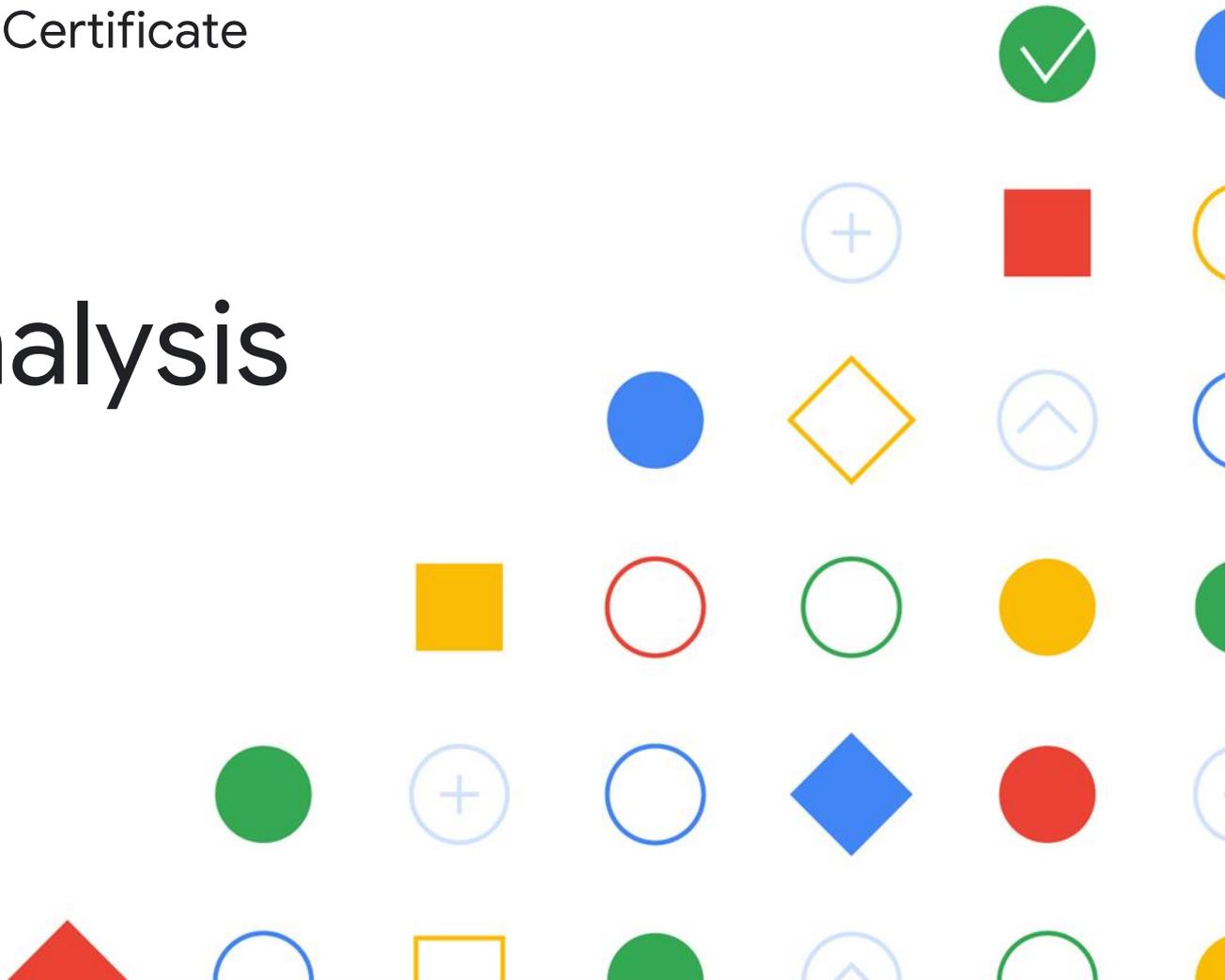


# 7. Data Analysis

*with R Programming*



# Overview:

01

Programming and Data Analytics

02

Programming using RStudio

03

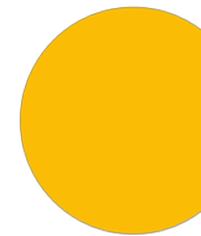
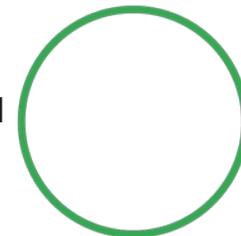
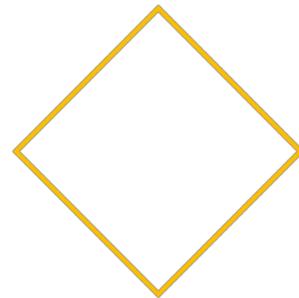
Working with Data in R

04

More about Visualizations, Aesthetics and Annotations

05

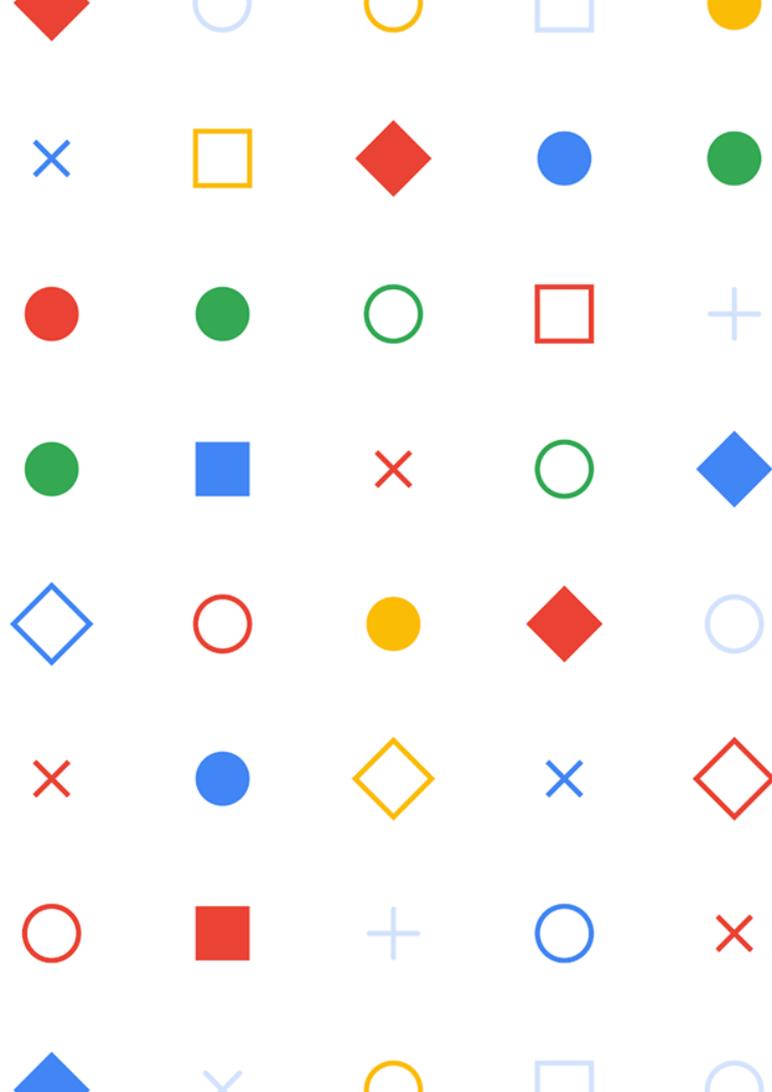
Documentation and Reports



# Programming and Data Analytics

— *Part1:*

1. Why Programming?
2. Why R?
3. Introducing RStudio



# Why Programming?

การเขียนโปรแกรมคือการเขียนคำสั่ง (Instructions) ให้คอมพิวเตอร์ทำตาม โดยคำสั่งที่เขียนต้องเป็นไปตาม “ไวยากรณ์” (Syntax) ของภาษาที่ใช้ เช่นภาษา Python กับ R มี Syntax ที่ต่างกัน

## Why Programming?

ช่วยให้ขั้นตอนการวิเคราะห์ข้อมูลมีความชัดเจน  
(Clarify the steps of your analysis)

ช่วยประหยัดเวลาในการวิเคราะห์ข้อมูล  
(Saves time)

ทำการวิเคราะห์ซ้ำกับข้อมูลใหม่ได้ง่าย และแชร์ code หรือผลลัพธ์ได้สะดวก  
(Reproduce and share your work)

## Examples

เราสามารถดูขั้นตอนการทำความสะอาดและวิเคราะห์ข้อมูลได้จาก Code ที่เราหรือคนอื่นเขียน ทำให้ตรวจสอบ Logic ได้ง่าย ต่างจาก Spreadsheets ที่เห็นแต่ผลลัพธ์สุดท้าย ไม่ได้มีสรุปขั้นตอนให้ชัดเจน

เราทำความสะอาดข้อมูลได้ด้วย Code ไม่กี่บรรทัด (บางทีบรรทัดเดียว) แต่การทำความสะอาดข้อมูลใน Spreadsheets อาจต้องทำหลายขั้นตอน

เราเขียนโปรแกรมให้นำเข้าข้อมูล ทำความสะอาดข้อมูล วิเคราะห์และสร้าง Report รายเดือนแบบอัตโนมัติได้ แต่ถ้าใช้ Spreadsheets เราต้องมานั่งทำแต่ละขั้นตอนด้วยตนเอง



# Why R?

## What is R?

- R คือภาษาคอมพิวเตอร์ (Programming language) ที่นิยมใช้สำหรับการวิเคราะห์ทางสถิติ (Statistical analysis) การสร้างกราฟและแผนภาพ (Visualizations) และการวิเคราะห์ข้อมูลทั่วไป (Data Analysis)

## Why R?

- Accessible เข้าถึงได้ง่ายสำหรับผู้เริ่มต้นเขียนโปรแกรม
- Data-centric ถูกออกแบบเพื่อเล่นกับข้อมูลโดยเฉพาะ
- Open source ใคร ๆ ก็ดาวน์โหลดไปใช้ ดัดแปลง แก้ Bugs หรือสร้าง Packages เสริมใหม่ ๆ แชร์ให้คนอื่นใช้ได้ฟรี
- Community มีสังคมออนไลน์คอยช่วยเหลือ แชร์ความรู้กัน

## What can R do?

- Reproducing your analysis ทำการวิเคราะห์ซ้ำเดิมได้อัตโนมัติแค่คลิกเดียว
- Processing lots of data ประมวลผลข้อมูลปริมาณมาก ๆ ได้
- Creating visualizations สร้างและปรับแต่งกราฟสวย ๆ ได้ในไม่กี่บรรทัด

