

# Lab 6

Sharleen Price spp2122

February 22nd, 2018

## Instructions

Make sure that you upload an RMarkdown file to the canvas page (this should have a .Rmd extension) as well as the PDF of HTML output after you have knitted the file (this will have a .pdf extension or .html). Note that since you have already knitted this file, you should see both a **Lab6\_UNI.pdf** and a **Lab6\_UNI.Rmd** file in your UN2102 folder. Click on the **Files** tab to the right to see this. The files you upload to the Canvas page should be updated with commands you provide to answer each of the questions below. You can edit this file directly to produce your final solutions. The lab is due 11:59pm on Tuesday, March 27th.

## Titanic

In this lab we will be studying a data set which provides information on the survival rates of passengers on the fatal voyage of the ocean liner *Titanic*. The dataset provides information on each passenger including, for example, economic status, sex, age, cabin, name, and survival status. This is a training dataset taken from the Kaggle competition website; for more information on Kaggle competitions, please refer to <https://www.kaggle.com>. Students should download the data set on Canvas.

## Tasks

- 1) Load the **Titanic** data set.

```
# Read in data
# setwd("~/Desktop/Data")
titanic <- read.table("Titanic.txt", header = TRUE, as.is = TRUE)
```

- 2) Look at the first 10 entries of the variable **Name**. Notice that each person has a *title*, i.e., Mr., Mrs. Miss., etc...

```
## R Code -----
head(titanic$Name, 10)

## [1] "Braund, Mr. Owen Harris"
## [2] "Cumings, Mrs. John Bradley (Florence Briggs Thayer)"
## [3] "Heikkinen, Miss. Laina"
## [4] "Futrelle, Mrs. Jacques Heath (Lily May Peel)"
## [5] "Allen, Mr. William Henry"
## [6] "Moran, Mr. James"
## [7] "McCarthy, Mr. Timothy J"
## [8] "Palsson, Master. Gosta Leonard"
## [9] "Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)"
## [10] "Nasser, Mrs. Nicholas (Adele Achem)"
```

- 3) Create a new variable of the **titanic** dataframe called **Title** that gives the appropriate title of each passenger. The variable **Title** should have 5 levels: **Miss**, **Mrs**, **Mr**, **Master**, and **Other**. Display the first 10 entries of the new variable **Title**.

```
## R Code -----
#head(titanic$Name)
name_split <- strsplit(titanic$Name, split = ",")
#length(name_split)
split_comma = c()
for (i in 1:length(name_split)){
  split_comma<-append(split_comma,name_split[[i]][2])
}
var<-strsplit(split_comma, split=" ")
Title = c()
for (i in 1:length(split_comma)){
  Title<-append(Title,var[[i]][2])
}

for (j in 1:length(Title)){
  if (Title[j]!="Mr." && Title[j]!="Mrs." && Title[j]!="Miss."&& Title[j]!="Master."){
    Title[j] = "Other"
  }
}

head(Title, 10)
```

```
## [1] "Mr."      "Mrs."      "Miss."     "Mrs."      "Mr."      "Mr."      "Mr."
## [8] "Master." "Mrs."      "Mrs."
```

```
titanic$TitleAll <- Title
```

- 4) Create a table showing the counts for each level of the variable **Title**. Also create a table showing the number of passengers that survived split by their *title*.

```
## R Code -----
Title_count <- table(Title)
Title_count <- sort(Title_count, decreasing = TRUE)
Title_count
```

```
## Title
##      Mr.   Miss.   Mrs. Master.   Other
##      517    182    125     40     27
```

```
Survivor <- titanic[titanic$Survived==1,]
```

