# Data Analytics Platform Operation Manual

Student

After lecturers subscribe to a case, students will receive a notification email including information on the link to the Data Analytics Platform, a username (email address of the user), and a temporary password. Students could click on the website link http://47.243.52.252/hku-dap-client/#/Signin to enter the Data Analytics Platform. After entering the email addresses and temporary password, students will be notified to change the temporary password to their password (8-20 digits). Log in to the Platform with the user's email address and password.

My Labs



*Screenshot of the user interface of the Data Analytics Platform*



*Screenshot of the Labs that students are going to use*

Instructions for using Rstudio

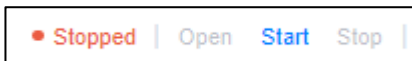*'My Labs'* show all the labs that teachers subscribed to for students.

1. For detailed instructions on using the virtual environments, click 'help'.

   RStudio ⑦ help

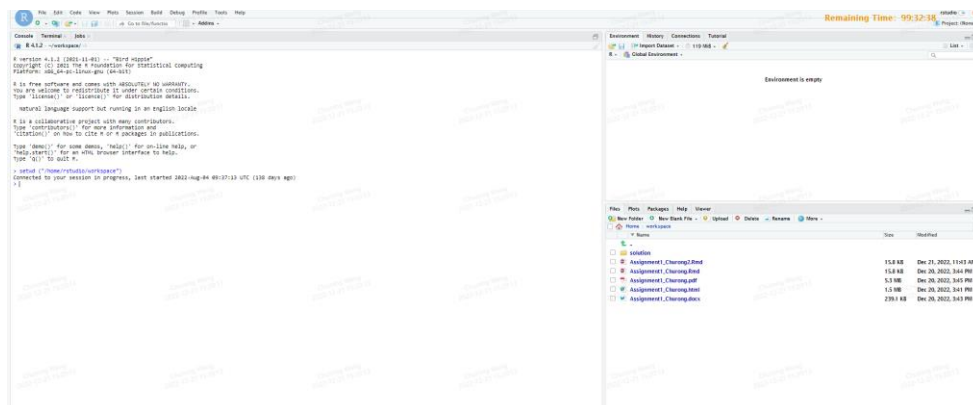2. For instructions on the case, click the upper right icon.

   

3. To start, open or stop the labs, click the buttons at the bottom of each lab. The left bottom icon shows the status of the lab. Wait for around 15 seconds for the lab to load.

   ● Stopped | Open | **Start** | Stop |

   When the status icon turns to "Running," click "Open" to access the virtual environment of Rstudio.

   ● Running

   The system will then automatically jump to the virtual environment of Rstudio's user interface. "Remaining Time" shows the time limit of your lab usage. If you use up the lab time, you will need to email the lecturer or the teaching assistant to renew the lab. The watermark shows the user name and the access time of users.

   

4. Students can view detailed information about each lab by clicking 'Details.'

   Details

   After entering the "Lab Details" page, you can find more info on the lab.

   

# Open a R script file

To create a new R script file, click 'File>New File > R script'.



## Dataset Access

The system creates the dataset folder by default, and all datasets are stored under the dataset folder. The following codes are examples to access the datasets:

```
1.  ## Set working directory
2.  setwd(dir= '/dataset')
3.  dir()
```

```
> ## Set working directory
> setwd(dir= '/dataset')
> dir()
[1] "Centaline_data"  "Centaline_test"  "Centaline_train" "HKG_adm1"
```