# Introducing RStudio

RStudio คือซอฟต์แวร์ที่รวมเอาเครื่องมือต่าง ๆ ในการช่วยเขียนโปรแกรมภาษา R มาไว้ในที่เดียว ซึ่งซอฟต์แวร์ลักษณะนี้เรียกว่า **Integrated Development Environment (IDE)** มีหลัก ๆ 4 หน้าต่าง (Panels) ดังนี้

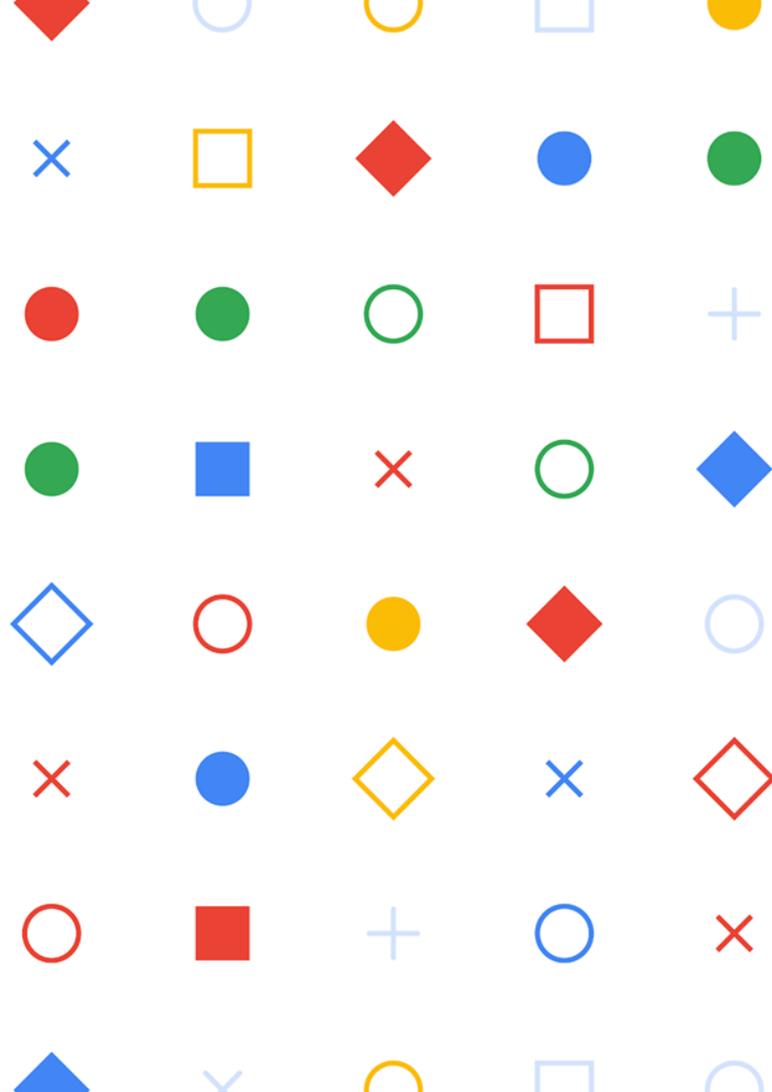
The screenshot shows the RStudio IDE interface with four panels highlighted in red boxes and labeled with Thai text:

- R Source editor**: สำหรับเขียน Code หลายๆ บรรทัด และบันทึกเป็นไฟล์ (Script) ไว้ใช้ทีหลังได้ เมื่อสั่ง Run แล้ว Code จะถูกนำไป Run ใน R Console (Note: ตอนสร้าง Project ใหม่ต้องกด File → New File → R Script ก่อนถึงจะเห็น Pane นี้)
- Environment / History**: Environment: แสดงข้อมูลที่เราโหลดไว้ในโปรเจกต์ History: ประวัติการ Run คำสั่ง (Commands)
- R Console**: R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details. R is a collaborative project with many contributors สำหรับ Run คำสั่ง (Commands) Type 'contributors()' for more information and 'citation()' on how to cite R or R packages in publications. และแสดงผลการคำนวณ Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'q()' to quit R. Connected to your session in progress, last started 2022-Jun-10 01:26:37 UTC (6 minutes ago)
- Files / Plots / Packages / Help**: Files: โครงสร้างไฟล์และโพลเดอร์ของโปรเจกต์ Plots: แสดงกราฟและแผนภาพที่สร้าง Packages: สรุบบทแพคเกจที่เราลงไว้ในโปรเจกต์ Help: คลังข้อมูลช่วยเหลือ

# Programming with RStudio

—— *Part2:*

1. RStudio Workflow
2. Basic Data Types
3. Basic Data Structures
4. Operators and Calculations
5. Tidyverse Packages
6. Nested Commands and Pipe



# RStudio Workflow

1

## Install and Load Packages

- สรุปลงแพ็คเกจที่ถูกลงไว้แล้วด้วยคำสั่ง `installed.packages()`
- ลงแพ็คเกจเพิ่มที่จะใช้ เช่น `install.packages("tidyverse")`
- โหลดแพ็คเกจเพื่อใช้ด้วยคำสั่ง `library(ชื่อแพ็คเกจ)` เช่น `library(tidyverse)`

2

## Import Datasets

- นำเข้า Dataset จากไฟล์ของเราเอง (เช่น ไฟล์ .csv ด้วยฟังก์ชัน `read_csv("ที่ตั้งไฟล์")` จากแพ็คเกจ `readr`)
- หรือใช้ Dataset ตัวอย่างที่มีอยู่แล้ว (เช่น `diamonds` จากแพ็คเกจ `ggplot2`)

3

## "Tidy" the Data

- ทำความสะอาด (Clean) และแปลง (Transform) ข้อมูลเพื่อเตรียมพร้อมสำหรับนำไปวิเคราะห์ต่อ (เช่น ด้วยฟังก์ชันต่าง ๆ ของแพ็คเกจ `tidyr` หรือ `janitor`)

4

## Analyze and Visualize

- วิเคราะห์ข้อมูล เช่น ด้วยฟังก์ชันต่าง ๆ ของแพ็คเกจ `dplyr`
- สร้างกราฟ เช่น ด้วยฟังก์ชันของแพ็คเกจ `ggplot2`

5

## Create and Share Reports

- ใช้ภาษา R Markdown ในการเขียนสรุปขั้นตอนและผลการวิเคราะห์ข้อมูล และบันทึกเป็นไฟล์ `.Rmd`, `.pdf` หรือ `.html` เพื่อส่งให้คนอื่นดูได้

# Basic Data Types

ประเภทของข้อมูลใน R ที่ควรรู้จักเบื้องต้นมีดังนี้

## Examples

## Remarks

### Logical

ข้อมูล ถูก/ผิด

TRUE, FALSE

ต้องตัวพิมพ์ใหญ่ล้วน ห้ามใช้ True, true, False, false

### Integer

จำนวนเต็ม

0, 1, 2, 3, 555, -123

เราเรียกข้อมูลที่เป็นตัวเลข (Integer กับ Double)  
โดยรวมว่า Numeric

### Double

จำนวนทศนิยม (Decimal)

31.44, 10.0, -999.11

### Character

ตัวอักษรหรือข้อความ (String)

"a", 'Great',  
"Jennie is cute!", "TRUE"

จะใช้เครื่องหมายคำพูดแบบ " " (Double quote)  
หรือ ' ' (single quote) ก็ได้

### Datetime\*

วันเวลา

"1991-01-07"  
"07:01:25 UTC"  
"2022-05-02 00:15:30 +07"

\*ต้องลงและโหลดแพ็คเกจที่ชื่อว่า lubridate ก่อน



# Basic Data Structures

โครงสร้างของข้อมูลใน R ที่ควรรู้จักเบื้องต้นมีดังนี้

## (Atomic) Vector

เวกเตอร์

```
c(1, 2, 3, -20)
c("Jennie", "is", "very cute!")
```

## List

ลิสต์

```
list("Great", 555, TRUE)
```

## Data frame

ตารางข้อมูลแบบมาตรฐาน

```
> head(orders)
  order_id customer_name order_date product_id
1   A0001         Great 2022-05-31     P0069
2   A0002         Great 2022-06-01     P0070
3   A0003          Elon 2022-06-01     P4234
4   A0004          Elon 2022-06-02     P4234
5   A0005         Great 2022-06-03     P4234
6   A0006          Larry 2022-05-31     P0069

> class(orders)
[1] "data.frame"
```

## Tibble

ตารางข้อมูลแบบพิเศษ  
ของแพ็คเกจ tidyverse

```
> diamonds
# A tibble: 53,940 × 10
  carat cut      color clarity depth table price
<dbl> <ord> <ord> <ord> <dbl> <dbl> <int>
1  0.23 Ideal  E   SI2    61.5  55  326
2  0.21 Premium E   SI1    59.8  61  326
3  0.23 Good  E   VS1    56.9  65  327
4  0.29 Premium I   VS2    62.4  58  334
5  0.31 Good  J   SI2    63.3  58  335
6  0.24 Very Good J   VS2    62.8  57  336
```

## Remarks

ข้อมูลใน Vector ต้องเป็นประเภทเดียวกันนั่นคือเป็น Numeric ล้วน, Character ล้วน หรือ Logical ล้วน

ข้อมูลใน List ไม่จำเป็นต้องเป็นประเภทเดียวกันก็ได้ และสามารถสร้าง List ของ List (Nested Lists) ได้ด้วย

ข้อมูลตารางที่ Import เข้ามาด้วยฟังก์ชัน read.csv("ที่ตั้งไฟล์") จะเป็น Data frame แต่ถ้าใช้ฟังก์ชัน read\_csv("ที่ตั้งไฟล์") จากแพ็คเกจ readr จะได้มาเป็น Tibble

Tibble คือ Data frame ที่มีการเพิ่มพฤติกรรมบางอย่างให้ใช้งานได้ง่ายขึ้นกับตารางที่มีขนาดใหญ่ เช่น เลือกบางคอลัมน์ (Subsetting) ได้ง่ายขึ้น แสดงผลสวยขึ้น



# Operators and Calculations

## Assignment operator <-

"ยัดค่า" ใส่ตัวแปรด้วยเครื่องหมาย <-

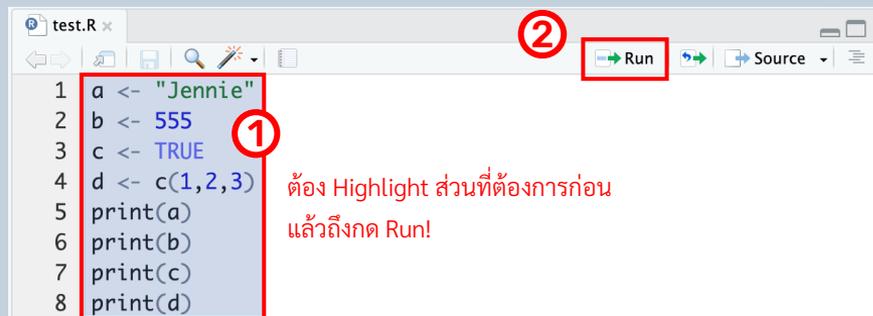
Syntax:

ชื่อตัวแปร <- อะไรก็ได้

Notes :

- ชื่อตัวแปรต้องไม่ขึ้นต้นด้วยตัวเลขหรือสัญลักษณ์แปลก ๆ เช่น lucky99 โอเค แต่ 99lucky ไม่โอเค
- ตัวพิมพ์เล็กใหญ่มีผล เช่น ถ้า a <- 5 และ b <- 10 หากเราคำนวณ A+B จะ error เพราะเรายังไม่ได้ยัดค่าใส่ตัวแปรชื่อ A หรือ B เลย
- เราใช้เครื่องหมาย = แทน <- ก็ได้ แต่ในสังคมคนใช้ R เขานิยมใช้ <- กันส่วนใหญ่

Ex. ลองยัดค่าใส่ตัวแปร แล้ว print ออกมา

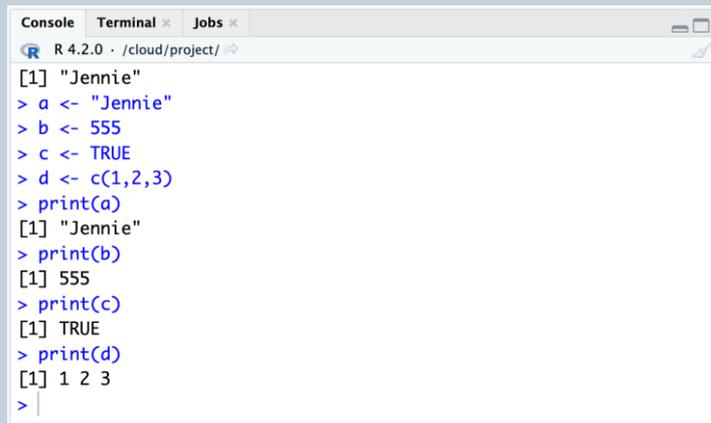


```
test.R x
1 a <- "Jennie"
2 b <- 555
3 c <- TRUE
4 d <- c(1,2,3)
5 print(a)
6 print(b)
7 print(c)
8 print(d)
```

① ต้อง Highlight ส่วนที่ต้องการก่อน แล้วถึงกด Run!

② Run

## Result



```
Console Terminal x Jobs x
R 4.2.0 · /cloud/project/ ↗
[1] "Jennie"
> a <- "Jennie"
> b <- 555
> c <- TRUE
> d <- c(1,2,3)
> print(a)
[1] "Jennie"
> print(b)
[1] 555
> print(c)
[1] TRUE
> print(d)
[1] 1 2 3
>
```

# Operators and Calculations

## Arithmetic operators + - \* /

เครื่องหมายสำหรับการคำนวณ บวก ลบ คูณ หาร

Syntax:

```
a + 10*b + (2/3)
```

Notes :

- สำหรับภาษา R เราไม่สามารถเอาข้อมูลแบบ String มาเชื่อมกันด้วยเครื่องหมาย + ได้เหมือน Python (ต้องใช้ฟังก์ชันชื่อว่า paste() แทน เช่น paste("Jennie","Narak") ได้ "Jennie Narak")
- TRUE คือ 1 และ FALSE คือ 0 ดังนั้น ถ้าลอง TRUE + FALSE + 5 จะเทียบเท่า 1 + 0 + 5 ซึ่งเท่ากับ 6

Ex. ลองคำนวณแบบต่าง ๆ แล้ว print ออกมา

```
test.R x
1 a <- 5
2 b <- 6
3 c <- TRUE
4 d <- c(1,2,3)
5 e <- d + a
6 print(a+b)
7 print(2*a - b/2)
8 print(a + b + c)
9 print(e)
```

## Result

```
Console Terminal x Jobs x
R 4.2.0 · /cloud/project/ ↗
> a <- 5
> b <- 6
> c <- TRUE
> d <- c(1,2,3)
> e <- d + a
> print(a+b)
[1] 11
> print(2*a - b/2)
[1] 7
> print(a + b + c)
[1] 12
> print(e)
[1] 6 7 8
> |
```

# Operators and Calculations

## Logical operators

เครื่องหมายสำหรับการผสมข้อมูลถูก/ผิด

Syntax:

<code>p &amp; q</code>	<code>#AND (และ)</code>
<code>p   q</code>	<code>#OR (หรือ)</code>
<code>!p</code>	<code>#NOT (ตรงข้าม a.k.a. นิเสธ)</code>
<code>a == 2</code>	<code>#a มีค่าเท่ากับ 2 หรือไม่?</code>
<code>a &lt;= b</code>	<code>#a มีค่าน้อยกว่าหรือเท่ากับ b ไหม?</code>

Ex. ลองคำนวณแบบต่าง ๆ แล้ว print ออกมา

```
test.R x
1 #Assign variables
2 a <- 5
3 b <- 6
4 c <- TRUE #Yo!
5 #Compute
6 print(a == b) #Is a equal to b?
7 print(a != b) #Is a NOT equal to b?
8 print((a < 2) & (b >= 6))
9 print((a < 2) | (b >= 6))
10 print(!c)
```

## Result

```
Console Terminal x Jobs x
R 4.2.0 · /cloud/project/
> #Assign variables
> a <- 5
> b <- 6
> c <- TRUE #Yo!
> #Compute
> print(a == b) #Is a equal to b?
[1] FALSE
> print(a != b) #Is a NOT equal to b?
[1] TRUE
> print((a < 2) & (b >= 6))
[1] FALSE
> print((a < 2) | (b >= 6))
[1] TRUE
> print(!c)
[1] FALSE
>
```

# Tidyverse Packages

R มีแพ็คเกจอยู่บนแพลตฟอร์มชื่อว่า CRAN (Comprehensive R Archive Network) ให้เราโหลดมาใช้ฟรี และที่นิยมสำหรับงานวิเคราะห์ข้อมูลคือกลุ่มแพ็คเกจชื่อ Tidyverse ซึ่งมีแพ็คเกจย่อยมากมาย แต่เบื้องต้นที่ควรรู้จักคือ

## Main Purposes

**readr**  
รีด-อาร์

นำเข้าข้อมูล

**tidyr**  
ไทดี-อาร์

ทำความสะอาด (Tidy) และจัดรูปข้อมูล

**dplyr**  
ดี-พลาย-อาร์

กรอง ตัดแปลง คำนวณทางสถิติแบบแบ่งกลุ่ม และอื่น ๆ

**ggplot2**  
จีจีพล็อต-ทู

สร้างและปรับแต่ง Visualizations

## Frequently-used Functions

- `table <- read_csv("ที่ตั้งไฟล์")`
- `as_tibble(df, ...)` #เปลี่ยน Data frame เป็น Tibble
- `unite(table, col1, col2, ...)` #เอาค่า 2 คอลัมน์มา concatenate กัน
- `separate(table, col, sep='/', ...)` #แยกค่าตามตัวคั่นเป็นคอลัมน์ย่อย
- `drop_na(table, คอลัมน์ที่สนใจ)` #เอาแถวที่มีค่าว่างในคอลัมน์ที่สนใจออก
- `filter(orders, price > 100)` #กรองข้อมูล
- `select(table, col1, col2, ...)` #เลือกมาเฉพาะบางคอลัมน์
- `table %>% group_by(col1) %>% summarise(...)` #สรุปข้อมูลแบบแบ่งกลุ่ม
- `geom_bar()`, `geom_histogram()`, `geom_line()`, `geom_boxplot()`, `geom_violin()`, `geom_point()` และอื่นๆ อีกมากมาย!

# Nested Commands and Pipes

## Nested functions

ฟังก์ชันที่ฝังอยู่ด้านในของอีกฟังก์ชันหนึ่ง

Syntax:

```
func2(func1(data, ...), ...)
```

Notes:

- การประมวลผลจะเริ่มต้นจากฟังก์ชันด้านในสุดก่อน แล้วผลของฟังก์ชันด้านใน จะเป็น Input ให้กับฟังก์ชันด้านนอก
- จากที่เห็น func1 จัดเป็น Nested function เพราะฝังอยู่ด้านใน func2 อีกทีหนึ่ง

Ex. จากตารางรายการสั่งซื้อสินค้า จงกรองเอาเฉพาะสินค้าราคา 5 บาท และแสดงผลเรียงลำดับรายได้ต่อการสั่งซื้อ จากน้อยไปมาก

```
1 orders <- read_csv('./data/orders_full.csv')
2
3 orders$revenue = orders$product_price * orders$quantity
```

```
> orders
# A tibble: 10 × 7
  order_id customer_name order_date product_id product_price quantity revenue
  <chr>      <chr>          <date>   <chr>          <dbl>   <dbl>   <dbl>
1 A0001     Great        2022-05-31 P0069           23000     1    23000
2 A0002     Great        2022-06-01 P0070           2000     2     4000
3 A0003     Elon         2022-06-01 P4234            5    10000    50000
4 A0004     Elon         2022-06-02 P4234            5    20000   100000
5 A0005     Great        2022-06-03 P4234            5    30000   150000
6 A0006     Larry        2022-05-31 P0069           23000     1    23000
7 A0007     Great        2022-06-04 P4234            5     300     1500
8 A0008     Elon         2022-06-04 P4234            5    40000   200000
9 A0009     Elon         2022-06-05 P4234            5    50000   250000
10 A0010     Great        2022-06-06 P4234            5   100000  500000
```

```
4
5 arrange(filter(orders, product_price == 5), revenue)
```

## Result

```
> arrange(filter(orders, product_price == 5), revenue)
# A tibble: 7 × 7
  order_id customer_name order_date product_id product_price quantity revenue
  <chr>      <chr>          <date>   <chr>          <dbl>   <dbl>   <dbl>
1 A0007     Great        2022-06-04 P4234            5     300     1500
2 A0003     Elon         2022-06-01 P4234            5    10000    50000
3 A0004     Elon         2022-06-02 P4234            5    20000   100000
4 A0005     Great        2022-06-03 P4234            5    30000   150000
5 A0008     Elon         2022-06-04 P4234            5    40000   200000
6 A0009     Elon         2022-06-05 P4234            5    50000   250000
7 A0010     Great        2022-06-06 P4234            5   100000  500000
```

# Nested Commands and Pipes

## Pipes %>%

ตัวดำเนินการของ R ที่ส่งต่อผลลัพธ์ของฟังก์ชันแรก ไปยังฟังก์ชันถัดไปเป็นทอด ๆ เหมือนไหลไปตามท่อ

Syntax:

```
data %>%  
  func1(...) %>%  
  func2(...)
```

Note: สังเกตว่า เราไม่ต้องเอา data ตั้งต้นไปใส่ในวงเล็บของฟังก์ชันอีกต่อไปแล้ว เพราะ %>% ทำการ "ยัดค่า" data ใส่ฟังก์ชันตัวถัดไปให้เราเอง!

