# Shao Yun Gao–Week 2: R Output

## Step 1: Read in the Data

```R
#! Step 1: Read in the Data

# List the structure of the data (str)
hmeq <- read.csv("../HMEQ_WK02/HMEQ_Loss.csv")

# Execute a summary of the data
str(hmeq)

# Execute a summary of the data
summary(hmeq)

# Print the first six records
head(hmeq)
```

```
$ REASON    : chr  "HomeImp" "HomeImp" "HomeImp" "" ...
$ JOB       : chr  "Other" "Other" "Other" "" ...
$ YOJ       : num  10.5 7 4 NA 3 9 5 11 3 16 ...
$ DEROG     : int  0 0 0 NA 0 0 3 0 0 0 ...
$ DELINQ    : int  0 2 0 NA 0 0 2 0 2 0 ...
$ CLAGE     : num  94.4 121.8 149.5 NA 93.3 ...
$ NINQ      : int  1 0 1 NA 0 1 1 0 1 0 ...
$ CLNO      : int  9 14 10 NA 14 8 17 8 12 13 ...
$ DEBTINC   : num  NA NA NA NA NA ...
>
> # Execute a summary of the data
> summary(hmeq)
 TARGET_BAD_FLAG TARGET_LOSS_AMT     LOAN          MORTDUE          VALUE
 Min.   :0.0000  Min.   :  224   Min.   : 1100   Min.   :  2063   Min.   :  8000
 1st Qu.:0.0000  1st Qu.: 5639   1st Qu.:11100   1st Qu.: 46276   1st Qu.: 66076
 Median :0.0000  Median :11003   Median :16300   Median : 65019   Median : 89236
 Mean   :0.1995  Mean   :13415   Mean   :18608   Mean   : 73761   Mean   :101776
 3rd Qu.:0.0000  3rd Qu.:17634   3rd Qu.:23300   3rd Qu.: 91488   3rd Qu.:119824
 Max.   :1.0000  Max.   :78987   Max.   :89900   Max.   :399550   Max.   :855909
                 NA's   :4771                    NA's   :518      NA's   :112
    REASON             JOB                 YOJ             DEROG            DELINQ
 Length:5960        Length:5960        Min.   : 0.000   Min.   : 0.0000   Min.   : 0.0000
 Class :character   Class :character   1st Qu.: 3.000   1st Qu.: 0.0000   1st Qu.: 0.0000
 Mode  :character   Mode  :character   Median : 7.000   Median : 0.0000   Median : 0.0000
                                       Mean   : 8.922   Mean   : 0.2546   Mean   : 0.4494
                                       3rd Qu.:13.000   3rd Qu.: 0.0000   3rd Qu.: 0.0000
