

## 1

- readfile rg\_train.csv to R. use stringsAsFactors=F ( make it a practice )
- break column age\_band into two columns a1 and a2( use - as separator with separate function from package tidyr)
- Replace value 71+ in a1 with 71
- convert a1 and a2 to numeric types
- create a column a3 which takes values using following logic . if a1=71 then a3=71, else a3=(a1+a2)/2
- Replace missing values in a3 by average of a3
- remove columns a1, a2 from data

## 2

- from family\_income columns remove ‘,’ and ‘<’ symbol
- separate family income into two columns f1 and f2( use >= as separator)
- convert f1 and f2 into numeric type
- when f1 is 4000, put 4000 in f2
- when f2 is 35000 put 35000 in f1
- create column f3 which is (f1+f2)/2
- replace missing values in f3 with its mean
- remove columns f1,f2 from the data

## 3

- load data hflights from package hflights
- Find out for each month count of flights which each unique carrier had
- Find out which carrier had the highest number of flights for each month
- Find out , considering the entire year which plane had the highest number of flights for each carrier ( planes are identified by tailnum)

## 4

- Find out the difference between airtime and actual elapsed time for each observation ( time spent on land)
- Find out the average land time for each month
- Find out how the average land time changed month over month ( make use of difference from lag)

## 5

use following code to create data

```
set.seed(2)
library(dplyr)
data=data.frame(week=sample(1:4,40,replace = T),
                 letr=sample(c("a","b","c","d","e","f","g"),40,replace = T))
data=data %>%
  arrange(week) %>%
```

```
distinct(week,letr, .keep_all= TRUE) %>%  
mutate(id=1:n())
```

create a new column called output . it takes value “new” if the correponding letr for the observation did not appear in the immediate previous week, otherwise it takes value “repeated”. Here are few initial steps to help you

- a. sort the data by (letr,week)
- b. group by letr
- c. create a new column which takes value week-lag(week)
- d. use this new column and any other information to create the output column as stated above.