  single
    3 E
##
                1
                          Navv
                                  single
                                               withchildren
                                                                male
                                                                          117
    4 E
                          Navy
                                               withchildren
                                                                female
                                                                           34
                1
                                  single
##
   5 E
                          Navy
                                  married
                                               jointservice
                                                                male
                                                                           30
##
   6 E
                1
                          Navy
                                 married
                                               jointservice
                                                                female
                                                                           57
    7 E
##
                1
                          Navy
                                  married
                                               civilian
                                                                female
                                                                          806
##
  8 E
                1
                                               civilian
                                                                male
                                                                          162
                          Navy
                                  married
##
  9 E
                                               withoutchildren male
                                                                        11198
                          Navy
                                  single
## 10 E
                2
                          Navy
                                  single
                                               withoutchildren female 2718
## # ... with 174 more rows
##
## [[5]]
## # A tibble: 192 x 7
                                                                gender count
##
      enlisted pay_grade branch relationship family_status
##
      <chr>
                                  <chr>>
                                                <chr>
                                                                        <dbl>
                <chr>>
                          <chr>
                                                                 <chr>
##
   1 E
                1
                          Army
                                  single
                                               withoutchildren male
                                                                         9456
## 2 E
                1
                          Army
                                  single
                                               withoutchildren female
                                                                         1309
##
   3 E
                                               withchildren
                                                                male
                                                                          365
                1
                          Army
                                  single
  4 E
##
                1
                          Army
                                  single
                                               withchildren
                                                                female
                                                                           80
                                 married
## 5 E
                1
                          Army
                                                jointservice
                                                                male
                                                                           40
## 6 E
                1
                          Army
                                 married
                                               jointservice
                                                                female
```

```
7 E
               1
                         Army
                                married
                                              civilian
                                                               female
                                                                       2579
   8 F.
##
               1
                                married
                                              civilian
                                                              male
                                                                        358
                         Army
                                                                      21600
## 9 E
               2
                         Army
                                 single
                                              withoutchildren male
                                              withoutchildren female 3324
## 10 E
               2
                                 single
                         Army
## # ... with 182 more rows
```

```
marital_tidy_all <- data_list %>%
  bind_rows()
marital_tidy_all
```

```
## # A tibble: 912 x 7
##
      enlisted pay_grade branch
                                   relationship family_status
                                                                 gender count
##
               <chr>
                                   <chr>
                                                 <chr>
      <chr>
                          <chr>
                                                                 <chr>
                                                                         <dbl>
##
    1 E
               1
                          TotalDoD single
                                                 withoutchildren male
                                                                         31229
    2 E
                          TotalDoD single
                                                 withoutchildren female
                                                                          5717
##
               1
    3 E
##
               1
                         TotalDoD single
                                                 withchildren
                                                                 male
                                                                           563
##
    4 E
                         TotalDoD single
                                                withchildren
                                                                 female
                                                                           122
               1
##
   5 E
               1
                          TotalDoD married
                                                 jointservice
                                                                 male
                                                                           139
   6 E
                          TotalDoD married
##
               1
                                                 jointservice
                                                                 female
                                                                           141
                                                                 female
##
    7 E
               1
                          TotalDoD married
                                                 civilian
                                                                          5060
##
   8 E
                         TotalDoD married
               1
                                                 civilian
                                                                 male
                                                                           719
## 9 E
               2
                         TotalDoD single
                                                withoutchildren male
                                                                         53094
               2
## 10 E
                          TotalDoD single
                                                 withoutchildren female 8388
## # ... with 902 more rows
```

# Conclusion

You now know more about reading in messy data files to R. This is by no means the only solution to this problem. You may have even found a way that's more efficient!