Transforming Data

Pedram Navid

September 01, 2016



Overview

Why transform?

Data rarely comes in the format you want it Even if it does, you always want to do something with it Good work upfront saves headache down the road, so don't skimp!

Tidy Data

Data should be tidy when working with R. We have a tendency to create untidy data, especially in finance.

"tidy datasets are all alike but every messy dataset is messy in its own way"

Tidy data is:

- Each variable is its own column
- Each observation is in its own row

First Steps

Use dplyr and tidyr to help clean and manipulate data. Check out the cheatsheet: Help > Cheatsheets > Data Manipulation with dplyr,tidyr

Syntax

```
# Load required packages
library(dplyr)
library(tidyr)

# Let's say we have a dataset
iris

## Sepal.Length Sepal.Width Petal.Length Petal.Width
```

##		sebar.rengun	Sepai.width	retal.Length	retar.width
##	1	5.1	3.5	1.4	0.2
##	2	4.9	3.0	1.4	0.2
##	3	4.7	3.2	1.3	0.2
##	4	4.6	3.1	1.5	0.2
##	5	5.0	3.6	1.4	0.2
##	6	5.4	3.9	1.7	0.4
##	7	4.6	3.4	1.4	0.3
##	8	5.0	3.4	1.5	0.2
##	9	4.4	2.9	1.4	0.2
##	10	4.9	3.1	1.5	0.1

A better way: syntax continued

```
iris_easy <- tbl_df(iris)</pre>
iris_easy
```

```
## # A tibble: 150 × 5
##
      Sepal.Length Sepal.Width Petal.Length Petal.Width Spe
              /11-15
```

##	<ab1></ab1>	<ab1></ab1>	<ab1></ab1>	<ab1></ab1>	< ;
##	1 5.1	3.5	1.4	0.2	S
##	2 4.9	3.0	1.4	0.2	S
##	3 4 7	3.2	1.3	0.2	S

					· · -	-
##	2	4.9	3.0	1.4	0.2	S
##	3	4.7	3.2	1.3	0.2	S

## 2	4.9	3.0	1.4	0.2	S
## 3	4.7	3.2	1.3	0.2	S
## 4	4.6	3.1	1.5	0.2	S

## 3	4.1	3.2	1.5	0.2	D
## 4	4.6	3.1	1.5	0.2	S
## 5	5.0	3.6	1.4	0.2	s
## 6	5.4	3.9	1.7	0.4	S

					-
## 6	5.4	3.9	1.7	0.4	s
## 7	4.6	3.4	1.4	0.3	S
## 8	5.0	3 4	1 5	0.2	9

##	3	4.7	3.2	1.3	0.2
##	4	4.6	3.1	1.5	0.2
##	5	5.0	3.6	1.4	0.2
	0	_ 4	0 0	4 7	A 4

S

S

0.2

4.4 2.9 1.4 1.5 0.1

10 4.9 3.1 ... with 140 more rows

View

View(iris_easy)

$\Diamond \Diamond$	Æ Filter				
	Sepal.Lengtĥ	Sepal.Widtĥ	Petal.Lengtĥ	Petal.Widtĥ	Species ‡
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa
5	5.0	3.6	1.4	0.2	setosa
6	5.4	3.9	1.7	0.4	setosa
7	4.6	3.4	1.4	0.3	setosa
8	5.0	3.4	1.5	0.2	setosa

Figure 1:

Tidy Data

An example:

Is this data tidy? Why/why not?

