

1 Wrangling categorical data in R

2 Amelia McNamara¹ and Nicholas J Horton²

3 ¹Statistical and Data Sciences, Smith College

4 ²Mathematics and Statistics, Amherst College

5 ABSTRACT

6 Data wrangling is a critical foundation of data science. Wrangling of categorical data is an important
7 component of the analysis cycle. Aspects of these operations can sometimes be tricky, particularly for
8 complex transformations that arise in real-world settings. This paper discusses aspects of categorical
9 variable transformations in R. We consider several motivating examples, suggest defensive coding
10 strategies, and outline principles for data wrangling to help ensure data quality and sound analysis.

11 Keywords:

12 INTRODUCTION

13 Wrangling skills provide an intellectual and practical foundation for data science. Because of the com-
14 plexity of some transformations, careless data derivation operations can lead to errors or inconsistencies.
15 The wrangling of categorical data presents particular challenges and is highly relevant because so many
16 variables are categorical (e.g., gender, income bracket, and state). Data that are ingested into R from
17 spreadsheets can lead to many problems.

18 In this paper, we consider a number of common idioms related to categorical data that often arise in
19 data cleaning and preparation, propose some guidelines for defensive coding, and discuss some settings
20 where analysts often get tripped up when working with categorical variables and factors (R's data type for
21 categorical data).

22 To ground our work, we will compare and contrast how categorical data are treated in **base R** versus
23 the so-called tidyverse (?). Tools from the tidyverse, discussed in another paper in this special issue (see
24 <https://github.com/dsscollection/tidyflow>), are designed to make analysis purer, more predictable, and
25 pipeable. They help facilitate a reproducible workflow where a new version of the data could be supplied
26 in the code with updated results produced (?). Wrangling of categorical data can make this task even
27 more complex (e.g., if a new level of a categorical variable is added in an updated dataset or inadvertently
28 introduced by a careless error in a spreadsheet to be ingested into R).

29 THE IMPORTANCE OF TOOLING

30 It's important that statistical and data science tools foster good practice and provide a robust environment
31 for data wrangling and data management. In this section we will how factors XX AM define? are used in
32 base R and enumerate problems that arise with their use.

33 FACTORS IN R

34 Consider a variable describing gender that includes categories `male`, `female` and `non-conforming`.
35 In R, there are two ways to store this information. One is to use a series of character strings, and the other
36 is to store it as a factor.

37 Historically, storing categorical data as a factor variable was more efficient than storing the same data
38 as strings, because factor variables only store the factor labels once (?). However, R uses hashed versions
39 of all character strings, so the storage issue is no longer a consideration (?). For historical reasons, many
40 functions store variables by default as factors.

41 Factors can be very tricky to deal with, since many operations applied to them return different values
42 than when applied to character vectors. As an example, consider a set of decades,

```
x1 <- c(10, 10, 20, 20, 40)
x1f <- factor(x1)
ds <- data.frame(x1, x1f)
library(dplyr)
ds <- mutate(ds, x1recover = as.numeric(x1f))
ds

##   x1 x1f x1recover
## 1 10  10         1
## 2 10  10         1
## 3 20  20         2
## 4 20  20         2
## 5 40  40         3
```

43 Instead of creating a new variable with a numeric version of the value of the factor variable `x1f`
 44 the variable is created with a factor number (e.g., 10 is mapped to 1, 20 is mapped to 2). This result is
 45 unexpected and unfortunate because `as.numeric()` is intended to be recover numeric information in
 46 the **base R** paradigm when coercing a character variable. Compare the following:

```
as.numeric(c("hello"))

## [1] NA

as.numeric(factor(c("hello")))

## [1] 1
```

47 The unfortunate behavior of factors in R has led to an online movement against the default behavior of
 48 many data import functions to take any variable composed as strings and automatically convert the variable
 49 to a factor. The tidyverse moves away from this default behavior, with functions from the **readr** package de-
 50 faulting to leaving strings as-is. (Others have chosen to add `options(stringAsFactors=FALSE)`
 51 into their startup commands.)

