# Getting and Cleaning Data

Daniel Reynaga Gil                 201504086

```
# Append two more function calls to accomplish the following:
#
# 1. Use group_by() (from dplyr) to group the data by part and
# sex, in that order.
#
# 2. Use mutate to add two new columns, whose values will be
# automatically computed group-by-group.
#
#    * total = sum(count)
#    * prop = count / total
#
sat %>%
        select(-contains("total")) %>%
        gather(part_sex, count, -score_range) %>%
        separate(part_sex, c("part", "sex")) %>%
        group_by(part, sex) %>%
        mutate(total = sum(count),
               prop = count / total
        ) %>% print
```

Console

```
...

  |=============================================================| 100%

| Would you like to receive credit for completing this course on Coursera.org?

1: No
2: Yes

Selection:
```

**Environment**

| | |
|---|---|
| failed | 6 obs. of 4 variables |
| galton | 928 obs. of 2 variables |
| gradebook | 10 obs. of 4 variables |
| mydf | 225468 obs. of 11 variables |
| pack_sum | 6023 obs. of 5 variables |
| passed | 4 obs. of 4 variables |
| res | 20 obs. of 3 variables |
| result1 | 46 obs. of 5 variables |

R: Spread a key-value pair across multiple columns.

spread {tidyr}          R Documentation

## Spread a key-value pair across multiple columns.

### Description

Spread a key-value pair across multiple columns.

### Usage

```
spread(data, key, value, fill = NA, convert = FALSE, sep = NULL)
```

### Arguments

---

```
structure(1478960666.06322,
class = c("POSIXct",
"POSIXt"), tzone =
"Asia/Hong_Kong")
```

Console

```
| savings time, the length of any given minute, day, month, week, or year is relative to when it occurs. In
| contrast, the length of a second is always the same, regardless of when it occurs.

...

  |===========================================================| 97%

| To address these complexities, the authors of lubridate introduce four classes of time related objects:
| instants, intervals, durations, and periods. These topics are beyond the scope of this lesson, but you can
| find a complete discussion in the 2011 Journal of Statistical Software paper titled 'Dates and Times Made
| Easy with lubridate'.

...

  |============================================================| 98%

| This concludes our introduction to working with dates and times in lubridate. I created a little timer that
| started running in the background when you began this lesson. Type stopwatch() to see how long you've been
| working!

> stopwatch()
[1] "35M 49.0137670040131S"

| Excellent job!


  |=============================================================| 100%

| Would you like to receive credit for completing this course on Coursera.org?

1: No
2: Yes

Selection:
```

**Environment**

**Values**

| | |
|---|---|
| arrive | 2016-11-12 22:24:26 |
| depart | 2016-11-11 17:34:26 |
| dt1 | "2014-08-23 17:23:02" |
| dt2 | chr [1:3] "2014-05-14" "2014-0... |
| how_long | Formal class Interval |
| last_time | |
| my_date | 1989-05-17 |

R: Utilities for creation and manipulation of 'Interval'...

interval {lubridate}          R Documentation

## Utilities for creation and manipulation of Interval objects.

### Description

interval creates an Interval-class object with the specified start and end dates. If the start date occurs before the end date, the interval will be positive. Otherwise, it will be negative.

%--% Creates an interval that covers the range spanned by two dates. It replaces the original behavior of lubridate, which created an interval by default whenever two date-times were