

1 Wrangling categorical data in R

2 Amelia McNamara¹

3 ¹Smith College

4 ABSTRACT

5 Wrangling of categorical data is an important part of the analysis cycle. Many aspects of these operations
6 can be tricky, particularly for complex transformations. This paper discusses aspects of transformation of
7 categorical variables in R. We suggest defensive coding strategies and principles for data wrangling to
8 ensure data quality and sound analysis.

9 Keywords:

10 INTRODUCTION

11 The wrangling of categorical data is an important component in data science because so many variables
12 are categorical. Gender, income bracket, and state are all examples of categorical data. While defensive
13 coding is important for any analysis, categorical data presents particular problems that can slip in without
14 the analyst noticing.

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

18 In particular, we are considering how categorical data is treated in **base R** versus the so-called
19 tidyverse (Wickham, 2014). Tools from the tidyverse are discussed in another paper in this special issue,
20 but briefly they aim to make analysis more pure, predictable, and pipeable. One canonical thought exercise
21 is whether, after the analysis, a new version of the data could be supplied in the code and have parallel
22 results come out the other end. Again, categorical data can make this even more complex.

23

24 FACTORS IN R

25 Consider a gender variable including the categories `male`, `female` and `gender non-conforming`.
26 In R, there are two ways to store this information. One is to use a series of character strings, and the other
27 is to store it as a factor. Historically, storing categorical data as a factor variable was more efficient than
28 storing the same data as strings, because factor variables only store the factor labels once (Peng, 2015).
29 However, R has changed to use hashed versions of all character strings, so the storage issue is no longer a
30 consideration (Peng, 2015).

31 Factors can be very tricky to deal with, since many operations applied to them return different values
32 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
```

33 This is unexpected because `as.numeric()` feels like the way to recover numeric information in
34 the **base R** paradigm. Compare the following:

```
as.numeric(c("hello"))  
  
## [1] NA  
  
as.numeric(factor(c("hello")))  
  
## [1] 1
```

35 This behavior has led to an online movement against the default behavior of many of R's data import
36 functions to take any variable composed as strings and automatically convert the variable to a factor. The
37 tidyverse moves away from this default behavior, with functions from the **readr** package defaulting to
38 leaving strings as-is. (Others have chosen to add `options(stringAsFactors=FALSE)` into their
39 startup commands.)

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

47 So, factors are important. But, they can often be hard to deal with. Because of the way the group
48 numbers are stored separately from the factor labels, it can be easy to overwrite data in such a way that
49 the original data is lost. In this paper, we will consider the best practices for working with factor data.

50 To do this, we will consider data from the General Social Survey. There are some import issues
51 inherent to the data which are not particular to categorical data, so that processing is in Appendices ??.
52 We'll work with the data that has cleaned variable names.

```
library(readr)  
GSS <- read_csv("../data/GSScleaned.csv")  
# Maybe I should use read.csv so I don't have to convert things to factors?
```

53 CONSIDERING SOME FACTOR VARIABLES

```
GSS <- GSS %>%  
  mutate(LaborStatus = factor(LaborStatus),  
         PolParty = factor(PolParty))  
  
levels(GSS$LaborStatus) # I wish I had a piece of dplyr code for this  
  
## [1] "Keeping house"      "No answer"          "Other"  
## [4] "Retired"            "School"             "Temp not working"  
## [7] "Unempl, laid off"   "Working fulltime"   "Working parttime"  
  
levels(GSS$PolParty)  
  
## [1] "Don't know"          "Ind,near dem"       "Ind,near rep"  
## [4] "Independent"         "No answer"          "Not str democrat"  
## [7] "Not str republican"  "Other party"        "Strong democrat"  
## [10] "Strong republican"
```

54 TASK 1: CHANGING THE LABELS OF FACTORS (BASE R)

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

57 Here is how I do that in base R. Typically, I end up ruining something in the process of doing this, so I
58 **always** start with a summary call, to check after I have done my attempt.