Ex. จากตารางรายการสั่งซื้อสินค้า จงกรองเอาเฉพาะสินค้าราคา 5 บาท และแสดงผลเรียงลำดับรายได้ต่อการสั่งซื้อ จากน้อยไปมาก

```
5 arrange(filter(orders, product_price == 5), revenue)  
6
```

```
7 orders %>% filter(product_price == 5) %>% arrange(revenue)
```

## Results

```
> arrange(filter(orders, product_price == 5), revenue)  
# A tibble: 7 × 7  
  order_id customer_name order_date product_id product_price quantity revenue  
  <chr>      <chr>          <date>   <chr>          <dbl>   <dbl>   <dbl>  
1 A0007      Great          2022-06-04 P4234             5         300     1500  
2 A0003      Elon            2022-06-01 P4234             5      10000    50000  
3 A0004      Elon            2022-06-02 P4234             5      20000   100000  
4 A0005      Great           2022-06-03 P4234             5      30000   150000  
5 A0008      Elon            2022-06-04 P4234             5      40000   200000  
6 A0009      Elon            2022-06-05 P4234             5      50000   250000  
7 A0010      Great           2022-06-06 P4234             5     100000  500000
```

```
> orders %>% filter(product_price == 5) %>% arrange(revenue)  
# A tibble: 7 × 7  
  order_id customer_name order_date product_id product_price quantity revenue  
  <chr>      <chr>          <date>   <chr>          <dbl>   <dbl>   <dbl>  
1 A0007      Great          2022-06-04 P4234             5         300     1500  
2 A0003      Elon            2022-06-01 P4234             5      10000    50000  
3 A0004      Elon            2022-06-02 P4234             5      20000   100000  
4 A0005      Great           2022-06-03 P4234             5      30000   150000  
5 A0008      Elon            2022-06-04 P4234             5      40000   200000  
6 A0009      Elon            2022-06-05 P4234             5      50000   250000  
7 A0010      Great           2022-06-06 P4234             5     100000  500000
```

# Nested Commands and Pipes

## Pipes %>%

ตัวดำเนินการของ R ที่ส่งต่อผลลัพธ์ของฟังก์ชันแรก ไปยังฟังก์ชันถัดไปเป็นทอด ๆ เหมือนไหลไปตามท่อ

Syntax:

```
data %>%  
  func1(...) %>%  
  func2(...)
```

Note: สังเกตว่า เราไม่ต้องเอา data ตั้งต้นไปใส่ในวงเล็บของฟังก์ชันอีกต่อไปแล้ว เพราะ %>% ทำการ "ยัดค่า" data ใส่ฟังก์ชันตัวถัดไปให้เราเอง!

Ex. จากตารางรายการสั่งซื้อสินค้า จงกรองเอาเฉพาะสินค้าราคา 5 บาท จากนั้นหาผลรวมรายได้แบ่งตามลูกค้าแต่ละคน

```
9 summarize(group_by(filter(orders, product_price == 5), customer_name), sum_revenue=sum(revenue))  
10
```

```
11 orders %>%  
12   filter(product_price == 5) %>%  
13   group_by(customer_name) %>%  
14   summarize(sum_revenue = sum(revenue))
```

## Results

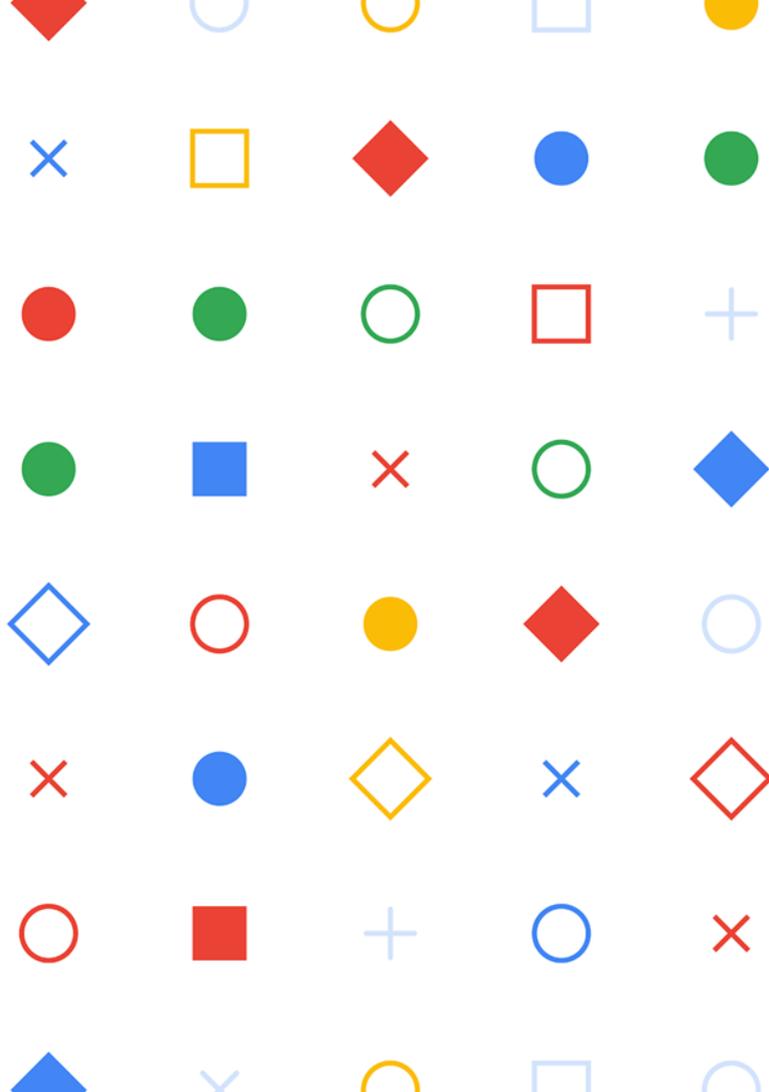
```
> summarize(group_by(filter(orders, product_price == 5), customer_name), sum_revenue=sum(revenue))  
# A tibble: 2 × 2  
  customer_name sum_revenue  
  <chr>          <dbl>  
1 Elon           600000  
2 Great          651500
```

```
> orders %>%  
+   filter(product_price == 5) %>%  
+   group_by(customer_name) %>%  
+   summarize(sum_revenue = sum(revenue))  
# A tibble: 2 × 2  
  customer_name sum_revenue  
  <chr>          <dbl>  
1 Elon           600000  
2 Great          651500
```

# Working with Data in R

## — *Part3:*

1. Data Frames vs. Tibbles
2. Tidy data standards
3. Clean: Inspecting Data Frames
4. Clean: Selecting and Renaming Columns
5. Organize: Sorting and Filtering
6. Organize: Summarize and Group By
7. Transform: Separate, Unite and Mutate
8. Warning: Anscombe's Quartet
9. The Bias Function



# Data Frames vs. Tibbles

## Data frames

โครงสร้างเก็บข้อมูลตารางแบบมาตรฐานของ R

Notes:

- คอลัมน์ควรมีชื่อที่สอดคล้องกับข้อมูลในแต่ละคอลัมน์  
(Columns should be named)
- ข้อมูลในตารางมีได้หลาย Types ไม่ว่าจะเป็น Numeric, Factor หรือ Character
- แต่ละคอลัมน์ควรมีจำนวนข้อมูลที่เท่ากัน  
(Each column should contain the same number of data items) ซึ่งรวมถึงค่าว่าง (Missing values) ด้วย

Ex. โหลดตารางรายการสั่งซื้อสินค้าของเดือนมิถุนายน 2022 ใส่ตัวแปร orders

```
> orders <- read.csv('./data/orders_2022_06.csv')
> class(orders)
[1] "data.frame"
> orders
```

	Order.ID	Customer.Name	Band	Order.Date	Product.ID	Product.Price	QTT
1	A0001	Great	SKD69	2022-05-31	P0069	23000	1
2	A0002	Great	SKD69	2022-06-01	P0070	2000	2
3	A0003	Elon	SHIBA	2022-06-01	P4234	5	10000
4	A0004	Elon	SHIBA	2022-06-02	P4234	5	20000
5	A0005	Great	SKD69	2022-06-03	P4234	5	30000
6	A0006	Larry	GOOGL	2022-05-31	P0069	23000	1
7	A0007	Great	SKD69	2022-06-04	P4234	5	300
8	A0008	Elon	SHIBA	2022-06-04	P4234	5	40000
9	A0009	Elon	SHIBA	2022-06-05	P4234	5	50000
10	A0010	Great	SKD69	2022-06-06	P4234	5	100000
11	A0011	Sergey	GOOGL	2022-06-08	P4234	5	10000
12	A0012	Jeff	AMAZE	2022-06-10	P4234	5	20000
13	A0013	Jeff	AMAZE	2022-06-11	P4234	5	30000
14	A0014	Sergey	GOOGL	2022-06-12	P4234	5	40000
15	A0015	Sergey	GOOGL	2022-06-13	P4234	5	50000
16	A0016	Sergey	GOOGL	2022-06-14	P4234	5	60000
17	A0017	Jeff	AMAZE	2022-06-15	P4234	5	70000
18	A0018	Mark	FB555	2022-06-16	P4234	5	10
19	A0019	Mark	FB555	2022-06-17	P4234	5	20
20	A0020	Mark	FB555	2022-06-18	P4234	5	30
21	A0021	Sergey	GOOGL	2022-06-19	P4234	5	10000
22	A0022	Sergey	GOOGL	2022-06-20	P4234	5	20000
23	A0023	Sergey	GOOGL	2022-06-24	P4234	5	30000
24	A0024	Sergey	GOOGL	2022-06-27	P4234	5	40000
25	A0025	Sergey	GOOGL	2022-06-30	P4234	5	50000

# Data Frames vs. Tibbles

## Tibbles

โครงสร้างเก็บข้อมูลตารางแบบพิเศษของ tidyverse

Notes:

- Tibble คือ Data frame ที่ถูกออกแบบมาให้ทำงานได้สะดวกคล่องตัวขึ้น (Streamlined) โดยมีสมบัติต่อไปนี้
  - ประเภทของข้อมูลแต่ละคอลัมน์จะไม่เปลี่ยนจากต้นฉบับ เช่น คอลัมน์ตัวเลขก็ต้องเป็นตัวเลขไปตลอด
  - ชื่อตัวแปรของแต่ละคอลัมน์จะไม่เปลี่ยน (เว้นแต่จะใช้ฟังก์ชันเพื่อเปลี่ยนชื่อโดยเฉพาะ)
  - ไม่มีชื่อแถว มีแต่หมายเลขแถว
  - เวลา Print จะแสดงเฉพาะ 10 แถวแรกเพื่อให้ไม่รกหน้าจอ console

Ex. เปลี่ยนตาราง orders จาก data.frame เป็น tibble โดยบันทึกใส่ตัวแปร orders\_tb

```
> orders_tb <- as_tibble(orders)
> class(orders_tb)
[1] "tbl_df"      "tbl"        "data.frame"
> orders_tb
# A tibble: 25 × 7
  Order.ID Customer.Name Band Order.Date Product.ID Product.Price QTT
  <chr>    <chr>          <chr> <chr>    <chr>          <int> <int>
1 A0001    Great           SKD69 2022-05-31 P0069          23000 1
2 A0002    Great           SKD69 2022-06-01 P0070          2000 2
3 A0003    Elon            SHIBA 2022-06-01 P4234           5 10000
4 A0004    Elon            SHIBA 2022-06-02 P4234           5 20000
5 A0005    Great           SKD69 2022-06-03 P4234           5 30000
6 A0006    Larry           GOOGL 2022-05-31 P0069          23000 1
7 A0007    Great           SKD69 2022-06-04 P4234           5 300
8 A0008    Elon            SHIBA 2022-06-04 P4234           5 40000
9 A0009    Elon            SHIBA 2022-06-05 P4234           5 50000
10 A0010   Great           SKD69 2022-06-06 P4234           5 100000
# ... with 15 more rows
```

# Tidy Data Standards

Tidy data คือมาตรฐานการจัดการข้อมูลในระบบนิเวศของ R ซึ่งมี 3 ข้อดังนี้

order_id	name	order_date	product
A0001	Great	2022-05-31	PS5
A0002	Great	2022-06-01	LAC
A0003	Elon	2022-06-01	PS5
A0004	Elon	2022-06-02	LAC
A0005	Great	2022-06-03	LAC
A0006	Larry	2022-05-31	PS5

order_id	name	order_date	product
<del>A0001</del>	Great	2022-05-31	PS5
<del>A0002</del>	Great	2022-06-01	LAC
<del>A0003</del>	Elon	2022-06-01	PS5
<del>A0004</del>	Elon	2022-06-02	LAC
<del>A0005</del>	Great	2022-06-03	LAC
<del>A0006</del>	Larry	2022-05-31	PS5

order_id	name	order_date	product
A0001	Great	2022-05-31	PS5
A0002	Great	2022-06-01	LAC
A0003	Elon	2022-06-01	PS5
A0004	Elon	2022-06-02	LAC
A0005	Great	2022-06-03	LAC
A0006	Larry	2022-05-31	PS5

1

Variables are organized into columns

ตัวแปรถูกจัดให้อยู่ในรูปคอลัมน์

2

Observations are organized into rows

1 แถวคือ 1 รายการ

3

Each value must have its own cell

1 ช่อง (cell) ในตารางมีข้อมูลได้ 1 ค่า

Source: <https://r4ds.had.co.nz/tidy-data.html>

# Clean: Inspecting Data Frames

เราสำรวจตารางเบื้องต้นได้ด้วยฟังก์ชันต่อไปนี้

Syntax:

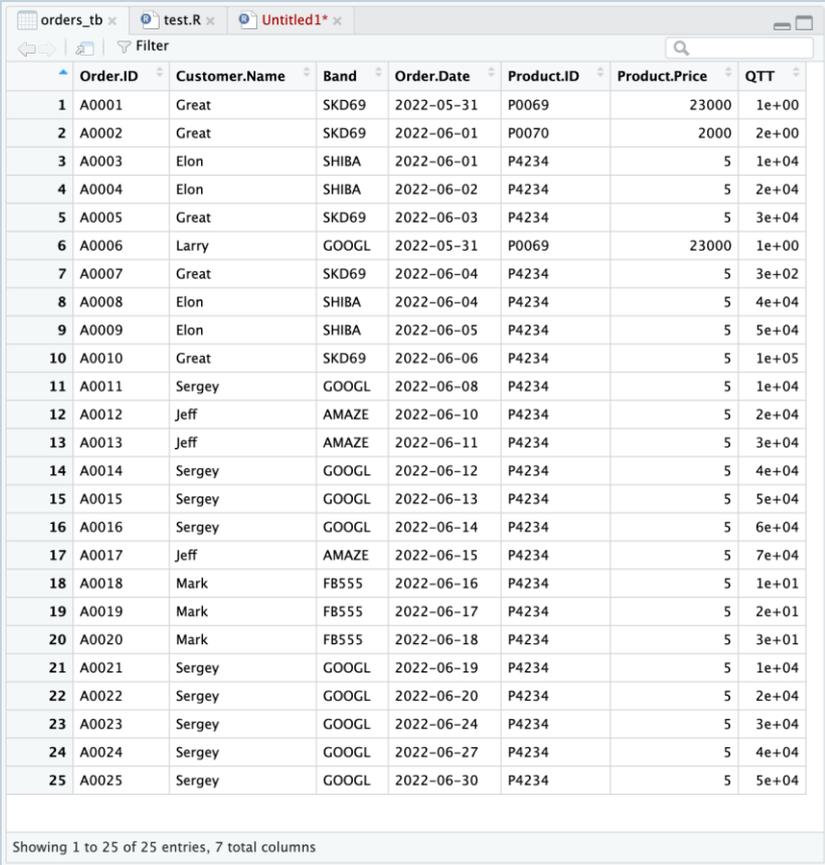
```
View(table)           #ดูตารางแบบเต็มๆ
head(table)           #ดู 6 แถวบน
colnames(table)       #ดูชื่อคอลัมน์ทั้งหมด
#ดูโครงสร้างตาราง
str(table)
glimpse(table)
#ดูสรุปข้อมูลในตารางคร่าว ๆ
skim_without_charts(table)
```

Note: glimpse() เป็นของแพ็คเกจ dplyr และ  
skim\_without\_charts() เป็นของแพ็คเกจ skimr

Ex. ดูตาราง orders\_tb แบบเต็มๆ

```
> View(orders_tb)
```

← V ใน View ต้องตัวพิมพ์ใหญ่!!!



Order.ID	Customer.Name	Band	Order.Date	Product.ID	Product.Price	QTT	
1	A0001	Great	SKD69	2022-05-31	P0069	23000	1e+00
2	A0002	Great	SKD69	2022-06-01	P0070	2000	2e+00
3	A0003	Elon	SHIBA	2022-06-01	P4234	5	1e+04
4	A0004	Elon	SHIBA	2022-06-02	P4234	5	2e+04
5	A0005	Great	SKD69	2022-06-03	P4234	5	3e+04
6	A0006	Larry	GOOGL	2022-05-31	P0069	23000	1e+00
7	A0007	Great	SKD69	2022-06-04	P4234	5	3e+02
8	A0008	Elon	SHIBA	2022-06-04	P4234	5	4e+04
9	A0009	Elon	SHIBA	2022-06-05	P4234	5	5e+04
10	A0010	Great	SKD69	2022-06-06	P4234	5	1e+05
11	A0011	Sergey	GOOGL	2022-06-08	P4234	5	1e+04
12	A0012	Jeff	AMAZE	2022-06-10	P4234	5	2e+04
13	A0013	Jeff	AMAZE	2022-06-11	P4234	5	3e+04
14	A0014	Sergey	GOOGL	2022-06-12	P4234	5	4e+04
15	A0015	Sergey	GOOGL	2022-06-13	P4234	5	5e+04
16	A0016	Sergey	GOOGL	2022-06-14	P4234	5	6e+04
17	A0017	Jeff	AMAZE	2022-06-15	P4234	5	7e+04
18	A0018	Mark	FB555	2022-06-16	P4234	5	1e+01
19	A0019	Mark	FB555	2022-06-17	P4234	5	2e+01
20	A0020	Mark	FB555	2022-06-18	P4234	5	3e+01
21	A0021	Sergey	GOOGL	2022-06-19	P4234	5	1e+04
22	A0022	Sergey	GOOGL	2022-06-20	P4234	5	2e+04
23	A0023	Sergey	GOOGL	2022-06-24	P4234	5	3e+04
24	A0024	Sergey	GOOGL	2022-06-27	P4234	5	4e+04
25	A0025	Sergey	GOOGL	2022-06-30	P4234	5	5e+04

Showing 1 to 25 of 25 entries, 7 total columns

# Clean: Inspecting Data Frames

เราสำรวจตารางเบื้องต้นได้ด้วยฟังก์ชันต่อไปนี้

Syntax:

```
View(table)           #ดูตารางแบบเต็มๆ
head(table)           #ดู 6 แถวบน
colnames(table)       #ดูชื่อคอลัมน์ทั้งหมด
#ดูโครงสร้างตาราง
str(table)
glimpse(table)
#ดูสรุปข้อมูลในตารางคร่าว ๆ
skim_without_charts(table)
```

Note: glimpse() เป็นของแพ็คเกจ dplyr และ  
skim\_without\_charts() เป็นของแพ็คเกจ skimr

Ex. ดู 6 แถวแรกของตาราง orders\_tb และสรุปรายชื่อคอลัมน์ทั้งหมดของตาราง

```
> head(orders_tb)
# A tibble: 6 × 7
  Order.ID Customer.Name Band Order.Date Product.ID Product.Price QTT
  <chr>    <chr>          <chr> <chr>    <chr>          <int> <int>
1 A0001    Great            SKD69 2022-05-31 P0069          23000 1
2 A0002    Great            SKD69 2022-06-01 P0070          2000 2
3 A0003    Elon              SHIBA 2022-06-01 P4234           5 10000
4 A0004    Elon              SHIBA 2022-06-02 P4234           5 20000
5 A0005    Great            SKD69 2022-06-03 P4234           5 30000
6 A0006    Larry            GOOGL 2022-05-31 P0069          23000 1

> colnames(orders_tb)
[1] "Order.ID"      "Customer.Name" "Band"          "Order.Date"    "Product.ID"
[6] "Product.Price" "QTT"
```

# Clean: Inspecting Data Frames

เราสำรวจตารางเบื้องต้นได้ด้วยฟังก์ชันต่อไปนี้

Syntax:

```
View(table)           #ดูตารางแบบเต็มๆ
head(table)           #ดู 6 แถวบน
colnames(table)       #ดูชื่อคอลัมน์ทั้งหมด
#ดูโครงสร้างตาราง
str(table)
glimpse(table)
#ดูสรุปข้อมูลในตารางคร่าว ๆ
skim_without_charts(table)
```

