

# stringr solutions

*Jumping Rivers*

## Question 1

We'll start by loading the necessary packages and data sets

```
library("tidyverse")
data(names, package = "jrTidyverse2")
```

Here we have a data set containing 800 people with the names: "Abigail", "Alexander", "Ava", "Benjamin", "Charlotte", "Emily", "Emma", "Ethan", "Harper", "Isabella", "Jacob", "James", "Liam", "Mason", "Mia", "Michael", "Noah", "Olivia", "Sophia" and "William".

Using various functions from **stringr** and `count()` from **dplyr**, work out the frequency of each name. Which name occurs the most?

```
names %>%
  mutate(name = str_trim(name)) %>%
  mutate(name = str_to_title(name)) %>%
  count(name) %>%
  arrange(n)
## Warning: Detecting old grouped_df format, replacing `vars` attribute by
## `groups`
## # A tibble: 20 x 3
## # Groups:   sex [2]
##   sex  name      n
##   <chr> <chr>    <int>
## 1 F    Sophia      8
## 2 F    Harper     13
## 3 M    Mason      15
## 4 M    Benjamin   18
## 5 M    Noah       19
## 6 F    Emma       20
## 7 F    Emily      28
## 8 M    Alexander  30
## 9 M    Michael    34
## 10 F   Ava       38
## 11 F   Charlotte  42
## 12 M   James     42
## 13 M   Jacob     45
## 14 F   Abigail   47
## 15 M   Liam     58
## 16 F   Isabella  65
## 17 M   Ethan    68
## 18 F   Olivia   69
## 19 F   Mia      70
## 20 M   William  71
```

## Question 2

We'll start off by loading the data

```
data(beer, package = "jrTidyverse2")
```

Let's inspect the data

```
head(beer)
## # A tibble: 6 x 3
##   URL                                     ABV Color
##   <chr>                                <dbl> <dbl>
## 1 /homebrew/recipe/view/61925/the_devil_is_in_the_details_duve~ 8.4   3.11
## 2 /homebrew/recipe/view/110195/fat_head_s_headhunter_ipa_clone 7.52  7.3
## 3 /homebrew/recipe/view/244064/kronenbourg_1664_blanc_klone    5.36  3.26
## 4 /homebrew/recipe/view/213882/delirium_tremens_clone         8.97  5.64
## 5 /homebrew/recipe/view/392676/bakke-brygg-juleale-2016-20-l    6.54  20.9
## 6 /homebrew/recipe/view/106240/hunahpu_clone_                 12.4   40
```

Here we have a data set of beers with their alcohol percentage and colour. Colour is ranked from 1-50 with 1 being pale and 50 being black. The only problem is that the names of the beers have been scraped from the internet and so are contained in an url. To do any analysis on the beers we are going to need to extract the names. The names of the beers are always after the last forward slash in the url. For example, the first url, /homebrew/recipe/view/61925/the-devil-is-in-the-details-duvel-clone-

would become

The Devil Is In The Details Duvel Clone

It's going to be a bit easier to extract the vector of urls, work with it that way, then reattach it once we are done.

```
url = beer$URL
```

- 1) Extract the last part of the url.  
Hint: Your regex should start with \ and should end with a \$.

```
url = url%>%
  str_extract("/[a-zA-Z-_0-9]*$")
```

- 2) Going with the first example, your beer name should now look like  
/the\_devil\_is\_in\_the\_details\_duvel\_clone\_  
Get rid of the forward slash.

```
url = url %>%
  str_replace("/", "")
```

- 3) The beer names should now look like  
the\_devil\_is\_in\_the\_details\_duvel\_clone\_  
Replace the underscores with spaces. Careful, some of the urls have dashes instead of underscores inbetween words.  
Hint: Use a group, (), in your regex

```
url = url %>%
  str_replace_all("(\\-|\\_)", " ")
```

- 4) The beer names will now look like  
the devil is in the details duvel clone

Trim the surrounding whitespace and give all words capital letters. Once that is complete, overwrite the urls with the extracted names within the data.

```
url = url %>%
  str_to_title() %>%
```

```
str_trim()
beer$URL = url
```

5) We want to do some analysis on the beers based on whether they are an IPA, stout or pale ale. To do this we're going to introduce a new function called `if_else()` from **dplyr**. For example

```
df = data.frame(x = c(2,4,6,8))
```

Here we have made a data frame called `df` containing a column of numbers called `x`.

```
df = df %>%
  mutate(y = if_else(condition = x > 5, true = 1, false = 0))
df
##   x y
## 1 2 0
## 2 4 0
## 3 6 1
## 4 8 1
```

In this step we are mutating a new column called `y` that will be the value 1 when `x > 5` and the value 0 otherwise. We can do the same for the beers. Notice that if we run the code

```
str_detect(beer$URL, "Ale")
```

We get a TRUE when the beer name contains `Ale` and FALSE otherwise. We can use this inside `if_else()` as a condition

```
beer = beer %>%
  mutate(
    Ale = if_else(str_detect(URL, "Ale"), 1, 0)
  )
```

So here we would be creating a column called `Ale`, that is 1 when the beer name contains `Ale` and 0 otherwise. Create a column called `Ipa` that is 1 when the name contains `Ipa` and 0 otherwise. Do the same for `Stout`.

```
beer = beer %>%
  mutate(Ale = if_else(str_detect(URL, "Ale"), 1, 0),
         Ipa = if_else(str_detect(URL, "Ipa"), 1, 0),
         Stout = if_else(str_detect(URL, "Stout"), 1, 0)
  )
```

6) Under the principles of tidy data, this is no longer tidy, we should have one column containing whether the beer is an "Ale", "Ipa" or "Stout". We can do this using `gather()` from **tidyr**

```
beer = beer %>%
  gather(Type, Yes, Ale, Ipa, Stout) %>%
  filter(Yes != 0) %>%
  select(-Yes)
```

Here we are gathering the Ale, Ipa and Stout columns into two columns called Type and Yes. We're not interested in the beers with a value of 0 so we filter them out. Then we delete the Yes column using `select()`. Which type of beer has the highest average alcohol percentage and color? Hint: Use **dplyr**

7) Plot the data using **ggplot2**, assigning a different colour to each type of beer

```
ggplot(beer, aes(x = ABV, y = Color)) +
  geom_point(aes(colour = Type))
```