# Four Page dplyr

Andy Grogan-Kaylor

2019-02-05

## Contents

## 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 |
|------|-------|---------|-------|
| 2005 | NA | Group B | 89.48 |
| 2009 | 54.51 | Group A | 110.3 |
| 2010 | 41.54 | Group B | 94.54 |
| 2015 | 38.86 | Group A | 106.7 |
| 2007 | 46.99 | Group B | 89.39 |

## 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 |
| 54.51 | Group A |
| 41.54 | Group B |
| 38.86 | Group A |
| 46.99 | Group B |

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

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

| year | x | y | z |
|---|---|---|---|
| 2015 | 38.86 | Group A | 106.7 |

## 6   Create New Variables: `mutate()`

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

| year | x | y | z | myscale |
|---|---|---|---|---|
| 2005 | NA | Group B | 89.48 | NA |
| 2009 | 54.51 | Group A | 110.3 | 164.8 |
| 2010 | 41.54 | Group B | 94.54 | 136.1 |
| 2015 | 38.86 | Group A | 106.7 | 145.6 |
| 2007 | 46.99 | Group B | 89.39 | 136.4 |

## 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 |
|------|------|---------|-------|------------|
| 2005 | NA | Group B | 89.48 | low |
| 2009 | 54.51 | Group A | 110.3 | high |
| 2010 | 41.54 | Group B | 94.54 | low |
| 2015 | 38.86 | Group A | 106.7 | high |
| 2007 | 46.99 | Group B | 89.39 | 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 |
|------|------|---------|-------|-------------|
| 2005 | NA | Group B | 89.48 | Blue Group |
| 2009 | 54.51 | Group A | 110.3 | Red Group |
| 2010 | 41.54 | Group B | 94.54 | Blue Group |
| 2015 | 38.86 | Group A | 106.7 | Red Group |
| 2007 | 46.99 | Group B | 89.39 | Blue Group |

## 8    Rename Variables: `rename()`

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

| year | age | y | mental_health |
|------|-------|---------|-------------|
| 2005 | NA | Group B | 89.48 |
| 2009 | 54.51 | Group A | 110.3 |
| 2010 | 41.54 | Group B | 94.54 |

| year | age | y | mental_health |
|------|------|---------|---------------|
| 2015 | 38.86 | Group A | 106.7 |
| 2007 | 46.99 | Group B | 89.39 |

## 9   Drop Missing Values: `filter()`

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

| year | x | y | z |
|------|-------|---------|-------|
| 2009 | 54.51 | Group A | 110.3 |
| 2010 | 41.54 | Group B | 94.54 |
| 2015 | 38.86 | Group A | 106.7 |
| 2007 | 46.99 | Group B | 89.39 |

## 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))
```