

# HW2\_\_Liang\_\_Dan

*Dan Liang*

*9/5/2018*

## Problem 1

## Problem 2

## Problem 3

complete

## Problem 4

1. Version control can help us handling changes, for example, undo the content we mistakenly deleted.
2. Version control storing member's code online for the whole group of members to have access the code. Once one member change the code or submit new function, the version will be updated and the online system will keep the old version in the repository and have the new one, so all members of course project group can get to know that you update your work and have access to check your work.
3. I can try new code at the same time save the older version. If the new one does not work, just revert it.
4. Through checking the old versions of the code to find out when and where bugs were introduced, so helps to avoid those next time.

## Problem 5

```
mydata5a <- read.csv('~/Desktop/5a.csv')
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

cran1 <- tbl_df(mydata5a)
cran1

## # A tibble: 150 x 3
##   Item Operator Data
##   <int>      <int> <dbl>
## 1     1         1     4.3
## 2     1         1     4.3
## 3     1         1     4.1
## 4     2         1     6
## 5     2         1     4.9
```

```
## 6      2      1      6
## 7      3      1     2.4
## 8      3      1     3.9
## 9      3      1     1.9
## 10     4      1     7.4
## # ... with 140 more rows
```

```
summary(cran1)
```

```
##      Item      Operator      Data
## Min.   : 1.0   Min.   :1   Min.   :0.700
## 1st Qu.: 3.0   1st Qu.:2   1st Qu.:3.025
## Median : 5.5   Median :3   Median :4.700
## Mean   : 5.5   Mean   :3   Mean   :4.657
## 3rd Qu.: 8.0   3rd Qu.:4   3rd Qu.:6.000
## Max.   :10.0   Max.   :5   Max.   :9.400
```

```
mydata5b <- read.csv('~/Desktop/5b.csv')
library(dplyr)
cran2 <- tbl_df(mydata5b)
cran2
```

```
## # A tibble: 22 x 2
##   year long.jump
##   <int>   <dbl>
## 1    -4    250.
## 2     0    283.
## 3     4    289.
## 4     8    294.
## 5    12    299.
## 6    20    282.
## 7    24    293.
## 8    28    305.
## 9    32    301.
## 10   36    317.
## # ... with 12 more rows
```

```
summary(cran2)
```

```
##      year      long.jump
## Min.   :-4.00   Min.   :249.8
## 1st Qu.:21.00   1st Qu.:295.4
## Median :50.00   Median :308.1
## Mean   :45.45   Mean   :310.3
## 3rd Qu.:71.00   3rd Qu.:327.5
## Max.   :92.00   Max.   :350.5
```

```
mydata5c <- read.csv('~/Desktop/5c.csv')
library(dplyr)
cran3 <- tbl_df(mydata5c)
cran3
```

```
## # A tibble: 62 x 2
##   Body.Wt Brain.Wt
##   <dbl>   <dbl>
## 1  3.38    44.5
## 2  0.48    15.5
```

```
## 3 1.35      8.1
## 4 465      423
## 5 36.3     120.
## 6 27.7     115
## 7 14.8     98.2
## 8 1.04      5.5
## 9 4.19      58
## 10 0.425    6.4
## # ... with 52 more rows
```

```
summary(cran3)
```

```
##      Body.Wt      Brain.Wt
## Min.   : 0.005  Min.   : 0.10
## 1st Qu.: 0.600  1st Qu.: 4.25
## Median : 3.342  Median : 17.25
## Mean   : 198.790 Mean   : 283.13
## 3rd Qu.: 48.203  3rd Qu.: 166.00
## Max.   :6654.000 Max.   :5712.00
```

```
mydata5d <- read.csv('~/.Desktop/5d.csv')
library(dplyr)
cran4 <- tbl_df(mydata5d)
cran4
```

```
## # A tibble: 6 x 4
##   X          X10000 X20000 X30000
##   <fct>      <dbl> <dbl> <dbl>
## 1 "Ife\\#1"    16.1  16.6  20.8
## 2 "Ife\\#1"    15.3  19.2   18
## 3 "Ife\\#1"    17.5  18.5   21
## 4 PusaEarlyDwarf 8.1   12.7  14.4
## 5 PusaEarlyDwarf 8.6   13.7  15.4
## 6 PusaEarlyDwarf 10.1  11.5  13.7
```

```
summary(cran4)
```

```
##           X          X10000      X20000      X30000
## Ife\\#1      :3  Min.   : 8.100  Min.   :11.50  Min.   :13.70
## PusaEarlyDwarf:3  1st Qu.: 8.975  1st Qu.:12.95  1st Qu.:14.65
##              Median :12.700  Median :15.15  Median :16.70
##              Mean   :12.617  Mean   :15.37  Mean   :17.22
##              3rd Qu.:15.900  3rd Qu.:18.02  3rd Qu.:20.10
##              Max.   :17.500  Max.   :19.20  Max.   :21.00
```

