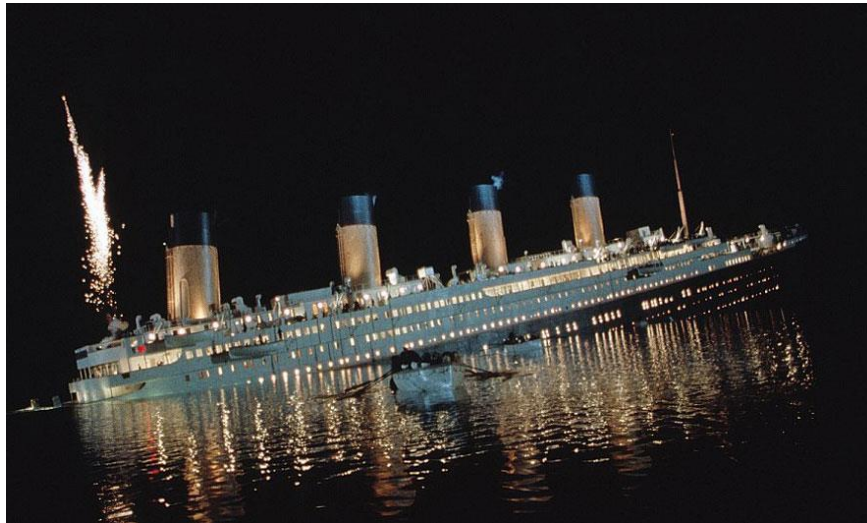


# Big Data and Cloud Computing

## Lab 1 – RStudio and Introduction to R

\* The dataset in this assignment is adopted from the famous Kaggle competition *Titanic: Machine Learning from Disaster*.



*The sinking of the RMS Titanic is one of the most famous shipwrecks in history. On April 15, 1912, during her maiden voyage, the Titanic sank after colliding with an iceberg, killing 1502 out of 2224 passengers and crew. This sensational tragedy shocked the international community and led to better safety regulations for ships.*

*One of the reasons that the shipwreck led to such loss of life was that there were not enough lifeboats for the passengers and crew. Although there was some element of luck involved in surviving the sinking, some groups of people were more likely to survive than others.*

*In this lab, we will analyze this disaster, get more insights into it and come up with conclusions about which sorts of people were more likely to survive than others.*

## Dataset:

Titanic dataset is the most famous dataset for beginners in Data Science. Let's have a view on the data dictionary.

★ **Tip 1: It's very important to spend some time understanding data. Take a closer look at the data available, assess and explore it using tables and graphics.**

Variable	Definition
<b>survival</b>	Survival (0= No, 1=Yes).
<b>pClass</b>	Passenger Class (1=1 <sup>st</sup> , 2=2 <sup>nd</sup> , 3=3 <sup>rd</sup> ).
<b>name</b>	Name.
<b>gender</b>	Gender.
<b>age</b>	Age in years.
<b>sibsp</b>	Number of siblings/spouses aboard the Titanic.
<b>parch</b>	Number of parents/children aboard the Titanic.
<b>ticket</b>	Ticket Number.
<b>fare</b>	Passenger Fare.
<b>cabin</b>	Cabin.
<b>embarked</b>	Port of embarkation (C= Cherbourg, Q=Queenstown, S= Southampton).

## Variable notes:

**pclass:** A proxy for socio-economic status (SES)  
1st = Upper, 2nd = Middle, 3rd = Lower

**age:** Age is fractional if less than 1. If the age is estimated, is it in the form of xx.5

**sibsp:** The dataset defines family relations in the following way:  
*Sibling* = brother, sister, stepbrother, stepsister  
*Spouse* = husband, wife (mistresses and fiancés were ignored)

**parch:** The dataset defines family relations in the following way:  
*Parent* = mother, father  
*Child* = daughter, son, stepdaughter, stepson  
\*some children travelled only with a nanny, therefore parch=0 for them.

## Requirements:

★ **Tip 2: Before asking others for help, it's generally a good idea for you to try to help yourself. R includes extensive facilities for accessing documentation and searching for help.**

1.	First of all, start by cleaning the workspace and setting the working directory.
2.	Import the dataset <b>titanic.csv</b> into a data frame.
3.	It's time to explore the dataset as a whole.
a.	Show the dimensions of the data frame. Hint: Use dim()
b.	Show the structure of the data frame. Hint: Use str()
c.	Get more insight into data by exploring the <b>first</b> and the <b>last TEN</b> rows in the dataset.
d.	Show summary of all variables in the data frame.

4.	Let's explore some variables in the dataset.
a.	Show a summary for the variable <i>age</i> <u>only</u> .
b.	What are the first and third quartile values for this variable? What do these values mean?
c.	Are there any missing values in the variable <i>age</i> ? (i.e. written as <NA>) <b>✦ Hint: Read the documentation for <i>is.na()</i> and <i>anyNA()</i> to find out how to know if a certain variable has missing values. What are the differences between them? Which one is better to use in this case?</b>
d.	What is the type of the variable <i>embarked</i> ? Show the levels of this variable. Is that what you were expecting?
e.	Can you conclude what's needed at this step in the data analysis cycle?
5.	As you probably might have answered in (4.e), preprocessing is needed. Data preprocessing is a very important step in any data analytics project.
a.	Remove the rows containing <NA> in the <i>age</i> variable from the data frame.
b.	Remove the rows containing any unexpected value in the <i>embarked</i> variable from the dataset.
c.	Now, check that no NA values exist in the <i>age</i> variable. Also, factor the <i>embarked</i> variable and display its levels. Is that what you are expecting?
d.	Some variables are not very interesting and provide no real indicative value. Remove columns <i>Cabin</i> and <i>Ticket</i> from the dataset.
6.	An important step also in any data analysis project is <b><u>statistical description</u></b> and <b><u>visualization</u></b> . We will now visualize some variables, and try to get insights out of them.  In this part, you will practice visualization and slicing/indexing in data frames.
a.	Show the number of <b>males</b> and <b>females</b> aboard the Titanic.
b.	Plot a pie chart showing the number of males and females aboard the Titanic. (Hint: use <b><i>pie()</i></b> function).
c.	Indicate males with a blue color and females with a red color in the above plot. (Hint: There is a color parameter in any plot function).
c.	Show the number of people who <b><u>survived</u></b> and <b><u>didn't survive</u></b> from each gender.
d.	Plot a pie chart showing the number of males and females who survived <b><u>only</u></b> .
e.	What do you conclude from that?
f.	Show the relationship between social class and survival i.e. show how many people survived and how many people didn't survive from each class.
g.	Plot this relationship as a <b>stacked bar plot</b> . (Hint: use <b><i>barplot()</i></b> function)
h.	Indicate survived passengers with a blue color and un-survived passengers with a red color in the above plot.
i.	What do you conclude from that?
j.	Plot a <b>box and whiskers</b> plot for the variable <i>age</i> (Hint: use <b><i>boxplot()</i></b> function) Read about the <b>box and whiskers</b> plot and understand it properly.
k.	What does this plot mean?
l.	Plot a density distribution for the variable <i>age</i> .
7.	Remove all columns but passenger name and whether they survived or not. Export the new dataset to a file named <b>"titanic_preprocessed.csv"</b>