Note: glimpse() เป็นของแพ็คเกจ dplyr และ  
skim\_without\_charts() เป็นของแพ็คเกจ skimr

Ex. ดูโครงสร้างตาราง orders\_tb

```
> str(orders_tb)
tibble [25 x 7] (S3: tbl_df/tbl/data.frame)
 $ Order.ID      : chr [1:25] "A0001" "A0002" "A0003" "A0004" ...
 $ Customer.Name: chr [1:25] "Great" "Great" "Elon" "Elon" ...
 $ Band         : chr [1:25] "SKD69" "SKD69" "SHIBA" "SHIBA" ...
 $ Order.Date   : chr [1:25] "2022-05-31" "2022-06-01" "2022-06-01" "2022-06-02" ...
 $ Product.ID   : chr [1:25] "P0069" "P0070" "P4234" "P4234" ...
 $ Product.Price: int [1:25] 23000 2000 5 5 5 23000 5 5 5 5 ...
 $ QTT         : int [1:25] 1 2 10000 20000 30000 1 300 40000 50000 100000 ...
> glimpse(orders_tb)
Rows: 25
Columns: 7
 $ Order.ID      <chr> "A0001", "A0002", "A0003", "A0004", "A0005", "A0006", "A0007...
 $ Customer.Name <chr> "Great", "Great", "Elon", "Elon", "Great", "Larry", "Great",...
 $ Band         <chr> "SKD69", "SKD69", "SHIBA", "SHIBA", "SKD69", "GOOGL", "SKD69...
 $ Order.Date   <chr> "2022-05-31", "2022-06-01", "2022-06-01", "2022-06-02", "202...
 $ Product.ID   <chr> "P0069", "P0070", "P4234", "P4234", "P4234", "P0069", "P4234...
 $ Product.Price <int> 23000, 2000, 5, 5, 5, 23000, 5, 5, 5, 5, 5, 5, 5, 5, 5, 5...
 $ QTT         <int> 1, 2, 10000, 20000, 30000, 1, 300, 40000, 50000, 100000, 100...
```

# Clean: Inspecting Data Frames

เราสำรวจตารางเบื้องต้นได้ด้วยฟังก์ชันต่อไปนี้

Syntax:

```
View(table)           #ดูตารางแบบเต็มๆ
head(table)           #ดู 6 แถวบน
colnames(table)       #ดูชื่อคอลัมน์ทั้งหมด
#ดูโครงสร้างตาราง
str(table)
glimpse(table)
#ดูสรุปข้อมูลในตารางคร่าว ๆ
skim_without_charts(table)
```

Note: glimpse() เป็นของแพ็คเกจ dplyr และ  
skim\_without\_charts() เป็นของแพ็คเกจ skimr

Ex. ดูสรุปข้อมูลในตาราง orders\_tb แบบคร่าว ๆ

```
> skim_without_charts(orders_tb)
— Data Summary —
Name                values
Number of rows      25
Number of columns    7

-----
Column type frequency:
character            5
numeric              2

-----
Group variables      None

— Variable type: character —
skim_variable n_missing complete_rate min max empty n_unique whitespace
1 Order.ID     0             1 5 5 0 25 0
2 Customer.Name 0             1 4 6 0 6 0
3 Band         0             1 5 5 0 5 0
4 Order.Date   0             1 10 10 0 22 0
5 Product.ID   0             1 5 5 0 3 0

— Variable type: numeric —
skim_variable n_missing complete_rate mean sd p0 p25 p50 p75 p100
1 Product.Price 0             1 1924. 6356. 5 5 5 5 23000
2 QTT           0             1 27215. 26208. 1 300 20000 40000 100000
```

# Clean: Selecting and Renaming Columns

เราเลือกบางคอลัมน์มาใช้ได้ด้วยฟังก์ชัน select()

Syntax:

```
#เลือกคอลัมน์ col1 และ col2 จากตาราง
select(table, col1, col2, ...)
#ใช้วิธี Pipe ก็ได้
table %>%
  select(col1, col2, ...)
#เลือกทุกคอลัมน์ยกเว้น col1
table %>%
  select(-col1)
```

Note: หากต้องการเก็บผลไว้ใช้ต่อ ต้องบันทึกใส่ตัวแปรใหม่ เช่น bands <- select(orders\_tb, Band)

Ex. เลือกคอลัมน์ Customer.Name และ Band

```
> select(orders_tb, Customer.Name, Band)
# A tibble: 25 x 2
  Customer.Name Band
  <chr>         <chr>
1 Great        SKD69
2 Great        SKD69
3 Elon         SHIBA
4 Elon         SHIBA
5 Great        SKD69
6 Larry        GOOGL
7 Great        SKD69
8 Elon         SHIBA
9 Elon         SHIBA
10 Great       SKD69
# ... with 15 more rows
```

```
> orders_tb %>% select(Customer.Name, Band)
# A tibble: 25 x 2
  Customer.Name Band
  <chr>         <chr>
1 Great        SKD69
2 Great        SKD69
3 Elon         SHIBA
4 Elon         SHIBA
5 Great        SKD69
6 Larry        GOOGL
7 Great        SKD69
8 Elon         SHIBA
9 Elon         SHIBA
10 Great       SKD69
# ... with 15 more rows
```

Ex. เลือกทุกคอลัมน์ยกเว้น Band มาดู

```
> orders_tb %>% select(-Band)
# A tibble: 25 x 6
  Order.ID Customer.Name Order.Date Product.ID Product.Price QTT
  <chr>     <chr>         <chr>     <chr>         <int> <int>
1 A0001    Great        2022-05-31 P0069         23000 1
2 A0002    Great        2022-06-01 P0070         2000 2
3 A0003    Elon         2022-06-01 P4234          5 10000
4 A0004    Elon         2022-06-02 P4234          5 20000
5 A0005    Great        2022-06-03 P4234          5 30000
6 A0006    Larry        2022-05-31 P0069         23000 1
7 A0007    Great        2022-06-04 P4234          5 300
8 A0008    Elon         2022-06-04 P4234          5 40000
9 A0009    Elon         2022-06-05 P4234          5 50000
10 A0010   Great        2022-06-06 P4234          5 100000
# ... with 15 more rows
```

# Clean: Selecting and Renaming Columns

มีหลายฟังก์ชันที่ใช้เปลี่ยนชื่อคอลัมน์ให้เลือกใช้

Syntax:

```
#เปลี่ยนคอลัมน์เดียว
```

```
rename(table, ชื่อcolใหม่=ชื่อcolเดิม)
```

```
#เปลี่ยนทีละตัวทุกคอลัมน์ เช่น เป็นตัวพิมพ์เล็ก
```

```
rename_with(table, tolower)
```

```
#เปลี่ยนให้ "สะอาด" ตามมาตรฐาน Tidy data
```

```
clean_names(table)
```

Notes:

- rename() และ rename\_with() เป็นของแพ็คเกจ dplyr ส่วน clean\_names() เป็นของแพ็คเกจ janitor
- ถ้าลืมนับที่ก็ผลใส่ตัวแปร!  
เช่น orders\_tb <- clean\_names(orders\_tb)

Ex. เปลี่ยนชื่อคอลัมน์ QTT เป็น Quantity

```
> orders_tb <- rename(orders_tb, Quantity = QTT)
```

```
> head(orders_tb)
```

```
# A tibble: 6 x 7
```

	Order.ID	Customer.Name	Band	Order.Date	Product.ID	Product.Price	Quantity
	<chr>	<chr>	<chr>	<chr>	<chr>	<int>	<int>
1	A0001	Great	SKD69	2022-05-31	P0069	23000	1
2	A0002	Great	SKD69	2022-06-01	P0070	2000	2
3	A0003	Elon	SHIBA	2022-06-01	P4234	5	10000
4	A0004	Elon	SHIBA	2022-06-02	P4234	5	20000
5	A0005	Great	SKD69	2022-06-03	P4234	5	30000
6	A0006	Larry	GOOGL	2022-05-31	P0069	23000	1

Ex. เปลี่ยนชื่อคอลัมน์ทั้งหมดให้ "สะอาด" ด้วย clean\_names()

```
> orders_tb <- clean_names(orders_tb)
```

```
> head(orders_tb)
```

```
# A tibble: 6 x 7
```

	order_id	customer_name	band	order_date	product_id	product_price	quantity
	<chr>	<chr>	<chr>	<chr>	<chr>	<int>	<int>
1	A0001	Great	SKD69	2022-05-31	P0069	23000	1
2	A0002	Great	SKD69	2022-06-01	P0070	2000	2
3	A0003	Elon	SHIBA	2022-06-01	P4234	5	10000
4	A0004	Elon	SHIBA	2022-06-02	P4234	5	20000
5	A0005	Great	SKD69	2022-06-03	P4234	5	30000
6	A0006	Larry	GOOGL	2022-05-31	P0069	23000	1

# Organize: Sorting and Filtering

เรียงลำดับข้อมูลได้ด้วยคำสั่ง arrange()

Syntax: เรียงจากน้อยไปมาก (Ascending)

```
arrange(table, col)
```

```
table %>%
```

```
arrange(col)
```

Syntax: เรียงจากมากไปน้อย (Descending)

```
arrange(table, -col)
```

```
arrange(table, desc(col))
```

```
table %>%
```

```
arrange(desc(col))
```

Ex. เรียงลำดับการแสดงผลจาก quantity น้อยไปมาก

หรือ

```
> arrange(orders_tb, quantity)
# A tibble: 25 x 7
  order_id customer_name band order_date product_id product_price quantity
  <chr>    <chr>      <chr> <chr>    <chr>          <int>    <int>
1 A0001    Great      SKD69 2022-05-31 P0069          23000     1
2 A0006    Larry      GOOGL 2022-05-31 P0069          23000     1
3 A0002    Great      SKD69 2022-06-01 P0070           2000     2
4 A0018    Mark       FB555 2022-06-16 P4234           5        10
5 A0019    Mark       FB555 2022-06-17 P4234           5        20
6 A0020    Mark       FB555 2022-06-18 P4234           5        30
7 A0007    Great      SKD69 2022-06-04 P4234           5        300
8 A0003    Elon       SHIBA 2022-06-01 P4234           5       10000
9 A0011    Sergey     GOOGL 2022-06-08 P4234           5       10000
10 A0021    Sergey     GOOGL 2022-06-19 P4234           5       10000
# ... with 15 more rows
```