A tibble: 4 × 13

```
library(readxl)
wide_sales <- read_excel("Transforming_Data/data/wide_data
head(wide_sales)</pre>
```

```
##
    Region March April May
                                  July August September
                            June
##
     <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <
                                       <dbl>
                                                 <dbl>
## 1
     North 5317 2858 4920 5787
                                  2119
                                        2439
                                                 3041
## 2 East 4850 1023 6169 8815
                                  6042
                                        4366
                                                 3869
## 3
    West 4249 6761
                        100 5685 8408
                                        3427
                                                 8939
## 4
     South 4195 3285 1654 5257 4206
                                        2144
                                                 8361
## #
    ... with 3 more variables: December <dbl>, January <dl
```

Reshaping data: gather

```
tidy_sales <- wide_sales %>%
  gather("month", "sales", 2:13)
head(tidy_sales)
```

```
## # A tibble: 6 × 3
## Region month sales
## <chr> <chr> <chr> <chr> <dbl>
## 1 North March 5317
## 2 East March 4850
## 3 West March 4249
## 4 South March 4195
## 5 North April 2858
## 6 East April 1023
```

Reshaping data: spread

```
tidy_sales %>%
   spread(Region, sales)
```

```
A tibble: 12 × 5
          month
                  East North South
##
                                     West
## *
          <chr> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1
          April
                  1023
                        2858
                               3285
                                     6761
## 2
         August
                  4366
                        2439
                               2144
                                     3427
## 3
       December
                  9792
                       8880
                               9652
                                     4974
## 4
       February
                  9746
                       4526
                               6991
                                     5139
## 5
        January
                  4806
                        3644
                                983
                                      766
                  6042
                        2119
                               4206
                                     8408
## 6
            July
                  8815
                        5787
                               5257
                                     5685
## 7
            June
                                     4249
## 8
          March
                  4850
                        5317
                               4195
##
                  6169
                        4920
                               1654
                                       100
            May
##
   10
       November
                  2635
                        5128
                               7307
                                     2893
##
   11
        October
                  7819
                        3811
                               7400
                                     8959
                  3869
                        3041
                               8361
                                     2939
## 10 Sentember
```

Manipulate Data

Intro to Dplyr

You will spend more time with dplyr than any other R package. First thing to understand is the pipe operator: %>%

%>% passes the result of the previous line to the next line, making your code much more readable

Consider this fake example:

```
miss_muffet_sitting <- sit(miss_muffet, tuffet)
miss_muffet_eating <- eat(miss_muffet_eating, curds_and_who
muffet_spider <- sit(miss_muffet_eating, spider, where = 'I
muffet_away <- run(muffet_spider, reason = 'frightened')</pre>
```

Or..

```
miss_muffet %>%
  sit(tuffet) %>%
  eat(curds_and_whey) %>%
```

sit(spider, where = 'beside') %>%

run(reason = 'frightened')

dplyr: continued

8

Best way to learn it is to use it. . . . a lot.

```
library(dplyr)
parking <- read.csv("Transforming_Data/data/parking_data.cs</pre>
  tbl df()
parking
## # A tibble: 234,983 × 11
```

##	tag_number_masked	date_or_infraction	iniraction_code
##	<fctr></fctr>	<int></int>	<int></int>
## 1	***42781	20141123	210
## 2	***30955	20141123	210

##	<fctr></fctr>	<int></int>	<int></int>
##	1 ***42781	20141123	210
##	2 ***30955	20141123	210
##	3 ***57421	20141123	210