52 Although the storage issues have been solved, and there are problems with defaulting strings to factors,
 53 factors are still necessary for some data analytic tasks. The most salient case is in modeling. When
 54 you pass a factor variable into `lm()` or `glm()`, R automatically creates indicator (or more pejoratively
 55 ‘dummy’) variables for each of the levels and picks one as a reference group. This behavior is lost if the
 56 variable is stored as a character vector. Factor variables also allow for the possibility of ordering between
 57 classes. Text strings `low`, `medium`, `high` would not preserve the ordering inherent in the groups.
 58 Again, this can be important for modeling when doing ordinal logistic regression and multinomial logistic
 59 regression.

60 While factors are important, they can often be hard to deal with. Because of the way the group
 61 numbers are stored separately from the factor labels, it can be easy to overwrite data in such a way that
 62 the original data are lost. In this paper, we will suggest best practices for working with factor data.

63 To motivate this process, we will consider data from the General Social Survey XX AM cite?. There
 64 are some import issues inherent to the data which are not particular to categorical data, so that processing
 65 is processing in Appendices ??, ??. We’ll work with the data that has cleaned variable names.

```
library(dplyr)
```

```
GSS <- read.csv("../data/GSScleaned.csv")
glimpse(GSS)

## Observations: 2,540
## Variables: 17
## $ Gss.year.for.this.respondent..... <dbl> 2014, 2014...
## $ Respondent.id.number <dbl> 1, 2, 3, 4...
## $ LaborStatus <fctr> Working f...
## $ Rs.occupational.prestige.score...1970. <dbl> 0, 0, 0, 0...
## $ Marital.status <fctr> Divorced,...
## $ Number.of.children <dbl> 0, 0, 1, 2...
## $ Age <fctr> 53.000000...
## $ Highest.year.of.school.completed <dbl> 16, 16, 13...
## $ Respondents.sex <fctr> Male, Fem...
## $ Race.of.respondent <fctr> White, Wh...
## $ Rs.family.income.when.16.yrs.old <fctr> Below ave...
## $ Total.family.income <fctr> $25000 or...
## $ Respondents.income <fctr> $25000 or...
## $ Total.family.income.1 <fctr> Not appli...
## $ PolParty <fctr> Not str r...
## $ Opinion.of.family.income <fctr> Above ave...
## $ Sexual.orientation <fctr> Heterosex...
```

66 The rest of this paper is organized around case studies related to particular tasks:

- 67 1. Changing the labels of factor levels,
- 68 2. Reordering factor levels,
- 69 3. Combining several levels into one (both string-like labels and numeric, probably go together), and
- 70 4. Making derived factor variables.

71 CHANGING THE LABELS OF FACTOR LEVELS

72 For this example, we will be considering the labor status variable. It has 9 factor levels. Most of the labels
73 are spelled out fully, but a few are strangely formatted. We want to change this.

74 One action you might want to take is just to change the text of one (or more) of the factor labels, so it
75 appears more nicely formatted in a **ggplot2** plot, for example.

76 There are two typical approaches in base R. One is more compact, but depends on the levels of the
77 factor not changing in the data being fed in, and the other is more robust, but extremely verbose. In
78 contrast, the **dplyr** package offers a method that is much more human readable, while also supporting
79 reproducibility.

80 Compact but fragile (base R)

```
levels(GSS$LaborStatus)

## [1] "Keeping house" "No answer" "Other"
## [4] "Retired" "School" "Temp not working"
## [7] "Unempl, laid off" "Working fulltime" "Working parttime"

summary(GSS$LaborStatus)

## Keeping house No answer Other Retired
## 263 2 76 460
## School Temp not working Unempl, laid off Working fulltime
## 90 40 104 1230
## Working parttime NA's
## 273 2
```

```
levels(GSS$LaborStatus) <- c(levels(GSS$LaborStatus)[1:5],
```

```

"Temporarily not working",
"Unemployed, laid off",
"Working full time",
"Working part time")
summary(GSS$LaborStatus)

##           Keeping house           No answer           Other
##           263                2                76
##           Retired           School Temporarily not working
##           460                90                40
## Unemployed, laid off Working full time Working part time
##           104                1230                273
##           NA's
##           2

```

81 This method is less than ideal, because it depends on the data coming in with the factor levels ordered
82 in a particular way. If the data gets changed outside of R, for example so that “Working full time” becomes
83 the first level, the code will silently fail with invalid results.

