

Four Page dplyr

Andy Grogan-Kaylor

2018-10-18

Contents

1	Background	1
2	Sample Data	1
3	Piping	2
4	Select A Subset of Variables: <code>select()</code>	2
5	Filter A Subset of Rows: <code>filter()</code>	2
6	Create New Variables: <code>mutate()</code>	2
7	Recode Variables: <code>mutate()</code>	3
8	Rename Variables: <code>rename()</code>	3
9	Drop Missing Values: <code>filter()</code>	3
10	Connecting To Other Packages Like <code>ggplot</code>	4

1 Background

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

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.

```
library(dplyr)
```

2 Sample Data

year	x	y	z
2007	NA	Group B	87.25
2012	44.94	Group A	89.56
2017	46.34	Group A	78.24

year	x	y	z
2001	31.18	Group A	100.3
2010	36.65	Group A	81.47

3 Piping

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

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

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

x	y
NA	Group B
44.94	Group A
46.34	Group A
31.18	Group A
36.65	Group A

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

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

year	x	y	z
2012	44.94	Group A	89.56
2017	46.34	Group A	78.24

6 Create New Variables: `mutate()`

```
mynewdata <- mydata %>%
  mutate(myscale = x + z) # create a new variable e.g. a scale
```

year	x	y	z	myscale
2007	NA	Group B	87.25	NA
2012	44.94	Group A	89.56	134.5

year	x	y	z	myscale
2017	46.34	Group A	78.24	124.6
2001	31.18	Group A	100.3	131.5
2010	36.65	Group A	81.47	118.1

7 Recode Variables: mutate()

```
mynewdata <- mydata %>%
  mutate(zcategorical = cut(z, # use mutate to recode
                             breaks=c(-Inf, 100, Inf),
                             labels = c("low", "high")))
```

year	x	y	z	zcategorical
2007	NA	Group B	87.25	low
2012	44.94	Group A	89.56	low
2017	46.34	Group A	78.24	low
2001	31.18	Group A	100.3	high
2010	36.65	Group A	81.47	low

8 Rename Variables: rename()

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

year	age	y	mental_health
2007	NA	Group B	87.25
2012	44.94	Group A	89.56
2017	46.34	Group A	78.24
2001	31.18	Group A	100.3
2010	36.65	Group A	81.47

9 Drop Missing Values: filter()

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

year	x	y	z
2012	44.94	Group A	89.56
2017	46.34	Group A	78.24
2001	31.18	Group A	100.3
2010	36.65	Group A	81.47

10 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
  scale_x_continuous(breaks = scales::pretty_breaks(n = 10)) +
  theme(axis.text.x = element_text(size = 10, # theme() to tweak theme
                                     angle = 90))
```