You can see there are some datasets stored under the directory "/dataset". Select your desired dataset from the directory:

```
1.  ## load dataset
2.  df = read.csv('/dataset/Centaline_data/Centaline_data.csv', header=TRUE)
3.  head(df, 10)
```

```
> df = read.csv('/dataset/Centaline_data/Centaline_data.csv',header=TRUE)
> head(df,10)
   Transaction_price Transaction_year Transaction_month                      Location
1            2000000             2021                11               6 SHUN PING STREET
2            3000000             2020                 2 200 SHA TAU KOK ROAD SHEK CHUNG AU
3            9300000             2020                 3             31 SHUN LUNG STREET
4            9200000             2020                 3             31 SHUN LUNG STREET
5            6205200             2020                 1             31 SHUN LUNG STREET
6            6170900             2020                 8             31 SHUN LUNG STREET
7            6049800             2020                 9             31 SHUN LUNG STREET
8            5871900             2020                10             31 SHUN LUNG STREET
9            5800000             2020                 1             31 SHUN LUNG STREET
10           5711000             2021                 4             31 SHUN LUNG STREET
                      Estate          HMA Developer Gross_size Saleable_size No_of_rooms Floor      Region
1  Kam Tong Lau Sha Tau Kok        Other        -1           290          -1     1 New Territory
```

# Exploratory Data Analysis (EDA)

Explore the distribution of variables through simple visualizations.

```
1.  ## Distribution of Transaction_Price
2.  with(df, hist(Transaction_price,breaks=100))
```

## Histogram of Transaction_price



# Linear Regression

Linear regression is a statistical model that analyzes the relationship between a response variable (often called y) and one or more variables and their interactions (often called x or explanatory variables). Linear regression can be calculated in R with the command lm(). In the next example, use this command to calculate the property price based on relevant variables. Select a few variables that influence the property price the most. Use transaction_price as the response variable and the other variables as explanatory variables.

The lm command takes the variables in the format:

lm([target] ~ [predictor / features], data = [data source])

Sample code:
```
1.  lm.fit1=lm(Transaction_price ~
2.              Saleable_size +
3.              No_of_rooms +
4.              Floor +
5.              Age_of_property,
6.          data = df)
7.  summary(lm.fit1)
```