84 XX AM I’m not sure that this is true. See the following example:

```

x1 <- as.factor(c("Foo", "Bar", "Splat"))
x2 <- as.factor(c("Bar", "Splat", "Foo"))
levels(x1)

## [1] "Bar" "Foo" "Splat"

levels(x2)

## [1] "Bar" "Foo" "Splat"

```

85 The workflow will fail if additional factor levels are added after the fact. In our experience, both with
86 students and scientific collaborators, spreadsheet data can be easily changed in these ways (XX AM cite
87 Leek how to give data to a statistician?).

88 Robust but verbose (base R)

89 Another (more robust method) to recode this variable in **base R** is to use subsetting to overwrite particular
90 values in the data.

```
summary(GSS$PolParty)
```

```
##           Don't know           Ind,near dem           Ind,near rep
##           1              337              249
##           Independent           No answer           Not str democrat
##           502              25              406
## Not str republican           Other party           Strong democrat
##           292              62              419
## Strong republican           NA's
##           245              2

GSS$NewParty <- as.character(GSS$PolParty)
GSS$NewParty[GSS$PolParty=="Ind,near dem"] <- "Independent, near democrat"
GSS$NewParty[GSS$PolParty == "Ind,near rep"] <- "Independent, near republican"
GSS$NewParty[GSS$PolParty == "Not str democrat"] <- "Not strong democrat"
GSS$NewParty <- factor(GSS$NewParty)
summary(GSS$NewParty)

##           Don't know           Independent
##           1              502
## Independent, near democrat Independent, near republican
##           337              249
##           No answer           Not str republican
##           25              292
##           Not strong democrat           Other party
##           406              62
##           Strong democrat           Strong republican
##           419              245
##           NA's
##           2
```

91 This approach has a number of limitations in addition to being tedious and error prone. It is possible to
 92 miss cases and the potential for cut-and-paste errors to apply this same transformation to another variable.

93 Direct and robust (dplyr)

94 The `recode()` function in the **dplyr** package is a vector function,

```
GSS <- GSS %>%
  mutate(dplyrParty =
    recode(PolParty, `Not str republican` = "Not a strong republican",
                  `Ind,near dem` = "Independent, near democrat",
                  `Ind,near rep` = "Independent, near republican",
                  `Not str democrat` = "Not a strong democrat"))
summary(GSS$dplyrParty)

##           Don't know           Independent, near democrat
##           1              337
## Independent, near republican           Independent
##           249              502
##           No answer           Not a strong democrat
##           25              406
##           Not a strong republican           Other party
##           292              62
##           Strong democrat           Strong republican
##           419              245
##           NA's
##           2
```

95 Aside – Editing whitespace out of levels

96 Whitespace can be dealt with when data is read, or later using string manipulations. This can be done
 97 using the `trimws()` function in **base R**.

```
gender <- factor(c("male ", "male ", "male ", "male"))
```

```

levels(gender)

## [1] "male"      "male "     "male  "    "male   "

gender <- factor(trimws(gender))
levels(gender)

## [1] "male"

```

REORDERING FACTOR LEVELS

Often, factor levels have a natural ordering to them. However, the default in **base R** is to order levels alphabetically. Again, there is a fragile way to reorder the factor levels in base R, and a more robust method in the tidyverse.