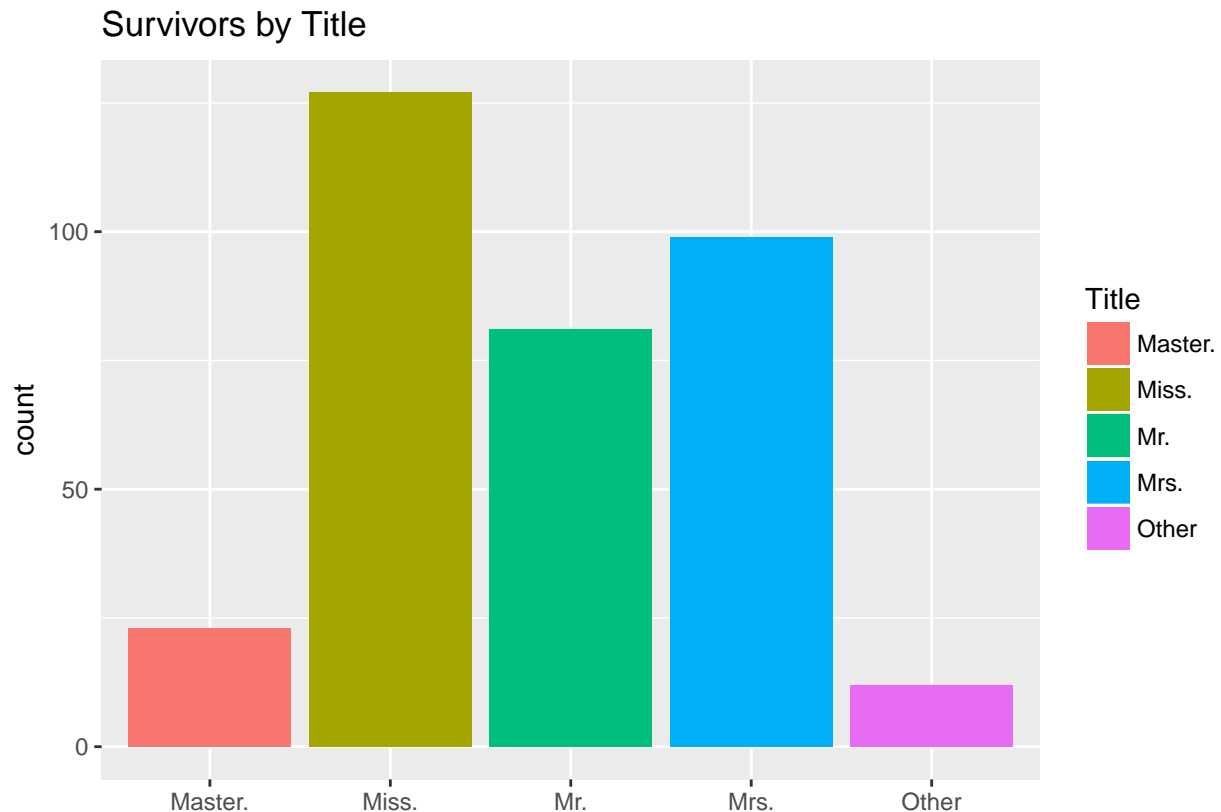
```
table(Survivor$TitleAll)
```

```
##
## Master.   Miss.    Mr.    Mrs.   Other
##         23     127     81     99     12
```

- 5) Plot the number of passengers that survived split by the variable **Title**. Use **ggplot** with the geometric object **geom\_bar**.

```
## R Code -----
library(ggplot2)
ggplot(data=Survivor) +
```

```
geom_bar(aes(x=TitleAll,fill=(TitleAll)))+
labs(title = "Survivors by Title",fill="Title",x="")
```



- 6) Display all of the names that correspond to **Other**. How many cases fall in the this category and what are the name titles corresponding to the level **Other**? Note: you can just identify the names by inspection.

```
## R Code -----
others <-titanic[Title=="Other",]
length(others$Name)
```

```
## [1] 27
```

```
others$Name
```

```
## [1] "Uruchurtu, Don. Manuel E"
## [2] "Byles, Rev. Thomas Roussel Davids"
## [3] "Bateman, Rev. Robert James"
## [4] "Minahan, Dr. William Edward"
## [5] "Carter, Rev. Ernest Courtenay"
## [6] "Moraweck, Dr. Ernest"
## [7] "Aubart, Mme. Leontine Pauline"
## [8] "Pain, Dr. Alfred"
## [9] "Reynaldo, Ms. Encarnacion"
## [10] "Peuchen, Major. Arthur Godfrey"
## [11] "Butt, Major. Archibald Willingham"
## [12] "Duff Gordon, Lady. (Lucille Christiana Sutherland) (\\"Mrs Morgan\\")"
## [13] "Duff Gordon, Sir. Cosmo Edmund (\\"Mr Morgan\\")"
## [14] "Kirkland, Rev. Charles Leonard"
## [15] "Stahelin-Maeglin, Dr. Max"
```

```
## [16] "Sagesser, Mlle. Emma"
## [17] "Simonius-Blumer, Col. Oberst Alfons"
## [18] "Frauenthal, Dr. Henry William"
## [19] "Weir, Col. John"
## [20] "Mayne, Mlle. Berthe Antonine (\"Mrs de Villiers\")"
## [21] "Crosby, Capt. Edward Gifford"
## [22] "Roths, the Countess. of (Lucy Noel Martha Dyer-Edwards)"
## [23] "Brewer, Dr. Arthur Jackson"
## [24] "Leader, Dr. Alice (Farnham)"
## [25] "Reuchlin, Jonkheer. John George"
## [26] "Harper, Rev. John"
## [27] "Montvila, Rev. Juozas"
```

27 names fall into the category "Other"

- 7) Create a new variable of the titanic data frame called **Last\_name** that gives the last name of passenger. Display the first 10 entries of the new variable **Last\_name**.

```
## R Code -----

first.element <- function(vec){
  return(vec[1])
}
Last_Name<-sapply(strsplit(titanic$Name,split = ", "),first.element)
head(Last_Name, 10)
```

```
## [1] "Braund"      "Cumings"     "Heikkinen"  "Futrelle"   "Allen"
## [6] "Moran"       "McCarthy"    "Palsson"    "Johnson"    "Nasser"
```

- 8) Display the first 8 most common last names.

```
## R Code -----

lastname_table<-table(Last_Name)
lastname_table<-sort(lastname_table, decreasing = TRUE)[1:8]
lastname_table
```

```
## Last_Name
## Andersson      Sage      Carter      Goodwin      Johnson      Panula      Skoog
##           9           7           6           6           6           6
##           Rice
##           5
```