# Data Manipulation and Visualization in R

Oliwier Przewlocki

Technologisches Gewerbemuseum

oprzewlocki@student.tgm.ac.at

# Contents

# 1. Introduction

We continue from the "Introduction to the R-Programming Language" and look at ways of optimizing the syntax to perform already familiar data operations and discover ways of visualizing the data.

# 2. Data Manipulation

This chapter covers ways to simplify already discussedways of manipulating data using tools like **Dplyr** (manipulating) and **Tidyr** (cleaning)

## 2.1. Dplyr

First, an installation using `install.packages('dplyr')` and an activation using `library(dplyr)` is required.

To allow for an easy way of showcasing the power of dplyr, the `nycflights13` package needs to be installed aswell using the same commands.

### 2.1.1. filter() & slice()

Allows for subset row selection. It's easier than the built-in subset function we've been using. To filter, we just pass in the data as the first argument and then any number of arguments that specify each column condition:

```
1    filter(flights, month==11, day==3, carrier=='AA')
```

The slice function allows for positional selection:

```
1    slice(flights, 1:10) # outputs the first 10 rows
```

### 2.1.2. arrange()

Allows for ordering and has a very similar structure to the `filter()` function.

```
1    arrange(flights, year, month, desc(day)) # first order by year, then month and
  then by day but in descending order.
```

### 2.1.3. select() & rename()

Selects a defined set of columns:

```
1    select(flights, carrier, day) # selects the carrier and the day columns.
```

And `rename()` quickly renames the columns (newname = oldname)

```
1    rename(flights, newdayname = day)
```

### 2.1.4. distinct()

Selects distinct row values. Powerful to use in combinations such as with `select()`

```
1    distinct(select(flights, carrier)) # selects the distinct carriers.
```

### 2.1.5. mutate() & transmutate()

Creates new columns as functions of existing columns

```
1    mutate(flights, new_col = arr_delay-dep_delay)
```

If you only want the newly created columns, and not the entire data frame, use `transmutate()`

```
1    transmutate(flights, new_col = arr_delay-dep_delay)
```

### 2.1.6. summarise()

Summarizes a column using an aggregate function

```
1    summarise(flights, avg_air_time=mean(air_time, na.rm = TRUE ))
```

### 2.1.7. sample_n() & sample_frac()

Sample_n() returns a random set of n-samples out of a data frame

```
1    sample_n(flights, 10) # 10 random sample rows
```

Sample_frac() returns a percentage of data. The range goes from 0 (0%) to 1 (100%)

```
1    sample_frac(flights, 0.3) # returns 30% of data.
```

## 2.2. The Pipe Operator

It allows better readability. An alternative to the pipe operator is nesting, which is barely readable beyond a certain point. You can also try to create multiple variables and assign them one by one, but then you're wasting memory space. The pipe operator allows you to chain an output of a function to the input of another function in a readable and clear way.

```
1    # nesting
2    result <- arrange(sample_n(filter(df, mpg>20),size=5),desc(mpg))
3
4    # multiple assignments
5    a <- filter(df, mpg>20)
6    b <- sample_n(a, size=5)
7    result <- arrange(b, desc(mpg))
8
9    # pipe operator
10   result <- df %>% filter(mpg>20) %>% sample_n(size=5) %>% arrange(desc(mpg))
```

## 2.3. Tidyr

First perform the installation & setup using `install.packages('tidyr')` and `library(tidyr)`. There's also a complimentary package that needs to be installed called `data.table` using the same commands. This package is very similar to the built-in data frames, although it increases computation speed substantially.

### 2.3.1. gather() & spread()

A prerequisite to the `gather()` and `spread()` functions is the knowledge of the wide and long formats. Long formats have repeating first columns, and wide formats have values that don't repeat in the first column.

The gather function converts a table into a long format and the spread function converts a table into a wide format. That's it.

`gather()` collapses a slice of columns into a key-value pair (the key is the column name and the value is the row data value.):

```
1    gather(df, key, value, Qtr1:Qtr4)
```

`spread()` widens the key-value pair passed into it to individual tables. The unique keys become columns and the values become the row values.

```
1    spread(df, key, value)
```

### 2.3.2. separate()

Separates a single column into multiple columns. As an example, there is a column that has data like "a-x" and it's required of us to separate "a" and "x" into their own respective columns

```
1    separate(data=df, col= col.name, into c("abc", "xyz"), sep = '-')
```

### 2.3.3. unite()

Pastes multiple columns into one.

```
1    unite(separated.df, new.joined.col, abc, xyz)
2
3    unite(separated.df, new.joined.col, abc, xyz, sep = "---")
```

# 3. Data Visualization

This chapter focuses on visualizing data using the **ggplot2** library.

## 3.1. Overview of ggplot2

The library is built on layers, namely the "**Data**" (the raw data), the "**Aesthetics**" (specify the columns and features you want to display) and "**Geometries**" (the type of plot) as the main layers.