***14417 20141123 ## 4

5 ***35411 20141123

210 207

6 ***49340 20141123

***17984 20141123 3 ## 7

20141123

207

***30956

dplyr: select

A tibble: 6×3

Pick columns, and optionally rename and reorder them

```
# Pick
parking %>%
select(date_of_infraction, infraction_description, 5) %>9
head()
```

```
##
    date of infraction
                                infraction_description set
##
                  <int>
                                                <fctr>
## 1
              20141123 PARK FAIL TO DISPLAY RECEIPT
## 2
              20141123 PARK FAIL TO DISPLAY RECEIPT
## 3
              20141123 PARK FAIL TO DISPLAY RECEIPT
## 4
              20141123 PARK FAIL TO DISPLAY RECEIPT
## 5
              20141123 PARK FAIL TO DEP. FEE MACHINE
               20141123 PARK-SIGNED HWY-PROHIBIT DY/TM
## 6
```

```
# Rename
parking %>%
select(inf_date = date_of_infraction, desc = infraction_of
head()
```

```
## # A tibble: 6 \times 4
     inf_date
                                             fine
                                                   time
##
                                       desc
##
        <int>
                                      <fctr> <int> <int>
## 1 20141123
               PARK FAIL TO DISPLAY RECEIPT
                                               30
                                                   1554
## 2 20141123
             PARK FAIL TO DISPLAY RECEIPT
                                               30
                                                   1554
## 3 20141123 PARK FAIL TO DISPLAY RECEIPT
                                               30
                                                   1554
## 4 20141123 PARK FAIL TO DISPLAY RECEIPT
                                               30
                                                   1554
## 5 20141123 PARK FATL TO DEP, FEE MACHINE
                                                   1554
                                               30
## 6 20141123 PARK-SIGNED HWY-PROHIBIT DY/TM
                                               40
                                                   1554
```

```
# More picking
parking %>%
  select(starts_with('location')) %>%
  head()
```

```
## # A tibble: 6 \times 4
                          location2 location3
                                                  location4
##
     location1
                                        <fctr>
##
        <fctr>
                             <fctr>
                                                      <fctr>
## 1
           N/S
                       LOURDES LANE
                                           W/O HOMEWOOD AVE
## 2
            NR.
                      391 KING ST W
            NR.
                   96 ST PATRICK ST
## 3
                     986 BLOOR ST W
## 4
           UPP
## 5
            NR.
                   111 ELIZABETH ST
## 6
            NR 56 BISHOP TUTU BLVD
```

```
parking %>%
  select(-starts_with('location')) %>%
  head()
```

A tibble: 6×7

<fctr>

***42781

##

##

1

```
***30955
                               20141123
## 2
                                                      210
## 3
             ***57421
                              20141123
                                                      210
## 4
             ***14417
                              20141123
                                                      210
## 5
              ***35411
                              20141123
                                                      207
              ***49340
                                20141123
## 6
                                                        5
## # ... with 4 more variables: infraction description <fc
## #
      set fine amount <int>, time of infraction <int>, pro
```

tag number masked date of infraction infraction code

<int>

20141123

<int>

210

Dplyr: Filter

2

3 ## 4

5

6 ## 7

8 ## 9

10

Subset data using logic

```
parking %>%
  filter(set_fine_amount > 300)
```

20141123

20141123

20141123

20141123

20141123

20141123

20141123

20141123

20141124

367

355

355

367

363

367

355

355

363

***49346

***61327

***30223

***13649

***55315

***49375

***88513

***03453

***53656

```
parking %>%
  filter(infraction_code == 355, province != "ON")
  # A tibble: 18 × 11
      tag_number_masked date_of_infraction infraction_code
##
##
                 <fctr>
                                      <int>
                                                      <int>
               ***45827
                                   20141127
## 1
## 2
               ***68860
                                   20141127
```

20141128

20141129

20141202

20141204

20141214

20141214

20141215

20141216

20141216

20141218

20141218

***69264

***11902

***79981

***11996

***42247

***45845

***80547

***44391

***84360

***69611

***48790

3

4

5

6 ## 7

8

9

10

11

12

13

355

355

355

355

355

355

355

355

355

355

355

355

355

Distinct

```
parking %>%
  select(infraction_code, infraction_description) %>%
  distinct()
```

```
# A tibble: 146 × 2
      infraction_code
                               infraction_description
##
##
                <int>
                                               <fctr>
                        PARK FAIL TO DISPLAY RECEIPT
## 1
                  210
## 2
                  207
                       PARK FAIL TO DEP. FEE MACHINE
## 3
                    5 PARK-SIGNED HWY-PROHIBIT DY/TM
## 4
                    3
                             PARK ON PRIVATE PROPERTY
                        STAND VEH.-PROHIBIT TIME/DAY
## 5
                   16
                       PARK-WITHIN 9M INTERSECT ROAD
## 6
                                 PARK IN A FIRE ROUTE
## 7
                  347
## 8
                  406 PARK-VEH. W/O VALID ONT PLATE
                        PARK IN ACCESSIBLE NO PERMIT
## 9
                  355
##
  10
                    6 PARK-SIGNED HWY-EXC PERMT TIME
##
     ... with 136 more rows
```