Check out the performance of the model with the summary() function.

```
> summary(lm.fit1)

Call:
lm(formula = Transaction_price ~ Saleable_size + No_of_rooms +
    Floor + Age_of_property, data = df)

Residuals:
      Min        1Q    Median        3Q       Max
-110922022  -2108672    191494   2098319 548393207

Coefficients:
                  Estimate Std. Error t value Pr(>|t|)
(Intercept)     -3.938e+06  4.131e+04  -95.33   <2e-16 ***
Saleable_size    2.785e+04  5.348e+01  520.86   <2e-16 ***
No_of_rooms     -5.081e+05  1.041e+04  -48.82   <2e-16 ***
Floor            3.087e+04  1.097e+03   28.16   <2e-16 ***
Age_of_property -7.545e+04  8.835e+02  -85.40   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6383000 on 216642 degrees of freedom
Multiple R-squared:   0.59,      Adjusted R-squared:   0.59
F-statistic: 7.794e+04 on 4 and 216642 DF,  p-value: < 2.2e-16
```
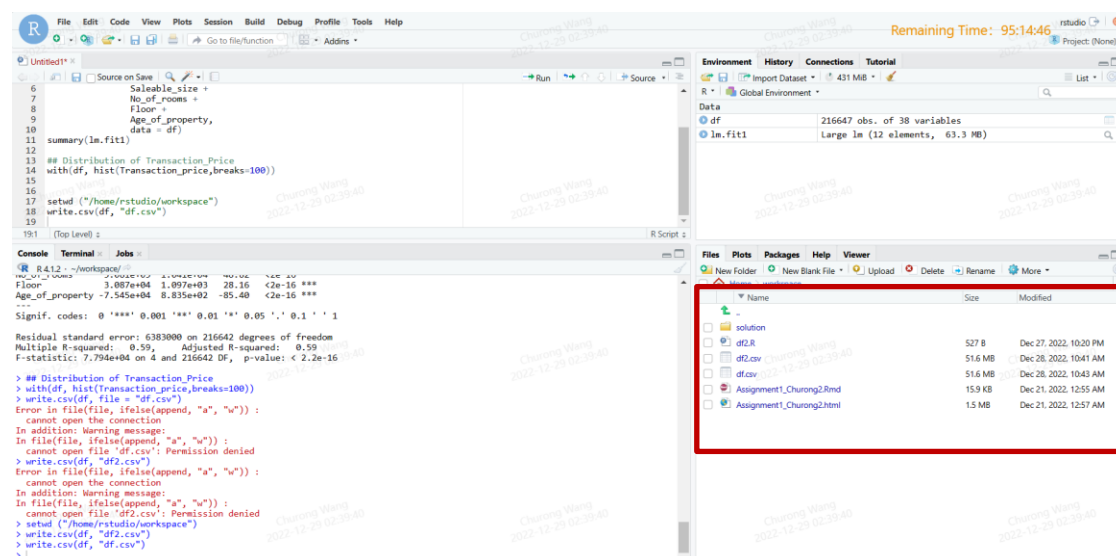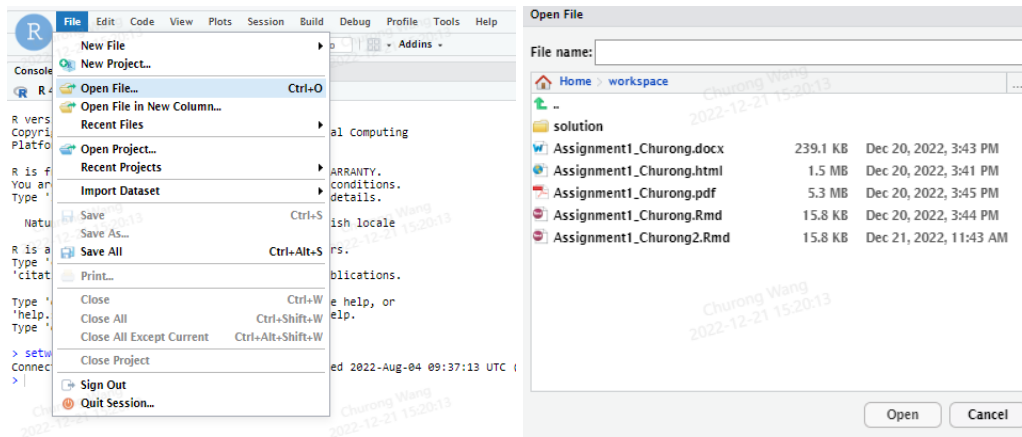
## Dataset Saving

To save a newly created dataset as a CSV file, use the following sample codes to save it in the 'workspace.' Then, you may find your dataset under "workspace" in the bottom right section of the user interface.

```
1. ## Save dataset to workspace
2. setwd ("/home/rstudio/workspace")
3. write.csv(df, file = "df.csv")
```
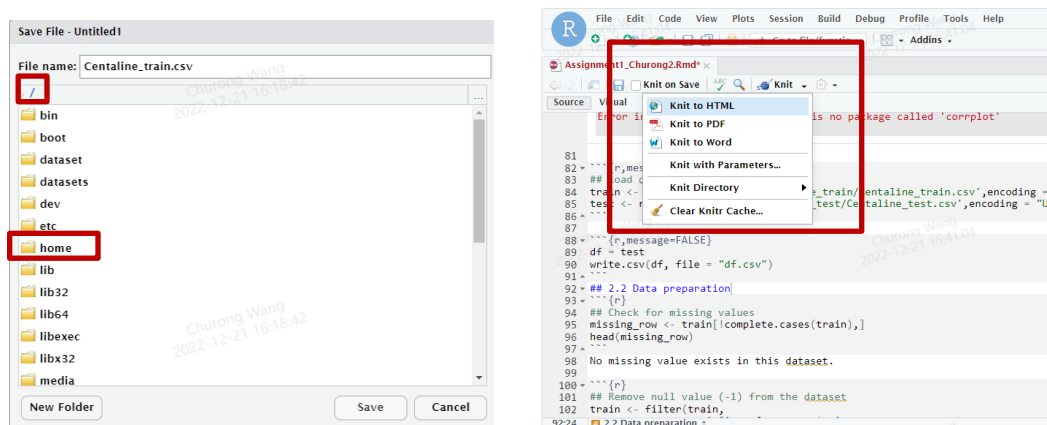
## File Access