                                       Max.   :41.000   Max.   :10.0000   Max.   :15.0000
                                       NA's   :515      NA's   :708       NA's   :580
     CLAGE             NINQ             CLNO            DEBTINC
 Min.   :   0.0   Min.   : 0.000   Min.   : 0.0   Min.   :  0.5245
 1st Qu.: 115.1   1st Qu.: 0.000   1st Qu.:15.0   1st Qu.: 29.1400
 Median : 173.5   Median : 1.000   Median :20.0   Median : 34.8183
 Mean   : 179.8   Mean   : 1.186   Mean   :21.3   Mean   : 33.7799
 3rd Qu.: 231.6   3rd Qu.: 2.000   3rd Qu.:26.0   3rd Qu.: 39.0031
 Max.   :1168.2   Max.   :17.000   Max.   :71.0   Max.   :203.3121
 NA's   :308      NA's   :510      NA's   :222    NA's   :1267
>
> # Print the first six records
> head(hmeq)
  TARGET_BAD_FLAG TARGET_LOSS_AMT LOAN MORTDUE  VALUE REASON     JOB  YOJ DEROG DELINQ
1              1             641 1100   25860  39025 HomeImp  Other 10.5     0      0
2              1            1109 1300   70053  68400 HomeImp  Other  7.0     0      2
3              1             767 1500   13500  16700 HomeImp  Other  4.0     0      0
4              1            1425 1500      NA     NA            NA   NA    NA     NA
5              0              NA 1700   97800 112000 HomeImp Office  3.0     0      0
6              1             335 1700   30548  40320 HomeImp  Other  9.0     0      0
     CLAGE NINQ CLNO  DEBTINC
1  94.36667    1    9       NA
2 121.83333    0   14       NA
3 149.46667    1   10       NA
4        NA   NA   NA       NA
5  93.33333    0   14       NA
6 101.46600    1    8 37.11361
> |
```
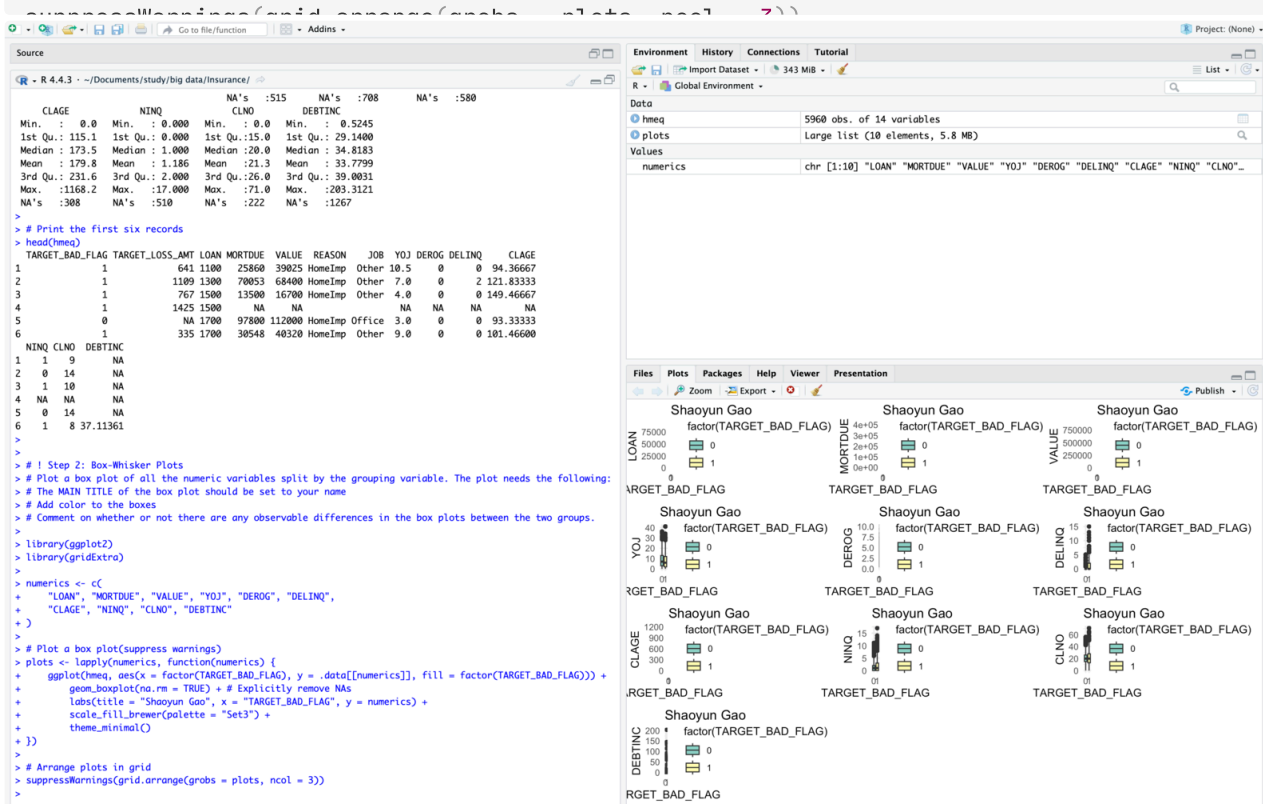
# Step 2: Box–Whisker Plots

```R
# ! Step 2: Box-Whisker Plots
# Plot a box plot of all the numeric variables split by the grouping variable. The plo
t needs the following:
# The MAIN TITLE of the box plot should be set to your name
# Add color to the boxes
# Comment on whether or not there are any observable differences in the box plots betw
een the two groups.