Samples

8

9

10

```
set.seed(42)
parking %>%
  select(location2, infraction_description, set_fine_amoun
  sample_n(10, replace = FALSE)
## # A tibble: 10 × 3
##
                location2
                                  infraction_description se
##
                   <fctr>
                                                   <fctr>
          ST CLAIR AVE W
## 1
                             PARK-SIGNED HWY-PUBLIC LANE
## 2
            225 KING ST W STOP-SIGNED HIGHWAY-RUSH HOUR
## 3
             16 NEWTON DR PARK-SIGNED HWY-PROHIBIT DY/TM
## 4
      700 HUMBERWOOD BLVD
                                PARK ON PRIVATE PROPERTY
## 5
           225 QUEBEC AVE PARK PROHIBITED TIME NO PERMIT
## 6
           BRENTWOOD RD N PARK FAIL TO DISPLAY RECEIPT
## 7
        20 EGLINTON AVE W STAND VEH.-PROHIBIT TIME/DAY
```

ST NICHOLAS ST PARK-SIGNED HWY-PROHIBIT DY/TM

14 COLLEGE ST PARK FAIL TO DEP. FEE MACHINE

256 DONIFA OR PARK-SIGNED HWY-PROHIRIT DY/TM

Slice

```
parking %>%
slice(12:15)
```

```
## # A tibble: 4 × 11
##
     tag_number_masked date_of_infraction infraction_code
##
                <fctr>
                                     <int>
                                                      <int>
## 1
              ***49341
                                  20141123
                                                          5
              ***42782
                                  20141123
                                                        207
## 2
## 3
              ***14419
                                  20141123
                                                        210
## 4
              ***14957
                                 20141123
## #
     ... with 8 more variables: infraction description <fc
## #
       set_fine_amount <int>, time_of_infraction <int>, loc
## #
       location2 <fctr>, location3 <fctr>, location4 <fctr>
```

Top

#

head(parking, n = 8)

```
## # A tibble: 8 × 11
##
     tag_number_masked date_of_infraction infraction_code
##
                <fctr>
                                     <int>
                                                      <int>
## 1
              ***42781
                                  20141123
                                                        210
                                  20141123
## 2
              ***30955
                                                        210
## 3
              ***57421
                                  20141123
                                                        210
## 4
              ***14417
                                 20141123
                                                        210
## 5
              ***35411
                                20141123
                                                        207
                                                          5
## 6
              ***49340
                                20141123
                                                          3
## 7
              ***17984
                                  20141123
## 8
              ***30956
                                  20141123
                                                        207
     ... with 8 more variables: infraction description <fc
## #
       set fine amount <int>, time of infraction <int>, loc
## #
```

location2 <fctr>, location3 <fctr>, location4 <fctr>

Bottom

tail(parking)

```
## # A tibble: 6 × 11
     tag_number_masked date_of_infraction infraction_code
##
##
                <fctr>
                                     <int>
                                                      <int>
                                  20141231
                                                          3
## 1
              ***48217
## 2
              ***67949
                                  20141231
                                                         29
              ***60555
                                  20141231
                                                          5
## 3
## 4
              ***45157
                                 20141231
                                                         29
## 5
              ***48218
                                  20141231
                                                          3
## 6
              ***87324
                                  20141231
                                                          9
## #
     ... with 8 more variables: infraction description <fc
## #
       set_fine_amount <int>, time_of_infraction <int>, loc
## #
       location2 <fctr>, location3 <fctr>, location4 <fctr>
```

```
Top X
   parking %>%
     filter(province != "ON") %>%
     top_n(n = 10, wt = set_fine_amount) %>%
     select(infraction_description, location2, set_fine_amoun-
   ## # A tibble: 55 × 3
   ##
                 infraction_description
                                                   location2 se
   ##
                                  <fctr>
                                                      <fctr>
   ## 1
         STND ONSTRT ACCESSIBLE NO PRMT
                                              150 SUDBURY ST
   ## 2
         PARK ONSTRT ACCESSIBLE NO PRMT
                                           143 GORE VALE AVE
   ## 3
         STND ONSTRT ACCESSIBLE NO PRMT
                                                  101 FLM ST
   ## 4
         STND ONSTRT ACCESSIBLE NO PRMT
                                              410 COLLEGE ST
   ## 5
         STND ONSTRT ACCESSIBLE NO PRMT
                                              410 COLLEGE ST
```