Fragile method (base R)

```

summary(GSS$Opinion.of.family.income)

##      Above average      Average      Below average      Don't know
##           483           1118           666           21
## Far above average Far below average      No answer      NA's
##           65           179           6           2

levels(GSS$Opinion.of.family.income)

## [1] "Above average"      "Average"           "Below average"
## [4] "Don't know"         "Far above average" "Far below average"
## [7] "No answer"

levels(GSS$Opinion.of.family.income) <-
  levels(GSS$Opinion.of.family.income)[c(5,1:3,6,4,7)]
levels(GSS$Opinion.of.family.income)

## [1] "Far above average" "Above average"      "Average"
## [4] "Below average"    "Far below average" "Don't know"
## [7] "No answer"

```

This is both verbose and depends on the number and order of the levels staying the same. If another factor level is added to the dataset, the above code will throw an error because the number of levels differs.

However, if the code gets run more than once, the order will be broken. Particularly when working dynamically, this is all too easy to do.

```

levels(GSS$Opinion.of.family.income) <-
  levels(GSS$Opinion.of.family.income)[c(5,1:3,6,4,7)]
levels(GSS$Opinion.of.family.income)

## [1] "Far below average" "Far above average" "Above average"
## [4] "Average"          "Don't know"       "Below average"
## [7] "No answer"

```

The more times the code is run, the worse it gets. (This example illustrates why it is sometimes dangerous to replace an old version of a data frame with a new version.)

But it gets worse! It is tempting for new analysts to write code such as the following, which completely ruins the data set.

```

test <- GSS$Opinion.of.family.income

```

```
summary(test)

## Far below average Far above average Above average Average
##          483          1118          666          21
##      Don't know    Below average    No answer    NA's
##          65          179           6           2

levels(test) <- c("Far above average", "Above average", "Average", "Below Average",
  "Far below average", "Don't know", "No answer")
summary(test)

## Far above average Above average Average Below Average
##          483          1118          666          21
## Far below average Don't know No answer NA's
##          65          179           6           2
```

111 Robust method

112 ??

```
# devtools::install_github("hadley/forcats")
library(forcats)
test <- GSS$Opinion.of.family.income
test <- fct_relevel(test, c("Far above average", "Above average"))
summary(test)

## Far above average Above average Far below average Average
##          1118          666          483          21
##      Don't know    Below average    No answer    NA's
##          65          179           6           2

# XX AM please review...
```

113 COMBINING SEVERAL LEVELS INTO ONE

114 Combining discrete levels

115 This is another common task. Maybe you want fewer coefficients to interpret in your model, or the process
116 that generated the data makes a finer distinction between categories than your research. For whatever the
117 reason, you want to group together levels that are currently separate.

118 Fragile method (base R)

119 This method overwrites the labels of factor levels with repeated labels in order to group levels together.

```
levels(GSS$LaborStatus) <- c("Not employed", "No answer",
  "Other", "Not employed",
  "Not employed", "Not employed",
  "Not employed", "Employed", "Employed")
summary(GSS$LaborStatus)

## Not employed No answer Other Employed NA's
##          957           2          76       1503          2
```

120 As before, this is fragile because it depends on the order of the factor levels not changing, and on a
121 human accurately counting the indices of all the levels they wish to change.

122 Robust method

```
levels(GSS$Race.of.respondent)
```

```
## [1] "Black" "Other" "White"

GSS <- GSS %>%
  mutate(Race.of.respondent = recode(Race.of.respondent,
    `Black` = "Nonwhite",
    `Other` = "Nonwhite"))
levels(GSS$Race.of.respondent)

## [1] "Nonwhite" "White"
```

COMBINING NUMERIC-TYPE LEVELS

Combining numeric-type levels is a problem that often arises even when `stringsasfactors=FALSE`. Often variables like age or income are right censored, so there is a final category containing the lumped remainder of people. This means the data is necessarily at least a character string if not a factor. However, it may be more natural to work with numeric expressions when recoding this data.

In this data, age is provided as an integer for respondents 18-88, but then also includes the possible answer “89 or older” as well as a possible “No answer” and NA values.

We might want to turn this into a factor variable with two levels: 18–65, and over 65. In this case, it would be much easier to deal with a conditional statement about the numeric values, rather than writing out each of the numbers as a character vector.