library(ggplot2)
library(gridExtra)

numerics <- c(
  "LOAN", "MORTDUE", "VALUE", "YOJ", "DEROG", "DELINQ",
  "CLAGE", "NINQ", "CLNO", "DEBTINC"
)

# Plot a box plot(suppress warnings)
plots <- lapply(numerics, function(numerics) {
  ggplot(hmeq, aes(x = factor(TARGET_BAD_FLAG), y = .data[[numerics]], fill = factor(T
ARGET_BAD_FLAG))) +
    geom_boxplot(na.rm = TRUE) + # Explicitly remove NAs
    labs(title = "Shaoyun Gao", x = "TARGET_BAD_FLAG", y = numerics) +
    scale_fill_brewer(palette = "Set3") +
    theme_minimal()
})
```
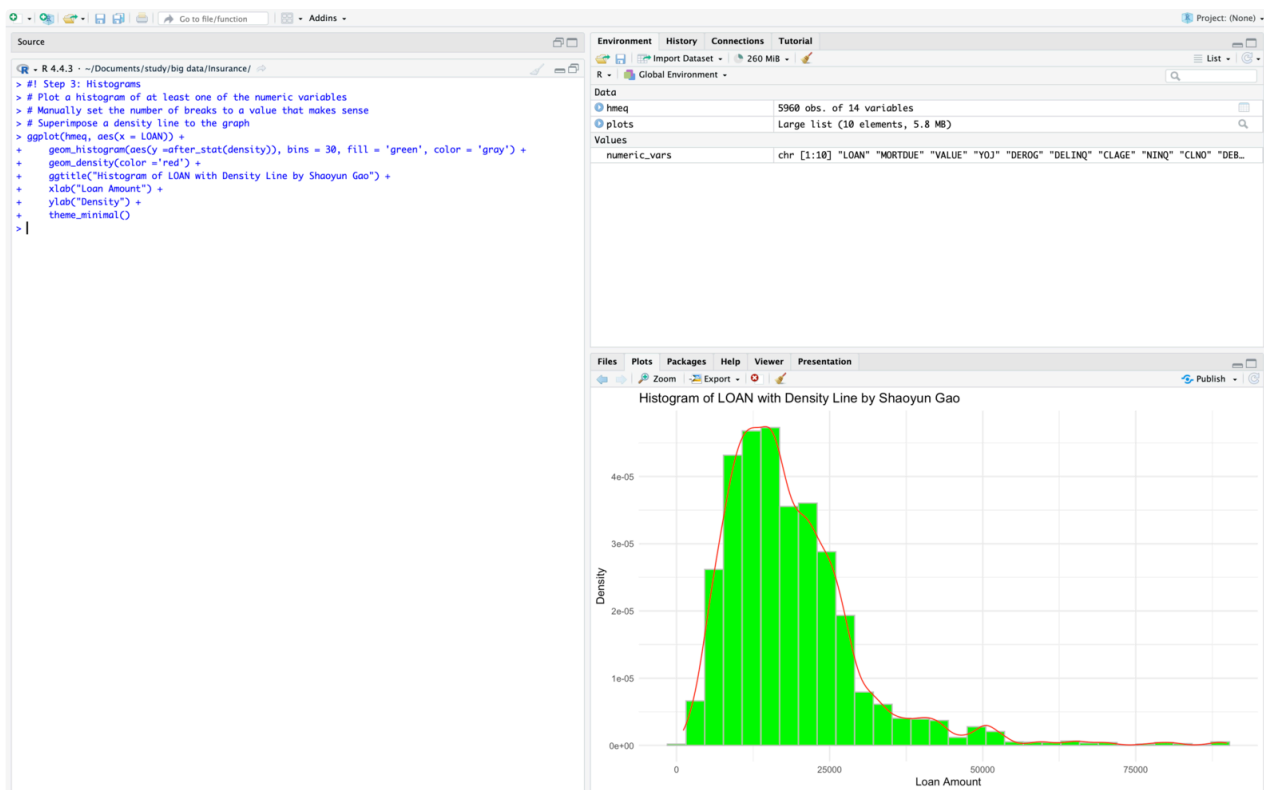
```
# Arrange plots in grid
```



## Step 3: Histograms

```r
#! Step 3: Histograms
# Plot a histogram of at least one of the numeric variables
# Manually set the number of breaks to a value that makes sense
# Superimpose a density line to the graph
ggplot(hmeq, aes(x = LOAN)) +
  geom_histogram(aes(y =after_stat(density)), bins = 30, fill = 'green', color = 'gra
y') +
  geom_density(color ='red') +
  ggtitle("Histogram of LOAN with Density Line by Shaoyun Gao") +
  xlab("Loan Amount") +
  ylab("Density") +
  theme_minimal()
```

# Step 4: Impute "Fix" all the numeric variables that have missing values

```r
                                                                          R
# ! Step 4: Impute "Fix" all the numeric variables that have missing values
# For the missing Target variables, simply set the missing values to zero
# For the remaining numeric variables with missing values, create two new variables. O
ne variable will have a name beginning with IMP_ and it will contained the imputed val
ue. The second value will have a name beginning with M_ and it will contain a 1 if the
record was imputed and a zero if it was not.
# You may impute with any method that makes sense. The median or mean value will be us
eful in most cases.
# Push yourself! Try one complex imputation like the one described in the lectures.
# Delete the original variable after it has been imputed.

impute_missing <- function(data, var) {