Then there are other layers, such as "**Facets**" (multiple plots), "**Statistics**", "**Coordinates**" and "**Themes**"

There exists a cheat sheet for `ggplot2` that explains the major aspects in around 2 pages. It can be found <u>here</u>.

## 3.2. Histograms

This section requires both the 'ggplot2' and 'ggplot2movies' libraries to be installed and activated.

```
1    # DATA & AESTHETICS
2    pl <- ggplot(movies, aes(x=rating)) # rating is a column in the movies data
   frame
3
4    # GEOMETRY
5    pl2 <- pl + geom_histogram(binwidth = 0.2, color='red', fill='pink',
   alpha=0.8)
6    print(pl2) # this actually displays the plot
7
8    pl3 <- pl2 + xlab('Movie Rating') + ylab("Count") # x & y labels
9    print(pl3 + ggtitle("MY TITLE"))
10
11
12   # CUSTOM COLOR GRADIENT
13   pl2 <- pl + geom_histogram(binwidth = 0.2, aes(fill=..count..)) # the higher
   the count the more blue it is.
```

## 3.3. Scatterplots

```
1    df <- mtcars
2
3    # DATA & AESTHETICS
4    ggplot(df, aes(x=wt,y=mpg))
5
6    # GEOMETRY
7    print(pl + geom_point(size=hp)) # size is based on the 'hp' column
8
9    # FACTOR
```

```
10   pl + geom_point(aes(size=factor(cyl))) # factor labels the 'cyl' column as
     categorical (not continuous) because it can only have either 4, 6 or 8
     cylinders, not 5 or 7.
11
12   # SHAPE
13   pl + geom_point(aes(shape=factor(cyl))) # this assigns a different shape for
     different cylinder sizes, like a square or a triangle.
14
15   # CUSTOM COLOR GRADIENT
16   pl <- pl + geom_point(aes(color=hp), size=5)
17   pl + scale_color_gradient(low='blue', high='red')
18
19   # ATTENTION: You can only use column values to define shapes, colors, etc.
     inside of the aes() block.
```

## 3.4. Barplots

```
1   df <- mpg
2
3   pl <- ggplot(df, aes(x=class))
4
5   print(pl + geom_bar(aes(fill=drv), position="dodge")) # The custom fill
    automatically creates a stacked bar plot. The position argument set to dodge
    makes the stacked bars appear next to eachother. There's also "fill" that makes
    it a normalized area graph.
```

## 3.5. Boxplots

```
1   df <- mtcars
2
3   ggplot(df, aes(x=factor(cyl), y=mpg))
4   print(pl + geom_boxplot(aes(fill=factor(cyl))) + coord_flip()) # the
    coord_flip() flips the coordinate axes. It essentially rotates the plot -90°.
```

## 3.6. 2 Variable Plotting

```
1   pl <- ggplot(movies, aes(x=year, y=rating))
2   print(pl + geom_bind2d() + scale_fill_gradient(high="red", low="green")) #
    creates a 2d-bin-chart. The bins change color based on their count. (Essentially
    a heatmap)
```

There's also a `hexbin`, which creates a heatmap where each bin is a hexagon. For that, install the `hexbin` package.

```
1
2   # HEXBIN
3   pl <- ggplot(movies, aes(x=year, y=rating))
4   pl2 <- pl + geom_hex()
5   print(pl2)
6
7   # DENSITY PLOT
8   pl2 <- pl + geom_density2d()
9   print(pl2)
```

## 3.7. Coordinates Faceting

`ggplot2` allows for the changing of coordinate systems and the adjustment of their paramters like x-lims and y-lims:

```
1    pl <- ggplot(mpg, aes(x=displ, y=hwy)) + geom_point()
2
3    # STANDARD COORDINATES
4    pl2 <- pl + coord_cartesian(xlim = c(1,4), ylim = c(15,30))
5
6    # FIXED RATIO COORDINATES
7    pl2 <- pl + coord_fixed(ratio = 1/3) # Default is 1:1.
8
9    # Consult the cheat sheet for more coordinate types.
10
11   print(pl2)
```

Faceting allows for the placement of multiple plots next to eachother.

```
1    pl <- ggplot(mpg, aes(x=displ, y=hwy)) + geom_point()
2
3    print(pl + facet_grid(. ~ cyl)) # we seperate the mpg plot (on the x-axis)
     into 3 subplots each seperated by the cylinder column. So there's one plot for
     cars with 4 cylinder, one for cars with 6 cylinder, etc.
4
5    print(pl + facet_grid(drv ~ .)) # this separates the plot along the y-axis.
6
7    # The general syntax is "what you want to faucet by the y-axis, tilde symbol,
     and then what you want to facet by on the x-axis." A dot means "everything
     else".
8
9    # So the x and y separations can be mixed:
10   print(pl + facet_grid(drv ~ cyl)) # this separates the plot along the y-axis.
11
```

## 3.8. Theming

You can set a theme by either setting it globaly using the `theme_set()` function or by just adding it
when printing:

```
1    theme_set(theme_minimal()) # Setting the theme globally
2
3    print(pl + theme_dark()) # Assigning the theme individually
```

To get even more themes, install the `ggthemes` library.