Fragile method (base R)

In order to break this data apart as simply as possible, we need to make it numeric. To start, we recode the label for “89 or older” to read “89”. Already, we are doing something fragile.

```
levels(GSS$Age) <- c(levels(GSS$Age)[1:71], "89", "No answer")
GSS$Age <- as.numeric(as.character(GSS$Age))
summary(GSS$Age)

##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's
##      18.00   34.00   49.00   49.01   62.00   89.00     11

# XX AM I'm not understanding this: can we talk?
# Need to flesh out a base R approach to this. It's going to need some conditional logic.
```

Of course, we’re cheating a little bit here– if we were going to use this as a numeric variable in an analysis, we wouldn’t necessarily want to turn all the “89 or older” cases into the number “89”. But, we’re just on our way to a two-category factor, so those cases would have gone to the “65 and up” category one way or the other.

Robust method

```
GSS <- GSS %>%
  mutate(Age = if_else(Age < 65, "18-65", "65 and up"),
    Age = factor(Age))
summary(GSS$Age)

##      18-65 65 and up    NA's
##      2011      518      11
```

CREATING DERIVED CATEGORICAL VARIABLE

Challenges often arise when data scientists need to create derived categorical variable. As an example, consider an indicator of moderate drinking status. The National Institutes of Alcohol Abuse and Alcoholism have published guidelines for moderate drinking. These state that women, or men aged 65 or older should drink no more than one drink per day on average and no more than three drinks at a sitting. The `HELPMiss` dataset from the `mosaicData` package includes baseline data from a randomized clinical trial (Health Evaluation and Linkage to Primary Care).

	variable	description
	sex	gender of subject female or male
148	i1	average number of drinks per day (in last 30 days)
	i2	maximum number of drinks per day (in past 30 days)
	age	age (in years)

149 These guidelines can be used to create a new variable called `abstinent` for those that reported no
 150 drinking based on the value of their `i1` variable and `moderate` for those that do not exceed the NIAAA
 151 guidelines, with all other non-missing values coded as `highrisk`.

```
library(dplyr)
library(mosaic)
library(readr)
```

152 We make one value missing as a pedagogical tool to check for misbehavior of missing values.

```
data(HELPmiss)
HELPsmall <- HELPmiss %>%
  mutate(i1 = ifelse(id==1, NA, i1)) %>% # make one value missing
  select(sex, i1, i2, age)
```

153 Fragile method (base R)

```
# create empty repository for new variable
drinkstat <- character(length(HELPsmall$i1))
# create abstinent group
drinkstat[HELPsmall$i1==0] = "abstinent"
# create moderate group
drinkstat[(HELPsmall$i1>0 & HELPsmall$i1<=1 &
  HELPsmall$i2<=3 & HELPsmall$sex=="female") |
  (HELPsmall$i1>0 & HELPsmall$i1<=2 &
  HELPsmall$i2<=4 & HELPsmall$sex=="male")] = "moderate"
# create highrisk group
drinkstat[((HELPsmall$i1>1 | HELPsmall$i2>3) & HELPsmall$sex=="female") |
  ((HELPsmall$i1>2 | HELPsmall$i2>4) & HELPsmall$sex=="male")] = "highrisk"
# do we need to account for missing values?
is.na(drinkstat) <- is.na(HELPsmall$i1) | is.na(HELPsmall$i2) |
  is.na(HELPsmall$sex)
tally(~ drinkstat)

## drinkstat
## abstinent highrisk moderate <NA>
##          69          372          28          1
```