  if (any(is.na(data[[var]]))) {
    if (var == "TARGET_LOSS_AMT") {
      data[[var]][is.na(data[[var]])] <- 0
    } else {
      if (var %in% c("INCOME", "HOME_VAL")) {
        job_groups <- split(data, data$JOB)
        for (group in names(job_groups)) {
          median_value <- median(job_groups[[group]][[var]], na.rm = TRUE)
          data[[paste0("IMP_", var)]][data$JOB == group & is.na(data[[var]])] <- media
n_value
        }
      } else {
```

```r
        median_value <- median(data[[var]], na.rm = TRUE)
        data[[paste0("IMP_", var)]][is.na(data[[var]])] <- median_value
      }
      data[[paste0("M_", var)]] <- as.numeric(is.na(data[[var]]))
      data[[var]] <- NULL
    }
  }
  return(data)
}


numerics <- names(hmeq)[sapply(hmeq, is.numeric)]
for (var in numerics) {
  hmeq <- impute_missing(hmeq, var)
}

# Run a summary to prove that all the variables have been imputed
summary(hmeq)

res_vars <- grep("^M_", names(hmeq), value = TRUE)
my_missing_count <- sapply(res_vars, function(var) sum(hmeq[[var]]))
my_missing_count
```



# Step 5: One Hot Encoding

```r
                                                                                R

# ! Step 5: One Hot Encoding
# For the character / category variables, perform one hot encoding. For this create a
Flag for each categories.
# Delete the original class variable
```

```r
# Run a summary to show that the category variables have been replaced by Flag variables.

hot_encode <- function(data, var) {
  if (is.factor(data[[var]]) || is.character(data[[var]])) {
    levels <- levels(as.factor(data[[var]]))
    for (level in levels) {
      data[[paste0("FLAG_", var, "_", level)]] <- as.numeric(data[[var]] == level)
    }
    data[[var]] <- NULL
  }
  return(data)
}


chars <- names(hmeq)[sapply(hmeq, function(x) is.factor(x) || is.character(x))]
for (var in chars) {
  hmeq <- hot_encode(hmeq, var)
}

summary(hmeq)
```