

Multi-Page dplyr

Andy Grogan-Kaylor

2019-10-24

Contents

1	Background	1
2	Simulated Data	1
3	Piping	2
4	Aggregate Data: group_by() & summarise()	2
5	Select A Subset of Variables: select()	2
6	Filter A Subset of Rows: filter()	2
7	Create New Variables: mutate()	3
8	Recode Variables: mutate()	3
8.1	Continuous Into Categorical: mutate() & cut()	3
8.2	Categorical Into Categorical: mutate() & recode()	3
9	Rename Variables: rename()	4
10	Drop Missing Values: filter()	4
11	Random Sample	4
12	Connecting To Other Packages Like ggplot	5

1 Background

dplyr is a very powerful R library for managing and processing data.¹

While dplyr is very powerful, learning to use dplyr can be very confusing. This guide aims to present some of the most common dplyr functions and commands in the form of a brief cheatsheet.

¹ The origins of the name dplyr seem somewhat obscure, but I sometimes think of this package as the *data plyers*.

```
library(dplyr)
```

2 Simulated Data

year	x	y	z
2016	NA	Group A	103
2018	30.17	Group C	110.9
2017	38.84	Group B	105.9

year	x	y	z
2017	50.45	Group B	100.4
2015	59.04	Group A	104.1

3 Piping

Pipes `%>%` connect pieces of a command e.g. *data* to *data wrangling* to a *graph command*.

4 Aggregate Data: `group_by()` & `summarise()`

```
mynewdata <- mydata %>%
  group_by(year) %>% # group by y
  summarise(mean_x = mean(x), # mean of x
            n = n()) # count up
```

year	mean_x	n
2015	59.04	1
2016	NA	1
2017	44.65	2
2018	30.17	1

5 Select A Subset of Variables: `select()`

```
mynewdata <- mydata %>% select(x, y) # select only x and y
```

x	y
NA	Group A
30.17	Group C
38.84	Group B
50.45	Group B
59.04	Group A

6 Filter A Subset of Rows: `filter()`

```
mynewdata <- mydata %>% filter(year > 2010) # filter on year
```



```
.default = "Other"))
```

year	x	y	z	yrecode
2016	NA	Group A	103	Red Group
2018	30.17	Group C	110.9	Other
2017	38.84	Group B	105.9	Blue Group
2017	50.45	Group B	100.4	Blue Group
2015	59.04	Group A	104.1	Red Group

9 Rename Variables: `rename()`

```
newdata <- mydata %>%
  rename(age = x, # rename
         mental_health = z)
```

year	age	y	mental_health
2016	NA	Group A	103
2018	30.17	Group C	110.9
2017	38.84	Group B	105.9
2017	50.45	Group B	100.4
2015	59.04	Group A	104.1

10 Drop Missing Values: `filter()`

```
newdata <- mydata %>% filter(!is.na(x)) # filter by x is not missing
```

year	x	y	z
2018	30.17	Group C	110.9
2017	38.84	Group B	105.9
2017	50.45	Group B	100.4
2015	59.04	Group A	104.1

11 Random Sample

```
newdata <- mydata %>% sample_frac(0.5) # fraction of data to sample
```

year	x	y	z
2018	30.17	Group C	110.9

year	x	y	z
2015	59.04	Group A	104.1

12 Connecting To Other Packages Like ggplot

Notice how, in the code below, I never actually create the new data set mynewdata. I simply pipe mydata into a dplyr command, and pipe the result directly to ggplot2.

```
library(ggplot2)

mydata %>% # my data
  mutate(myscale = x + z) %>% # dplyr command to make new variable
  ggplot(aes(x = year, # the rest is ggplot
             y = myscale)) +
  geom_point() + # points
  geom_smooth(se = FALSE) + # smoother without confidence interval
  labs(title = "My Scale By Year") + # labels
  theme(axis.text.x = element_text(size = 10, # tweak theme
                                     angle = 90))
```