154 Robust method (dplyr)

```
glimpse(HELPsmall)
```

```
## Observations: 470
## Variables: 4
## $ sex <fctr> male, male, male, female, male, female, female, male, fem...
## $ i1 <int> NA, 56, 0, 5, 10, 4, 13, 12, 71, 20, 0, 13, 20, 13, 51, 0,...
## $ i2 <int> 26, 62, 0, 5, 13, 4, 20, 24, 129, 27, 0, 13, 31, 20, 51, 0...
## $ age <int> 37, 37, 26, 39, 32, 47, 49, 28, 50, 39, 34, 58, 58, 60, 36...

HELPSmall <- with(HELPSmall, # this won't work unless HELPSmall is made accessible to mutate()
  mutate(HELPSmall,
    drink_stat = case_when(
      i1 == 0 ~ "abstinent",
      i1 <= 1 & i2 <= 3 & sex=='female' ~ "moderate",
      i1 <= 1 & i2 <= 3 & sex=='male' & age >= 65 ~ "moderate",
      i1 <= 2 & i2 <= 4 & sex=='male' ~ "moderate",
      TRUE ~ "highrisk"
    ))
tally( ~ drink_stat, data = HELPSmall)

## drink_stat
## abstinent highrisk moderate
##          69      373      28
```

DEFENSIVE CODING

It is always good practice to write conditional testing statements into code using factors. As an example, we can assert that there are three levels for the drinking status variable in the HELP dataset.

```
library(assertthat)
levels(as.factor(drinkstat))

## [1] "abstinent" "highrisk" "moderate"

assert_that(length(levels(as.factor(drinkstat)))==3)

## [1] TRUE
```

Similar code could be used to check other conditions to verify that recoding has been done successfully. Here is some code that doesn't work:

```
expect_equivalent(levels(GSS$Respondents.sex), c("Male", "Female"))
```

ACKNOWLEDGEMENTS

IDEAS

Two ways to do each thing (as long as one isn't totally stupid) Why is this hard? Why is this error-prone? Missing values Appendices for less interesting examples?

LOADING THE DATA

We provide several options for how to get this data. We could download it in SPSS or Stata formats and use the foreign package to read it in. The GSS download even provides an R file to do the translation for you. Here is the result of that:

```
source('../data/GSS.r')
```

```
glimpse(GSS)

## Observations: 2,538
## Variables: 17
## $ YEAR      <int> 2014, 2014, 2014, 2014, 2014, 2014, 2014, 2014, 2014, ...
## $ ID_       <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16...
## $ WRKSTAT   <int> 1, 1, 4, 2, 5, 1, 9, 1, 8, 1, 7, 8, 5, 1, 6, 2, 2, 1, ...
## $ PRESTIGE  <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ MARITAL   <int> 3, 1, 3, 1, 1, 1, 1, 1, 5, 1, 1, 5, 3, 1, 5, 1, 3, 5, ...
## $ CHILDS    <int> 0, 0, 1, 2, 3, 1, 2, 2, 4, 3, 2, 0, 5, 2, 0, 3, 3, 0, ...
## $ AGE       <int> 53, 26, 59, 56, 74, 56, 63, 34, 37, 30, 43, 56, 69, 4...
## $ EDUC      <int> 16, 16, 13, 16, 17, 17, 12, 17, 10, 15, 5, 11, 8, 11, ...
## $ SEX       <int> 1, 2, 1, 2, 2, 2, 1, 1, 2, 2, 2, 1, 1, 2, 2, 1, 2, 1, ...
## $ RACE      <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 3, 3, 1, 1, 1, 1, 1, 3, 1, ...
## $ INCOM16    <int> 2, 3, 2, 2, 4, 4, 2, 3, 3, 1, 1, 2, 2, 2, 2, 3, 2, 3, ...
## $ INCOME     <int> 12, 12, 12, 12, 13, 12, 13, 12, 10, 12, 9, 9, 10, 11, ...
## $ RINCOME    <int> 12, 12, 0, 9, 0, 12, 13, 12, 0, 12, 0, 0, 0, 11, 12, ...
## $ INCOME72   <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ...
## $ PARTYID    <int> 5, 5, 6, 5, 3, 6, 6, 8, 3, 3, 3, 3, 3, 1, 3, 6, 1, 3, ...
## $ FINRELA    <int> 4, 4, 2, 4, 3, 4, 9, 3, 2, 3, 8, 5, 1, 1, 3, 3, 2, 3, ...
## $ SEXORNT    <int> 3, 3, 3, 3, 3, 9, 0, 0, 3, 3, 3, 3, 3, 0, 3, 3, 0, 0, ...
```