Ex. เรียงลำดับการแสดงผลจาก quantity มากไปน้อย

หรือ

หรือ

```
> arrange(orders_tb, desc(quantity))
# A tibble: 25 x 7
  order_id customer_name band order_date product_id product_price quantity
  <chr>    <chr>      <chr> <chr>    <chr>          <int>    <int>
1 A0010    Great      SKD69 2022-06-06 P4234           5       100000
2 A0017    Jeff       AMAZE 2022-06-15 P4234           5        70000
3 A0016    Sergey     GOOGL 2022-06-14 P4234           5        60000
4 A0009    Elon       SHIBA 2022-06-05 P4234           5        50000
5 A0015    Sergey     GOOGL 2022-06-13 P4234           5        50000
6 A0025    Sergey     GOOGL 2022-06-30 P4234           5        50000
7 A0008    Elon       SHIBA 2022-06-04 P4234           5        40000
8 A0014    Sergey     GOOGL 2022-06-12 P4234           5        40000
9 A0024    Sergey     GOOGL 2022-06-27 P4234           5        40000
10 A0005    Great      SKD69 2022-06-03 P4234           5        30000
# ... with 15 more rows
```

# Organize: Sorting and Filtering

กรองข้อมูลได้ด้วยคำสั่ง filter()

Syntax:

```
filter(table, เงื่อนไขของสิ่งที่คงอยู่)
```

```
table %>%  
  filter(เงื่อนไขของสิ่งที่คงอยู่)
```

Ex. เอาเฉพาะการสั่งซื้อสินค้าที่มีราคาเกิน 10 บาท

```
> filter(orders_tb, product_price > 10)
```

หรือ

```
> orders_tb %>% filter(product_price > 10)
```

```
# A tibble: 3 x 7
```

	order_id	customer_name	band	order_date	product_id	product_price	quantity
	<chr>	<chr>	<chr>	<chr>	<chr>	<int>	<int>
1	A0001	Great	SKD69	2022-05-31	P0069	23000	1
2	A0002	Great	SKD69	2022-06-01	P0070	2000	2
3	A0006	Larry	GOOGL	2022-05-31	P0069	23000	1

Ex. เอาเฉพาะการสั่งซื้อสินค้าโดยคนชื่อ Great

```
> filter(orders_tb, customer_name == 'Great')
```

หรือ

```
> orders_tb %>% filter(customer_name == 'Great')
```

```
# A tibble: 5 x 7
```

	order_id	customer_name	band	order_date	product_id	product_price	quantity
	<chr>	<chr>	<chr>	<chr>	<chr>	<int>	<int>
1	A0001	Great	SKD69	2022-05-31	P0069	23000	1
2	A0002	Great	SKD69	2022-06-01	P0070	2000	2
3	A0005	Great	SKD69	2022-06-03	P4234	5	30000
4	A0007	Great	SKD69	2022-06-04	P4234	5	300
5	A0010	Great	SKD69	2022-06-06	P4234	5	10000

# Organize: Summarize and Group By

สรุปข้อมูลสถิติด้วย summarize() และแบบแบ่งกลุ่มด้วย group\_by() ตามด้วย summarize()

Syntax:

```
summarize(table, mean(col))
```

```
table %>%
```

```
  summarize(min(col), max(col))
```

Ex. สรุปค่าเฉลี่ย ค่ามัธยฐาน ค่าต่ำสุดและค่าสูงสุดของ quantity

```
> summarize(orders_tb, mean(quantity), median(quantity), min(quantity), max(quantity))
```

หรือ

```
> orders_tb %>% summarize(mean(quantity), median(quantity), min(quantity), max(quantity))
# A tibble: 1 × 4
  `mean(quantity)` `median(quantity)` `min(quantity)` `max(quantity)`
  <dbl>           <int>           <int>           <int>
1      27215.         20000             1         100000
```

Ex. ตั้งชื่อคอลัมน์ของผลลัพธ์ว่า mean\_qtt, median\_qtt, min\_qtt และ max\_qtt

```
1 orders_tb %>%
2   summarize(mean_qtt = mean(quantity),
3             median_qtt = median(quantity),
4             min_qtt = min(quantity),
5             max_qtt = max(quantity))
```

```
# A tibble: 1 × 4
  mean_qtt median_qtt min_qtt max_qtt
  <dbl>     <int>     <int> <int>
1    27215.         20000         1    100000
```

# Organize: Summarize and Group By

สรุปข้อมูลสถิติด้วย summarize() และแบบแบ่งกลุ่มด้วย group\_by() ตามด้วย summarize()

Syntax:

```
table %>%  
  group_by(col1) %>%  
  summarize(sum(col2), mean(col3))
```

Ex. สรุปผลรวมของ quantity สินค้ารหัส P4234 แบ่งตาม customer\_name และเรียงผลลัพธ์จากมากไปหาน้อย

```
1 orders_tb %>%  
2   filter(product_id == 'P4234') %>%  
3   group_by(customer_name) %>%  
4   summarize(total_quantity = sum(quantity)) %>%  
5   arrange(desc(total_quantity))
```

```
# A tibble: 5 × 2  
  customer_name total_quantity  
  <chr>          <int>  
1 Sergey          310000  
2 Great           130300  
3 Elon            120000  
4 Jeff            120000  
5 Mark              60
```

# Transform: Separate, Unite and Mutate

เราแยกข้อความออกเป็นส่วน ๆ ได้ด้วย separate()

Syntax:

```
separate(table, คอลัมน์ตั้งต้น,
         into=c(รายชื่อคอลัมน์ที่แยกตัว),
         sep='ตัวคั่น')
```

Notes:

- หากต้องการเก็บผลไว้ใช้ต่อ ต้องบันทึกใส่ตัวแปรใหม่ เช่น `orders_tb <- separate(orders_tb, order_date, ...)`
- หากไม่ต้องการลบคอลัมน์ตั้งต้นทิ้ง ให้ใส่ `remove=FALSE` เพิ่มเข้าไปใน arguments

Ex. แยกคอลัมน์ order\_date ออกเป็น 3 คอลัมน์ย่อยคือ year, month และ date

```
> head(orders_tb)
# A tibble: 6 x 7
  order_id customer_name band order_date product_id product_price quantity
  <chr>    <chr>      <chr> <chr>    <chr>      <int>    <int>
1 A0001    Great      SKD69 2022-05-31 P0069      23000     1
2 A0002    Great      SKD69 2022-06-01 P0070       2000     2
3 A0003    Elon       SHIBA 2022-06-01 P4234         5    10000
4 A0004    Elon       SHIBA 2022-06-02 P4234         5    20000
5 A0005    Great      SKD69 2022-06-03 P4234         5    30000
6 A0006    Larry      GOOGL 2022-05-31 P0069      23000     1
```

```
> separate(orders_tb, order_date, into=c('year','month','date'), sep='-')
# A tibble: 25 x 9
  order_id customer_name band year month date product_id product_price quantity
  <chr>    <chr>      <chr> <chr> <chr> <chr> <chr>      <int>    <int>
1 A0001    Great      SKD69 2022 05 31 P0069      23000     1
2 A0002    Great      SKD69 2022 06 01 P0070       2000     2
3 A0003    Elon       SHIBA 2022 06 01 P4234         5    10000
4 A0004    Elon       SHIBA 2022 06 02 P4234         5    20000
5 A0005    Great      SKD69 2022 06 03 P4234         5    30000
6 A0006    Larry      GOOGL 2022 05 31 P0069      23000     1
7 A0007    Great      SKD69 2022 06 04 P4234         5     300
8 A0008    Elon       SHIBA 2022 06 04 P4234         5    40000
9 A0009    Elon       SHIBA 2022 06 05 P4234         5    50000
10 A0010    Great      SKD69 2022 06 06 P4234         5    100000
# ... with 15 more rows
```

# Transform: Separate, Unite and Mutate

เรารวมข้อความย่อยเข้าด้วยกันด้วย unite()

Syntax:

```
unite(table, 'ชื่อคอลัมน์ใหม่',  
      col1, col2, ...,  
      sep='ตัวคั่น')
```

Notes:

- หากต้องการเก็บผลไว้ใช้ต่อ ต้องบันทึกใส่ตัวแปรใหม่ เช่น `orders_tb <- unite(orders_tb, 'stage_name', ...)`
- หากไม่ต้องการลบคอลัมน์ตั้งต้นทิ้ง ให้ใส่ `remove=FALSE` เพิ่มเข้าไปใน arguments

Ex. รวมคอลัมน์ `customer_name` และ `band` คั่นด้วยเว้นวรรค (whitespace) เรียกชื่อคอลัมน์ที่รวมแล้วว่า `stage_name`

```
> head(orders_tb)  
# A tibble: 6 x 7  
  order_id customer_name band order_date product_id product_price quantity  
  <chr>    <chr>      <chr> <chr>    <chr>          <int>    <int>  
1 A0001    Great      SKD69 2022-05-31 P0069          23000     1  
2 A0002    Great      SKD69 2022-06-01 P0070           2000     2  
3 A0003    Elon       SHIBA 2022-06-01 P4234            5    10000  
4 A0004    Elon       SHIBA 2022-06-02 P4234            5    20000  
5 A0005    Great      SKD69 2022-06-03 P4234            5    30000  
6 A0006    Larry      GOOGL 2022-05-31 P0069          23000     1
```

```
> unite(orders_tb, 'stage_name', customer_name, band, sep=' ')  
# A tibble: 25 x 6  
  order_id stage_name order_date product_id product_price quantity  
  <chr>    <chr>      <chr>    <chr>          <int>    <int>  
1 A0001    Great SKD69 2022-05-31 P0069          23000     1  
2 A0002    Great SKD69 2022-06-01 P0070           2000     2  
3 A0003    Elon SHIBA 2022-06-01 P4234            5    10000  
4 A0004    Elon SHIBA 2022-06-02 P4234            5    20000  
5 A0005    Great SKD69 2022-06-03 P4234            5    30000  
6 A0006    Larry GOOGL 2022-05-31 P0069          23000     1  
7 A0007    Great SKD69 2022-06-04 P4234            5     300  
8 A0008    Elon SHIBA 2022-06-04 P4234            5    40000  
9 A0009    Elon SHIBA 2022-06-05 P4234            5    50000  
10 A0010    Great SKD69 2022-06-06 P4234            5   100000  
# ... with 15 more rows
```

# Transform: Separate, Unite and Mutate

เรานำค่าในคอลัมน์ต่าง ๆ มาแปลงแล้วสร้างเป็นคอลัมน์ใหม่ได้ด้วย mutate()

Syntax:

```
mutate(table, new_col = ...)
```

```
table %>%
```

```
mutate(new_col = ...)
```

Note: หากต้องการเก็บผลไว้ใช้ต่อ ต้องบันทึกใส่ตัวแปรใหม่ เช่น orders\_tb <- mutate(orders\_tb, revenue=...)

