# SESSION 6: Visualization & Plotting

# Assignment 1

Problem Statement

1. Import the Titanic Dataset from the following link:

https://drive.google.com/file/d/1JTJCjdGuUxzKXYlwOavwovB01k6FWg3r/view?ts=5b42ea10

Perform the below operations:

a. Pre-process the passenger names to come up with a list of titles that represent families and represent using appropriate visualization graph.

```
# Grab title from passenger names
full$Title <- gsub('(.*, )|(\\..*)', '', full$Name)


# Titles with very low cell counts to be combined to "rare" level
rare_title <- c('Dona', 'Lady', 'the Countess','Capt', 'Col', 'Don',
                'Dr', 'Major', 'Rev', 'Sir', 'Jonkheer')
# Also reassign mlle, ms, and mme accordingly
full$Title[full$Title == 'Mlle']        <- 'Miss'
full$Title[full$Title == 'Ms']          <- 'Miss'
full$Title[full$Title == 'Mme']         <- 'Mrs'
full$Title[full$Title %in% rare_title]  <- 'Rare Title'
# Finally, grab surname from passenger name
full$Surname <- sapply(full$Name,
                function(x) strsplit(x, split = '[,.]')[[1]][1])
```

b. Represent the proportion of people survived by family size using a graph.

```r
ggplot(full[1:891,], aes(x = Fsize, fill = factor(Survived))) +
  geom_bar(stat='count', position='dodge') +
  scale_x_continuous(breaks=c(1:11)) +
  labs(x = 'Family Size') +
  theme_few()
```

c. Impute the missing values in Age variable using Mice library, create two different graphs showing Age distribution before and after imputation

```r
par(mfrow=c(1,2))
hist(full$Age, freq=F, main='Age: Original Data',
  col='darkgreen', ylim=c(0,0.04))
hist(mice_output$Age, freq=F, main='Age: MICE Output',
  col='lightgreen', ylim=c(0,0.04))
```