```
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

with(GSS, summary(LaborStatus)) # I prefer this to the $

##      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
```

59 This method is less than ideal, because it depends on the data coming in with the factor levels ordered
60 in a particular way. If the data gets changed outside of R (perhaps with an additional level), this code will
61 not perform as expected. (Additionally, if the code gets run more than once, it can also lead to unexpected
62 behavior since the result is being added back into the original variable.)

63 CHANGING THE LABELS OF FACTORS (DPLYR)

64 In the **dplyr** package, you can use the `recode` function to do the same thing. There are a few things to
65 remember with `recode`. The first is that it is a vector function, which means it must be used within a
66 `mutate` call or with a variable pulled out using `$`. The second is that you need to tell it which variable to
67 recode, even if you are overwriting an existing variable.

```
GSS <- GSS %>%
  mutate(PolParty = recode(PolParty, `Not str republican` = "Not a strong republican"))
```

68 COMBINING SEVERAL LEVELS INTO ONE

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

72 How I do this in base R:

```
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
```

73 1 MOSAIC COMBINING LEVELS

74 XX Is there any need to use mosaic now that case_when is included in dplyr?

```
library(mosaic)
data(Births78)
Births78 <- Births78 %>%
  mutate(weekend = derivedFactor(weekend = wday== "Sun" | wday == "Sat", .default="weekday"))
```

75 COMBINING MANY CATEGORIES INTO ONE

76 In this data, age is provided as an integer for respondents 18-88, but then also includes the possible answer
77 "89 or older" as well as a possible "No answer" and NA values.

```
GSS <- GSS %>%
  mutate(Age = factor(Age))
summary(GSS$Age)
```

##	18.000000	19.000000	20.000000	21.000000	22.000000	23.000000
##	6	25	26	24	28	30
##	24.000000	25.000000	26.000000	27.000000	28.000000	29.000000
##	31	48	47	41	31	51
##	30.000000	31.000000	32.000000	33.000000	34.000000	35.000000
##	57	49	55	47	46	40
##	36.000000	37.000000	38.000000	39.000000	40.000000	41.000000
##	40	54	47	52	46	54
##	42.000000	43.000000	44.000000	45.000000	46.000000	47.000000
##	35	54	39	41	34	43
##	48.000000	49.000000	50.000000	51.000000	52.000000	53.000000
##	32	39	54	45	37	60
##	54.000000	55.000000	56.000000	57.000000	58.000000	59.000000
##	53	52	60	43	60	47
##	60.000000	61.000000	62.000000	63.000000	64.000000	65.000000
##	46	38	44	42	38	40
##	66.000000	67.000000	68.000000	69.000000	70.000000	71.000000
##	35	41	21	23	32	28
##	72.000000	73.000000	74.000000	75.000000	76.000000	77.000000
##	20	22	25	21	24	17
##	78.000000	79.000000	80.000000	81.000000	82.000000	83.000000
##	28	26	16	14	8	11
##	84.000000	85.000000	86.000000	87.000000	88.000000	89 or older
##	13	6	9	8	11	19
##	No answer	NA's				
##	9	2				

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

81 But, in order to do that we need to make it numeric.

```
# GSS$Age[GSS$Age == "No answer"] <- NA # Do I really need this? Nope!
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

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

```
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
```

86 Another way to do this:

```
# young <- as.character(18:64)
# derivedVariable(Age %in% young = "18-65", Age )
```

87 2 OTHER EXAMPLES

88 Here's a placeholder for the other examples.

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

89 2.1 Task 2: creating derived categorical variable

90 XX THESE ARE STILL WORDED AS TASKS

91 The National Institutes of Alcohol Abuse and Alcoholism have published guidelines for moderate
 92 drinking. These state that women, or men aged 65 or older should drink no more than one drink per day
 93 on average and no more than three drinks at a sitting.

94 The HELPMiss dataset from the **mosaicData** package includes baseline data from a randomized
 95 clinical trial (Health Evaluation and Linkage to Primary Care).

variable	description
sex	gender of subject female or male
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)

97 Use these guidelines and the HELPSmall dataset to create a new variable called abstinent for
 98 those that reported no drinking based on the value of their i1 variable and moderate for those that do
 99 not exceed the NIAAA guidelines. All other non-missing values should be coded as highrisk.

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

```
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...

# I definitely want to remove these ASAP
#attach(HELPsmall)

HELPsmall <- with(HELPsmall, # this won't work unless HELPsmall is made accessible
  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)

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

2.2 Task 3: Creating derived categorical variables

XX move to appendix (since it duplicates the earlier example?)

Subjects in the HELP study were 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?