6 PARK IN ACCESSIBLE NO PERMIT 3401 DUFFERIN ST ## 7 PARK IN ACCESSIBLE NO PERMIT 2020 SHEPPARD AVE W ## 8 PARK IN ACCESSIBLE NO PERMIT 90 EGITNTON AVE. F. STND ONSTRT ACCESSIBLE NO PRMT 410 COLLEGE ST ## 10 STND ONSTRY ACCESSIBLE NOT D/O 8 SIIITAN ST

Mutate: make new variables

```
library(stringr)
parking %>%
  select(2:8, -infraction code) %>%
 mutate(lower location = str to lower(location2))
## # A tibble: 234,983 × 7
##
     date_of_infraction
                                infraction_description_set
##
                  <int>
                                                <fctr>
## 1
               20141123 PARK FAIL TO DISPLAY RECEIPT
## 2
               20141123 PARK FAIL TO DISPLAY RECEIPT
               20141123 PARK FAIL TO DISPLAY RECEIPT
## 3
```

... with 234,973 more rows, and 4 more variables:

20141123

20141123 PARK FAIL TO DISPLAY RECEIPT

20141123 PARK FAIL TO DEP. FEE MACHINE

20141123 PARK-SIGNED HWY-PROHIBIT DY/TM

20141123 PARK FATL TO DEP. FEE MACHINE

20141123 PARK-SIGNED HWY-PROHIBIT DY/TM

20141123 PARK FAIL TO DISPLAY RECEIPT

PARK ON PRIVATE PROPERTY

4

5

6

7

8

9

10

rank things

```
parking %>%
  select(tag_number_masked, set_fine_amount) %>%
  mutate(rank = min_rank(-set_fine_amount)) %>%
  arrange(rank)
## # A tibble: 234,983 × 3
##
      tag_number_masked set_fine_amount
                                          rank
##
                 <fctr>
                                   <int> <int>
## 1
               ***54560
                                     450
## 2
               ***49346
                                     450
## 3
               ***61327
                                     450
## 4
               ***30223
                                     450
## 5
               ***13649
                                     450
               ***55315
                                     450
## 6
## 7
               ***49375
                                     450
## 8
               ***88513
                                     450
## 9
               ***03453
                                     450
                                     450
## 10
               ***53656
```

Group By / Summary

##

```
parking %>%
  summarise(total fee = sum(set fine amount))
## # A tibble: 1 × 1
## total_fee
##
         <int>
## 1 10885905
parking %>%
  group by(province) %>%
  summarise(total_fee = sum(set_fine_amount)) %>%
  arrange(desc(total fee)) %>%
  mutate(rank = min rank(-total fee))
## # A tibble: 65 \times 3
```