168 Obviously, this is less than ideal. Now, all the factor variables are encoded as integers, but their level
 169 labels have been lost. We have to look at a codebook to determine if SEX == 1 indicates male or female.
 170 We would rather preserve the integrated level labels. In order to do this, our best option is to download
 171 the data as an Excel file and use the **readxl** package to load it.

```
library(readxl)
GSS <- read_excel("../data/GSS.xls")
names(GSS) <- make.names(names(GSS), unique=TRUE)
glimpse(GSS)

## Observations: 2,540
## Variables: 17
## $ Gss.year.for.this.respondent..... <dbl> 2014, 2014...
## $ Respondent.id.number              <dbl> 1, 2, 3, 4...
## $ Labor.force.status                 <chr> "Working f...
## $ Rs.occupational.prestige.score...1970. <dbl> 0, 0, 0, 0...
## $ Marital.status                    <chr> "Divorced"...
## $ Number.of.children                 <dbl> 0, 0, 1, 2...
## $ Age.of.respondent                  <chr> "53.000000...
## $ Highest.year.of.school.completed   <dbl> 16, 16, 13...
## $ Respondents.sex                   <chr> "Male", "F...
## $ Race.of.respondent                 <chr> "White", "...
## $ Rs.family.income.when.16.yrs.old   <chr> "Below ave...
## $ Total.family.income                 <chr> "$25000 or...
## $ Respondents.income                 <chr> "$25000 or...
## $ Total.family.income.1               <chr> "Not appli...
## $ Political.party.affiliation         <chr> "Not str r...
## $ Opinion.of.family.income           <chr> "Above ave...
## $ Sexual.orientation                 <chr> "Heterosex..."
```

172 That's a little better. Now we have preserved the character strings. But, the data is not yet usable in an
 173 analysis.

174 A RENAMING THE VARIABLES

175 One problem is that the variable names (while human readable) are full of spaces, so are hard to use. But,
 176 we can rename them.

177 There is a fragile way to do this in **base R**, but we'll use the more robust `rename()` function from
 178 the **dplyr** package. `rename()`

```
library(dplyr)
```

```
GSS <- GSS %>%
  rename(LaborStatus = Labor.force.status,
         PolParty = Political.party.affiliation,
         Age = Age.of.respondent)
write_csv(GSS, path="../data/GSScleaned.csv")
```

B CLOSING EXERCISE

We have included the following as a possible closing exercise.

Subjects in the HELP study were also categorized into categories of drug and alcohol involvement, as displayed in the following table.

```
HELPbase <- HELPfull %>%
  filter(TIME==0)
tally( ~ PRIM_SUB + SECD_SUB, data=HELPbase)
```

		SECD_SUB								
	PRIM_SUB	0	1	2	3	4	5	6	7	8
	1	99	0	57	13	1	3	11	0	1
	2	51	84	0	6	0	0	15	0	0
	3	57	28	29	0	0	6	5	1	2
	6	0	1	0	0	0	0	0	0	0

Note that the following codings of substance use involvement were used:

value	description
0	None
1	Alcohol
2	Cocaine
3	Heroin
4	Barbituates
5	Benzos
6	Marijuana
7	Methadone
8	Opiates

Create a new variable called 'primsub' that combines the primary and secondary substances into a categorical variable with values corresponding to primary and secondary substances of the form: alcohol only, cocaine only, 'heroin only', 'alcohol-cocaine', 'cocaine-alcohol', or 'other'. Code any group with fewer than 5 entries as 'alcohol-other', 'cocaine-other', or 'heroin-other'. If 'PRIM_SUB==6' make the 'primsub' variable missing.

How many subjects are there in the 'alcohol-none' group? How many subjects are there in the 'alcohol-other' group? What are the three most common groups?