## Problem 6

```
library(swirl)
```

```
##
## | Hi! I see that you have some variables saved in your workspace. To keep
## | things running smoothly, I recommend you clean up before starting swirl.
##
## | Type ls() to see a list of the variables in your workspace. Then, type
## | rm(list=ls()) to clear your workspace.
```

```
##
## | Type swirl() when you are ready to begin.
# Path to data
.datapath <- file.path(path.package('swirl'), 'Courses', 'R_Programming_E', 'Looking_at_Data','plant-da
# Read in data
plants <- read.csv(.datapath, strip.white=TRUE, na.strings="")
str (plants)

## 'data.frame':   5166 obs. of  12 variables:
## $ Accepted.Symbol      : Factor w/ 5166 levels "ABBA","ABBAB",...: 3 4 5 1 2 7 6 8 15 16 ...
## $ Synonym.Symbol       : logi  NA NA NA NA NA NA NA ...
## $ Scientific.Name      : Factor w/ 5166 levels "Abelmoschus",...: 1 2 3 4 5 6 7 8 9 10 ...
## $ Duration             : Factor w/ 8 levels "Annual","Annual, Biennial",...: NA 4 NA 7 7 NA 1 NA
## $ Active.Growth.Period : Factor w/ 8 levels "Fall, Winter and Spring",...: NA NA NA 4 NA NA NA NA
## $ Foliage.Color        : Factor w/ 6 levels "Dark Green","Gray-Green",...: NA NA NA 3 NA NA NA NA
## $ pH..Minimum.         : num  NA NA NA 4 NA NA NA NA 7 NA ...
## $ pH..Maximum.         : num  NA NA NA 6 NA NA NA NA 8.5 NA ...
## $ Precipitation..Minimum. : int  NA NA NA 13 NA NA NA NA 4 NA ...
## $ Precipitation..Maximum. : int  NA NA NA 60 NA NA NA NA 20 NA ...
## $ Shade.Tolerance      : Factor w/ 3 levels "Intermediate",...: NA NA NA 3 NA NA NA NA 2 NA ...
## $ Temperature..Minimum...F.: int  NA NA NA -43 NA NA NA NA -13 NA ...

# remove NA in ph_Min and ph_Max
plantcleaned<-filter(plants, !is.na(plants$pH..Minimum) & !is.na(plants$pH..Maximum))
# get mean of ph max and ph min
meanofph <- 2/(plantcleaned$pH..Maximum+plantcleaned$pH..Minimum)
# Use function lm to test for a relationship
fit <- lm(meanofph ~ plantcleaned$Foliage.Color)

summary(fit)

##
## Call:
## lm(formula = meanofph ~ plantcleaned$Foliage.Color)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.040933 -0.009038 -0.001509  0.008056  0.062297
##
## Coefficients:
##              Estimate Std. Error t value
## (Intercept)      0.168176   0.001607 104.628
## plantcleaned$Foliage.ColorGray-Green -0.010637   0.003325  -3.199
## plantcleaned$Foliage.ColorGreen      -0.005292   0.001700  -3.113
## plantcleaned$Foliage.ColorRed        -0.003293   0.007453  -0.442
## plantcleaned$Foliage.ColorWhite-Gray -0.011291   0.005111  -2.209
## plantcleaned$Foliage.ColorYellow-Green 0.002086   0.003630   0.575
##              Pr(>|t|)
## (Intercept)      < 2e-16 ***
## plantcleaned$Foliage.ColorGray-Green  0.00143 **
## plantcleaned$Foliage.ColorGreen      0.00192 **
## plantcleaned$Foliage.ColorRed        0.65877
## plantcleaned$Foliage.ColorWhite-Gray  0.02744 *
## plantcleaned$Foliage.ColorYellow-Green 0.56573
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.01456 on 826 degrees of freedom
## (7 observations deleted due to missingness)
## Multiple R-squared:  0.02368,    Adjusted R-squared:  0.01777
## F-statistic: 4.007 on 5 and 826 DF,  p-value: 0.001343
```

## Problem 7

### Problem 7a-d

```
# read data Personal:Personal car details; ODefects: observed defects; DetailsD: Defect Details (dataset)
# Personal <- read.csv('~/Desktop/Personal.csv')
# ODefects <- read.csv('~/Desktop/Defects.csv')
# DetailsD <- read.csv('~/Desktop/Details.csv')

# Merge first two datasets by license
# mergebylicense <- merge(Personal, ODefects, by="Kenteken")
# Merge the defects details in
# mergebycode <- merge(mergebylicense, DetailsD, by="Gebrek.identificatie")

# Remove all NA
# mergecars <- na.omit(mergebycode)

# 5c. count how many different makes and models of cars in 2017
# install.packages('plyr')
# library(plyr)

# get subset 2017
# defectssubset <- subset(mergecars, mergecars$Meld.datum.door.keuringsinstantie > 20170000 & mergecars$Meld.datum.door.keuringsinstantie < 20180000)
# problem 7d, get different makes and models of cars
# length(unique(defectssubset$Merk))
# length(unique(defectssubset$Handelsbenaming))
# There are 137 types of makes in 2017. There are 2938 types of models in 2017
```

### Problem 7e.

```
# summary(defectssubset)
```

1. According to the summary of defects in 2017, top five defects are AC1 205 K04 476 210, and the top five models have the defect are AC1 are PEUGEOT, OPEL, VOLKSWAGEN, CITROEN, VOLVO. The top models have the defect AC1 are GOLF, 207, CORSA, POLP, TOYOTA AYGO;

**Problem 7h.** clean each dataset remove the NA first, then merge would save some times; try each operation in R script before directly knit the R markdown