

Four Page dplyr

Andy Grogan-Kaylor

2019-01-01

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
7.1	Continuous Into Categorical: <code>mutate()</code> & <code>cut()</code>	3
7.2	Categorical Into Categorical: <code>mutate()</code> & <code>recode()</code>	3
8	Rename Variables: <code>rename()</code>	3
9	Drop Missing Values: <code>filter()</code>	4
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	109
2014	35.33	Group B	92.26
2017	48.21	Group B	110
2012	37.84	Group B	111.2
2015	41.37	Group B	86.58

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
35.33	Group B
48.21	Group B
37.84	Group B
41.37	Group B

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

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

year	x	y	z
2014	35.33	Group B	92.26
2017	48.21	Group B	110
2012	37.84	Group B	111.2
2015	41.37	Group B	86.58

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	109	NA
2014	35.33	Group B	92.26	127.6
2017	48.21	Group B	110	158.2
2012	37.84	Group B	111.2	149.1
2015	41.37	Group B	86.58	127.9

7 Recode Variables: mutate()

7.1 Continuous Into Categorical: mutate() & cut()

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

year	x	y	z	zcategorical
2007	NA	Group B	109	high
2014	35.33	Group B	92.26	low
2017	48.21	Group B	110	high
2012	37.84	Group B	111.2	high
2015	41.37	Group B	86.58	low

7.2 Categorical Into Categorical: mutate() & recode()

```
mynewdata <- mydata %>%
  mutate(yrecoded = dplyr::recode(y, # recode values
                                   "Group A" = "Red Group",
                                   "Group B" = "Blue Group"))
```

year	x	y	z	yrecoded
2007	NA	Group B	109	Blue Group
2014	35.33	Group B	92.26	Blue Group
2017	48.21	Group B	110	Blue Group
2012	37.84	Group B	111.2	Blue Group
2015	41.37	Group B	86.58	Blue Group

8 Rename Variables: rename()

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

year	age	y	mental_health
2007	NA	Group B	109
2014	35.33	Group B	92.26
2017	48.21	Group B	110

year	age	y	mental_health
2012	37.84	Group B	111.2
2015	41.37	Group B	86.58

9 Drop Missing Values: `filter()`

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

year	x	y	z
2014	35.33	Group B	92.26
2017	48.21	Group B	110
2012	37.84	Group B	111.2
2015	41.37	Group B	86.58

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
  theme(axis.text.x = element_text(size = 10, # tweak theme
                                    angle = 90))
```