SOLUTION:

```
HELPbase <- with(HELPbase,
```

```
mutate(HELPhbase,
  primary= recode(PRIM_SUB,
    `1`="alcohol",
    `2`="cocaine",
    `3`="heroin",
    `4`="barbituates",
    `5`="benzos",
    `6`="marijuana",
    `7`="methadone",
    `8`="opiates"),
  second=recode(SECD_SUB,
    `0`="none",
    `1`="alcohol",
    `2`="cocaine",
    `3`="heroin",
    `4`="barbituates",
    `5`="benzos",
    `6`="marijuana",
    `7`="methadone",
    `8`="opiates"),
  title=paste0(primary, "-", second)
))
```

```
tally(~ primary, data=HELPhbase)

## primary
##   alcohol   cocaine   heroin marijuana
##      185      156      128         1

tally(~ second, data=HELPhbase)

## second
##   alcohol barbituates   benzos   cocaine   heroin   marijuana
##      113         1         9       86      19         31
## methadone      none   opiates
##         1      207         3

counts <- HELPhbase %>%
  group_by(primary, second) %>%
  summarise(observed=n())

merged <- left_join(HELPhbase, counts, by=c("primary", "second"))
```

```
merged <- with(merged,
```

```
mutate(merged,
  title =
    case_when(
      observed < 5 & primary=="alcohol" ~ "alcohol-other",
      observed < 5 & primary=="cocaine" ~ "cocaine-other",
      observed < 5 & primary=="heroin" ~ "heroin-other",
      TRUE ~ title),
  title = ifelse(primary=="marijuana", NA, title)))
tally(~ title + observed, data=merged)

##           observed
## title
## alcohol-cocaine    1  2  3  5  6 11 13 15 28 29 51 57 84 99
## alcohol-heroin     0  0  0  0  0  0 13  0  0  0  0  0  0  0
## alcohol-marijuana  0  0  0  0  0 11  0  0  0  0  0  0  0  0
## alcohol-none       0  0  0  0  0  0  0  0  0  0  0  0  0 99
## alcohol-other      2  0  3  0  0  0  0  0  0  0  0  0  0  0
## cocaine-alcohol    0  0  0  0  0  0  0  0  0  0  0  0 84  0
## cocaine-heroin     0  0  0  0  6  0  0  0  0  0  0  0  0  0
## cocaine-marijuana  0  0  0  0  0  0  0 15  0  0  0  0  0  0
## cocaine-none       0  0  0  0  0  0  0  0  0  0 51  0  0  0
## heroin-alcohol      0  0  0  0  0  0  0  0 28  0  0  0  0  0
## heroin-benzos       0  0  0  0  6  0  0  0  0  0  0  0  0  0
## heroin-cocaine      0  0  0  0  0  0  0  0  0 29  0  0  0  0
## heroin-marijuana    0  0  0  5  0  0  0  0  0  0  0  0  0  0
## heroin-none         0  0  0  0  0  0  0  0  0  0  0 57  0  0
## heroin-other        1  2  0  0  0  0  0  0  0  0  0  0  0  0
## <NA>                1  0  0  0  0  0  0  0  0  0  0  0  0  0
```

```
tally(~ title=="alcohol-none", data=merged)

## title == "alcohol-none"
## TRUE FALSE <NA>
##    99   370     1

tally(~ title=="alcohol-other", data=merged)

## title == "alcohol-other"
## TRUE FALSE <NA>
##     5   464     1

sort(tally(~ title, data=merged), decreasing=TRUE)[1:3]

## title
## alcohol-none cocaine-alcohol alcohol-cocaine
##           99             84             57
```

REFERENCES

- Broman, K. (2015). Initial steps toward reproducible research. <http://kbroman.org/steps2rr/>.
- Peng, R. D. (2015). stringsAsFactors: An unauthorized biography. <http://simplystatistics.org/2015/07/24/stringsasfactors-an-unauthorized-biography/>.
- Wickham, H. (2014). Tidy data. *Journal of Statistical Software*, 59(10).