All files are stored in a folder called "workspace." Access the "workspace" by clicking "File > Open File > workspace."
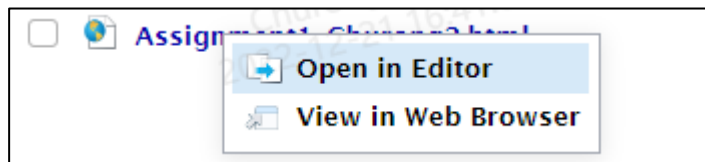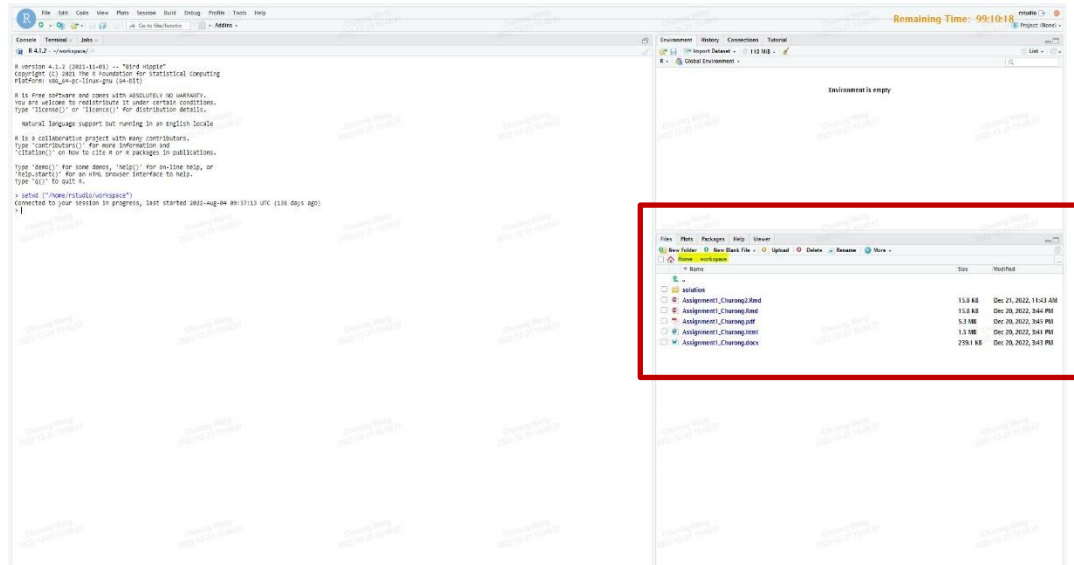


## File Saving

All files must be saved in the "workspace" for future access. Save R files to "workspace" by clicking "File > Save" or "File > Save As…" and browse to the directory "/ > home > rstudio > workspace." Export R codes to HTML/PDF/Word format by clicking 'Knit > Knit to …'.



To check files other than R format (e.g., PDF/HTML/Word format), check out files on the bottom right of the user interface. The "Files" section displays all the files under the folder 'workspace.' Double-click the files to see them on a pop-up web page. For HTML format files, click it and select "View in Web Browser" to open the HTML file in a new web page.

Final Note:

1. Please save all files into the folder "workspace," and ensure that your UID (student ID) and name are contained in the file name when saving your files. These two steps are very important because if you don't follow these instructions, your assignment may not be visible/identifiable to course instructors to grade your assignment.

2. Labs will automatically close after 30 mins of inactivity. Please always remember to save ALL your working files to "workspace" so you can retrieve your previous work in the future.

3. To save your lab time, remember to click 'Stop' to close the labs when you do not use the lab service.