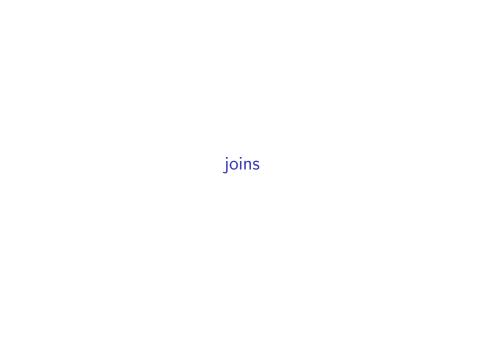
<fctr> <int> <int> <int> 10300590 1

province total_fee rank

more summary:

```
parking %>%
  group_by(province) %>%
  summarise(mean_amt = mean(set_fine_amount),
    max_amt = max(set_fine_amount),
    n = n(),
    sd = sd(set_fine_amount)) %>%
  arrange(desc(mean_amt))
```

```
## # A tibble: 65 × 5
    ##
                                    sd
##
      <fctr> <dbl> <int> <int> <dbl>
         TX 63.80368
                      450
                           163 78.17129
## 1
## 2
         VT 61.81818
                      150 11 47.92039
         OR 57.17391
                      450 46 68.74795
## 3
## 4
         NC 56.15385 450
                            39 72.60620
## 5
         AB 55.82815 450 1769 55.38691
## 6
         IN 55.07812
                      450
                            64 62.43404
## 7
         I.A 54 00000
                      150
                            10 40 33196
```



inner joins

##

A tibble: 552 × 4

<chr>

```
prov_names <- data.frame(province = c('AB', 'BC', 'MB', 'ND')</pre>
  'NT', 'NU', 'ON', 'PE', 'QC', 'SK', 'YT'),
  long_name = c('Alberta', 'British Columbia', 'Manitoba',
  'Nova Scotia', 'Northwest Territories', 'Nunavut', 'Onta:
    'Quebec', 'Saskatchewan', 'Yukon'))
small_parking <- parking %>%
  select(province, location2, set fine amount)
inner join(small parking, prov names) %>%
  filter(province %in% c('NS', 'NY'))
## Warning in inner_join_impl(x, y, by$x, by$y, suffix$x, a
## factors with different levels, coercing to character ve-
```

province location2 set fine amount long name

<fctr>

<fctr

<int>

left/right joins

3

4

5

6

7

8

9

NY

NS

NS

NY

```
left_join(small_parking, prov_names) %>%
  filter(province %in% c('NS', 'NY'))
## Warning in left_join_impl(x, y, by$x, by$y, suffix$x, st
```

```
## factors with different levels, coercing to character ve-
## # A tibble: 1,381 × 4
```

ong_nam	fine_amount lo	location2	province	##
<fctr< td=""><td><int></int></td><td><fctr></fctr></td><td><chr></chr></td><td>##</td></fctr<>	<int></int>	<fctr></fctr>	<chr></chr>	##
N.	30	444 FRONT ST W	L NY	## :
a Scotia	40 Nova	82 LINDYLOU RD	NS	## 2

##		<chr></chr>	<fctr></fctr>	<int></int>	<fctr< th=""></fctr<>
##	1	NY	444 FRONT ST W	30	N
##	2	NS	82 LINDYLOU RD	40 Nova	Scoti

BOGERT AVE

125 GEORGE ST

269 MUTUAL ST

30

30

30

30

30 Nova Scotia

30 Nova Scotia

40 Nova Scotia

N

N

399 SPADINA AVE

NY 92 YORKVILLE AVE

NY 228 DANFORTH AVE

NS 26 BROOKFIELD ST

-			-	
<fct< th=""><th><int></int></th><th><fctr></fctr></th><th><chr></chr></th><th>##</th></fct<>	<int></int>	<fctr></fctr>	<chr></chr>	##
	30	444 FRONT ST W	NY	## 1
Nova Scot	40	82 LINDYLOU RD	NS	## 2

anti joins

anti_join(small_parking, prov_names)

```
## # A tibble: 4,815 × 3
##
     province location2 set_fine_amount
       <fctr>
##
                       <fctr>
                                       <int>
## 1
          AR 1090 DON MILLS RD
                                         250
## 2
           AR.
             2305 QUEEN ST E
                                          30
## 3
          AR.
             19 CHESTER AVE
                                          40
## 4
    AR.
                    HARBORD ST
                                          30
    AR 1963 QUEEN ST E
                                          30
## 5
    AR 1963 QUEEN ST E
## 6
                                          30
## 7
          AR 420 DANFORTH AVE
                                          30
## 8
          AR
                 84 SAULTER ST
                                          40
                                          40
## 9
          AR 79 TORBARRIE RD
## 10
           AR
               1947 QUEEN ST E
                                          30
## # ... with 4.805 more rows
```