```

HELPPbase <- with(HELPPbase,
  mutate(HELPPbase,
    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=HELPPbase)

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

tally(~ second, data=HELPPbase)

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

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

merged <- left_join(HELPPbase, 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
```

113 **Answers:**

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

##
## TRUE FALSE <NA>
##    99   370     1

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

##
## TRUE FALSE <NA>
##     5   464     1

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

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

114 **ACKNOWLEDGEMENTS**

115 **IDEAS**

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

118 **LOADING THE DATA**

119 **??**

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


```

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

## 'data.frame': 2538 obs. of 17 variables:
## $ YEAR : int 2014 2014 2014 2014 2014 2014 2014 2014 2014 2014 2014 ...
## $ ID_ : int 1 2 3 4 5 6 7 8 9 10 ...
## $ WRKSTAT : int 1 1 4 2 5 1 9 1 8 1 ...
## $ PRESTIGE: int 0 0 0 0 0 0 0 0 0 0 ...
## $ MARITAL : int 3 1 3 1 1 1 1 1 5 1 ...
## $ CHILDS : int 0 0 1 2 3 1 2 2 4 3 ...
## $ AGE : int 53 26 59 56 74 56 63 34 37 30 ...
## $ EDUC : int 16 16 13 16 17 17 12 17 10 15 ...
## $ SEX : int 1 2 1 2 2 2 1 1 2 2 ...
## $ RACE : int 1 1 1 1 1 1 1 1 1 3 ...
## $ INCOM16 : int 2 3 2 2 4 4 2 3 3 1 ...
## $ INCOME : int 12 12 12 12 13 12 13 12 10 12 ...
## $ RINCOME : int 12 12 0 9 0 12 13 12 0 12 ...
## $ INCOME72: int 0 0 0 0 0 0 0 0 0 0 ...
## $ PARTYID : int 5 5 6 5 3 6 6 8 3 3 ...
## $ FINRELA : int 4 4 2 4 3 4 9 3 2 3 ...
## $ SEXORNT : int 3 3 3 3 3 9 0 0 3 3 ...
## - attr(*, "col.label")= chr "Gss year for this respondent" "Respondent

```

Obviously, this is less than ideal. Now, all the factor variables are encoded as integers, but their level labels have been lost. We have to look at a codebook to determine if SEX == 1 indicates male or female. We would rather preserve the integrated level labels. In order to do this, our best option is to download 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)
str(GSS)

## Classes 'tbl_df', 'tbl' and 'data.frame': 2540 obs. of 17 variables:
## $ Gss.year.for.this.respondent.....: num 2014 2014 2014 2014 2014 ...
## $ Respondent.id.number : num 1 2 3 4 5 6 7 8 9 10 ...
## $ Labor.force.status : chr "Working fulltime" "Working fulltime" ...
## $ Rs.occupational.prestige.score...1970. : num 0 0 0 0 0 0 0 0 0 0 ...
## $ Marital.status : chr "Divorced" "Married" "Divorced" "Married" ...
## $ Number.of.children : num 0 0 1 2 3 1 2 2 4 3 ...
## $ Age.of.respondent : chr "53.000000" "26.000000" "59.000000" ...
## $ Highest.year.of.school.completed : num 16 16 13 16 17 17 12 17 10 15 ...
## $ Respondents.sex : chr "Male" "Female" "Male" "Female" ...
## $ Race.of.respondent : chr "White" "White" "White" "White" ...
## $ Rs.family.income.when.16.yrs.old : chr "Below average" "Average" "Below average" ...
## $ Total.family.income : chr "$25000 or more" "$25000 or more" ...
## $ Respondents.income : chr "$25000 or more" "$25000 or more" ...
## $ Total.family.income.1 : chr "Not applicable" "Not applicable" ...
## $ Political.party.affiliation : chr "Not str republican" "Not str republican" ...
## $ Opinion.of.family.income : chr "Above average" "Above average" ...
## $ Sexual.orientation : chr "Heterosexual or straight" "Heterosexual or straight" ...

```

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

RENAMING THE VARIABLES

??

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

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

REFERENCES

- 135
- 136 Peng, R. D. (2015). stringsAsFactors: An unauthorized biography. <http://simplystatistics.org/2015/07/24/stringsasfactors-an-unauthorized-biography/>.
- 137
- 138 Wickham, H. (2014). Tidy data. *Journal of Statistical Software*, 59(10).