## 3.9. Data Visualization Assignment Solution

```
1       df <- fread("/home/ren/coding/r/R-Course-HTML-Notes/R-Course-HTML-Notes/R-
   for-Data-Science-and-Machine-Learning/Training Exercises/Capstone and Data Viz
   Projects/Data Visualization Project/Economist_Assignment_Data.csv", drop=1)
2  head(df)
3
4    pl <- ggplot(df, aes(x=CPI, y=HDI)) + geom_point(shape=1, size = 3,
   aes(color=Region))
5
6    pl2 <- pl + geom_smooth(aes(group=1), method = lm, formula = y ~ log(x), se =
   FALSE, color="red")
7
8    pointsToLabel <- c("Russia", "Venezuela", "Iraq", "Myanmar", "Sudan",
9                       "Afghanistan", "Congo", "Greece", "Argentina", "Brazil",
10                      "India", "Italy", "China", "South Africa", "Spane",
11                      "Botswana", "Cape Verde", "Bhutan", "Rwanda", "France",
12                      "United States", "Germany", "Britain", "Barbados",
   "Norway", "Japan",
13                      "New Zealand", "Singapore")
14
15   pl3 <- pl2 + geom_text(aes(label = Country), color = "gray20",
16                      data = subset(df, Country %in%
   pointsToLabel),check_overlap = TRUE)
17
18   print(pl3 + theme_bw() + scale_x_continuous(name = "Corruption Perception
   Index (CPI)", limits = c(1, 10), breaks = 1:10 ) + scale_y_continuous(name =
   "Human Development Index, 2011 (1=Best)", limits = c(0, 1), breaks = seq(0, 1,
   0.1)) + ggtitle("Corruption and Human development"))
```

## 3.10. Interactive Visualization with Plotly

Plotly allows for the interactive visualization of data in multiple languages like R, Python or Matlab. It's entirely open source, free and self-hosted.

A recent addition to their package is the ability to convert `ggplot2` plots into interactive plotly plots directly inside of the R environment.

To get started, install `plotly` activate it.

```
1   pl <- ggplot(mtcars, aes(mpg, wt)) + geom_point()
2
3   gpl <- ggplotly(pl) # now you can zoom in and get more informationwithin the
  plot.
```

Their ggplot2 plotly documentation has alot of great tutorials on how to create various types of plots that can be accessed here.

## 3.11. Moneyball Assignment Solution

```r
  batting <- read.csv('/home/ren/coding/r/R-Course-HTML-Notes/R-Course-HTML-
Notes/R-for-Data-Science-and-Machine-Learning/Training Exercises/Capstone and
Data Viz Projects/Capstone Project/Batting.csv')

  head(batting)

  str(batting)

  batting$BA <- batting$H / batting$AB

  tail(batting$BA, 5)

  batting$OBP <- (batting$H + batting$BB + batting$HBP) / (batting$AB +
batting$BB + batting$HBP + batting$SF)

  batting$SLG <- ((batting$H - batting$X2B - batting$X3B - batting$HR) + (2 *
batting$X2B) + (3 * batting$X3B) + (4 * batting$HR)) / batting$AB

  str(batting)

  sal <- read.csv("/home/ren/coding/r/R-Course-HTML-Notes/R-Course-HTML-Notes/R-
for-Data-Science-and-Machine-Learning/Training Exercises/Capstone and Data Viz
Projects/Capstone Project/Salaries.csv")

  summary(batting)
  batting <- subset(batting, yearID >= 1985)
  summary(batting)

  merged <- merge(batting, sal, by = c("playerID", "yearID"))

  summary(merged)

# option 1
  lost_players <- filter(merged, playerID == "giambja01" | playerID ==
"damonjo01" | playerID == "saenzol01")

# option 2
  lost_players <- subset(merged, playerID %in% c("giambja01", "damonjo01",
"saenzol01"))

  distinct(lost_players, playerID)
  count(lost_players, playerID)

  lost_players <- subset(lost_players, yearID == 2001)
  lost_players <- select(lost_players, H, X2B, X3B, HR, OBP, SLG, BA, AB)

  print(lost_players)

  sum(lost_players$AB) # 1469
  mean(lost_players$OBP) # 0.363

  replacement_players <- filter(merged, yearID == 2001, salary <= 5000000)
  replacement_players

  pl <- ggplot(replacement_players, aes(x = AB, y = OBP, name = playerID, salary
= salary))
  pl <- pl + geom_point()
```

```r
50
51  ply <- ggplotly(pl, tooltip = c("x", "y", "name", "salary"))
52  ply
53
54  # I'd choose "berkmla01", "heltoto01" and "gonzalu01" as replacements for the
    lost players
```