Ex. สร้างคอลัมน์ใหม่ชื่อ revenue ซึ่งเกิดจากการนำเอา product\_price \* quantity พร้อมทั้งบันทึกใส่ตารางเดิมด้วย

```
> orders_tb <- mutate(orders_tb, revenue = product_price * quantity)
> orders_tb
# A tibble: 25 x 8
  order_id customer_name band order_date product_id product_price quantity revenue
  <chr>    <chr>      <chr> <chr>    <chr>      <int>    <int>    <int>
1 A0001    Great      SKD69 2022-05-31 P0069      23000     1    23000
2 A0002    Great      SKD69 2022-06-01 P0070      2000     2     4000
3 A0003    Elon       SHIBA 2022-06-01 P4234         5    10000    50000
4 A0004    Elon       SHIBA 2022-06-02 P4234         5    20000   100000
5 A0005    Great      SKD69 2022-06-03 P4234         5    30000   150000
6 A0006    Larry     GOOGL 2022-05-31 P0069     23000     1    23000
7 A0007    Great      SKD69 2022-06-04 P4234         5     300     1500
8 A0008    Elon       SHIBA 2022-06-04 P4234         5    40000   200000
9 A0009    Elon       SHIBA 2022-06-05 P4234         5    50000   250000
10 A0010   Great      SKD69 2022-06-06 P4234         5   100000   500000
# ... with 15 more rows
```

Ex. เปลี่ยนคอลัมน์ revenue จาก int เป็น double (ทศนิยม) คราวนี้ลองใช้วิธี Pipe

```
> orders_tb <- orders_tb %>% mutate(revenue = as.double(revenue))
> orders_tb
# A tibble: 25 x 8
  order_id customer_name band order_date product_id product_price quantity revenue
  <chr>    <chr>      <chr> <chr>    <chr>      <int>    <int>    <dbl>
1 A0001    Great      SKD69 2022-05-31 P0069      23000     1    23000
2 A0002    Great      SKD69 2022-06-01 P0070      2000     2     4000
3 A0003    Elon       SHIBA 2022-06-01 P4234         5    10000    50000
4 A0004    Elon       SHIBA 2022-06-02 P4234         5    20000   100000
5 A0005    Great      SKD69 2022-06-03 P4234         5    30000   150000
6 A0006    Larry     GOOGL 2022-05-31 P0069     23000     1    23000
7 A0007    Great      SKD69 2022-06-04 P4234         5     300     1500
8 A0008    Elon       SHIBA 2022-06-04 P4234         5    40000   200000
9 A0009    Elon       SHIBA 2022-06-05 P4234         5    50000   250000
10 A0010   Great      SKD69 2022-06-06 P4234         5   100000   500000
# ... with 15 more rows
```

# Warning: Anscombe's quartet

Anscombe's quartet คือชุดข้อมูล 4 ชุดที่แตกต่างกันสิ้นเชิง แต่กลับให้ค่า mean, sd และ correlation ที่เหมือนกันเป๊ะ เพื่อเป็นการย้ำเตือนว่า เราไม่ควรดูแค่ค่าสรุปทางสถิติ แต่ควรดูธรรมชาติข้อมูลดิบด้วย

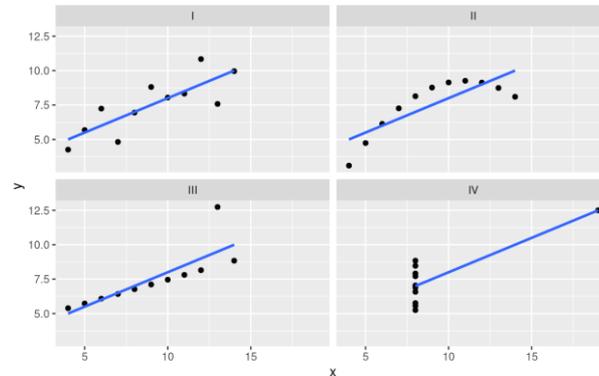
```
1 install.packages("Tmisc")
2 library("Tmisc")
3 data(quartet)
4 View(quartet)
```

set	x	y
1 I	10	8.04
2 I	8	6.95
3 I	13	7.58
4 I	9	8.81
5 I	11	8.33
6 I	14	9.96
7 I	6	7.24
8 I	4	4.26
9 I	12	10.84
10 I	7	4.82
11 I	5	5.68
12 II	10	9.14
13 II	8	8.14
14 II	13	8.74
15 II	9	8.77
16 II	11	9.26
17 II	14	8.10
18 II	6	6.13
19 II	4	3.10
20 II	12	9.13
21 II	7	7.26
22 II	5	4.74
23 III	10	7.46
24 III	8	6.77
25 III	13	12.74
26 III	9	7.11
27 III	11	7.81
28 III	14	8.84
29 III	6	6.08
30 III	4	5.39
31 III	12	8.15
32 III	7	6.42
33 III	5	5.73
34 IV	8	6.58
35 IV	8	5.76
36 IV	8	7.71
37 IV	8	8.84
38 IV	8	8.47
39 IV	8	7.04
40 IV	8	5.25
41 IV	19	12.50
42 IV	8	5.56
43 IV	8	7.91
44 IV	8	6.89

```
6 quartet %>%
7   group_by(set) %>%
8   summarize(mean(x),sd(x),mean(y),sd(y),cor(x,y))
9
```

```
# A tibble: 4 x 6
  set   `mean(x)` `sd(x)` `mean(y)` `sd(y)` `cor(x, y)`
  <fct>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
1 I       9     3.32    7.50    2.03    0.816
2 II      9     3.32    7.50    2.03    0.816
3 III     9     3.32    7.5    2.03    0.816
4 IV      9     3.32    7.50    2.03    0.817
```

```
10 ggplot(quartet,aes(x,y)) +
11   geom_point() +
12   geom_smooth(method=lm, se=FALSE) +
13   facet_wrap(~set)
```



# The Bias Function

เราดูว่าค่าจริง (Actual values) กับค่าทำนาย (Predicted values) ต่างกันมากแค่ไหนด้วยฟังก์ชัน bias()

Syntax:

```
bias(actual_vec, predicted_vec)
```

Notes:

- ถ้าผลลัพธ์เป็น 0 แปลว่า actual กับ predicted ตรงกันเปะ นั่นคือ ไม่ลำเอียง (unbiased)
- ยิ่งผลลัพธ์มีค่าต่างจาก 0 มาก แปลว่า predicted ยิ่งต่างจาก actual มาก ก็คือยิ่ง biased มาก
- ฟังก์ชัน bias() เป็นของแพ็คเกจ SimDesign ซึ่งต้อง install.packages("SimDesign") และโหลดด้วย library("SimDesign") ก่อน

Ex. ดูว่าค่าจริงกับค่าทำนายของยอด likes และ shares ในช่วง 5 วัน ต่างกันมากน้อยแค่ไหน

```
1 install.packages("SimDesign")
2 library("SimDesign")
3
4 actual_likes = c(100, 200, 300, 400, 500)
5 predicted_likes = c(100, 200, 300, 400, 500)
```

```
> bias(actual_likes, predicted_likes)
[1] 0
```

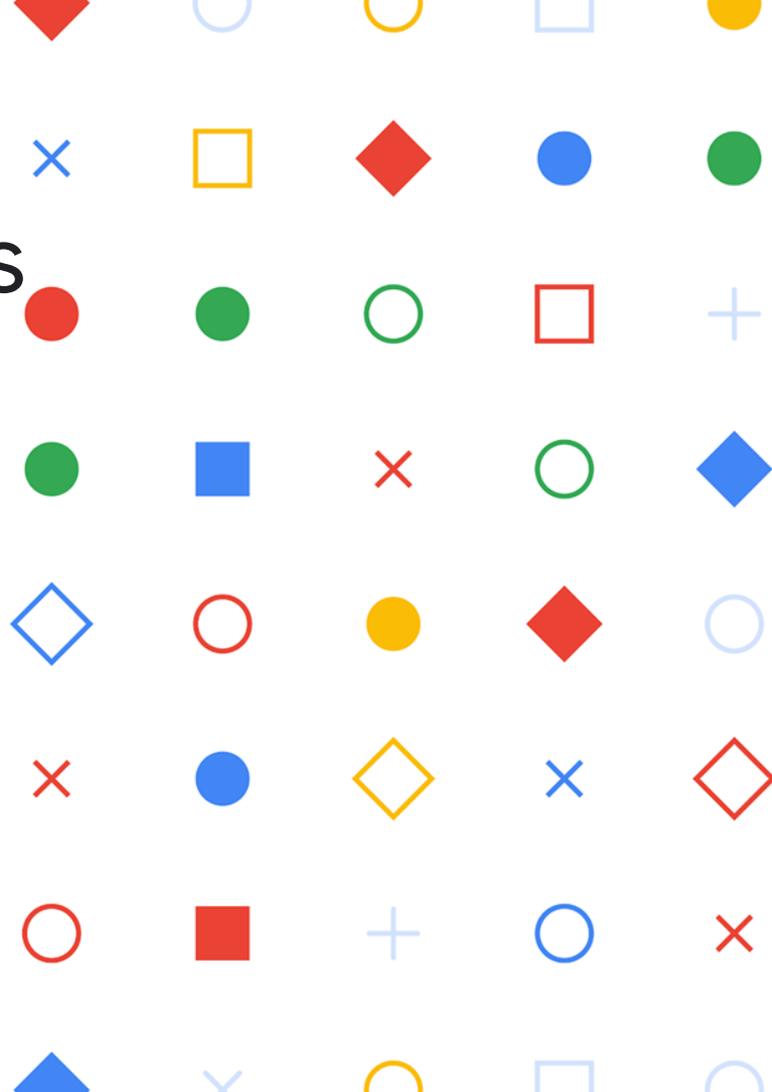
```
8 actual_shares = c(10, 20, 30, 40, 50)
9 predicted_shares = c(100, 200, 300, 400, 500)
```

```
> bias(actual_shares, predicted_shares)
[1] -270
> bias(predicted_shares, actual_shares)
[1] 270
```

# More about Visualizations, Aesthetics and Annotations

## — *Part4:*

1. Basic Components of ggplot2
2. Example: Basic Scatter Plots
3. Example: Scatter Plots with Trend Line
4. Example: Bar Charts
5. Example: Facet Functions
6. Example: Labels and Annotations
7. Saving Your Visualizations



# Basic Components of ggplot2

ggplot2 คือ library ที่ใช้สร้าง Visualizations ซึ่งมีโครงสร้างเบื้องต้นดังนี้

Syntax:

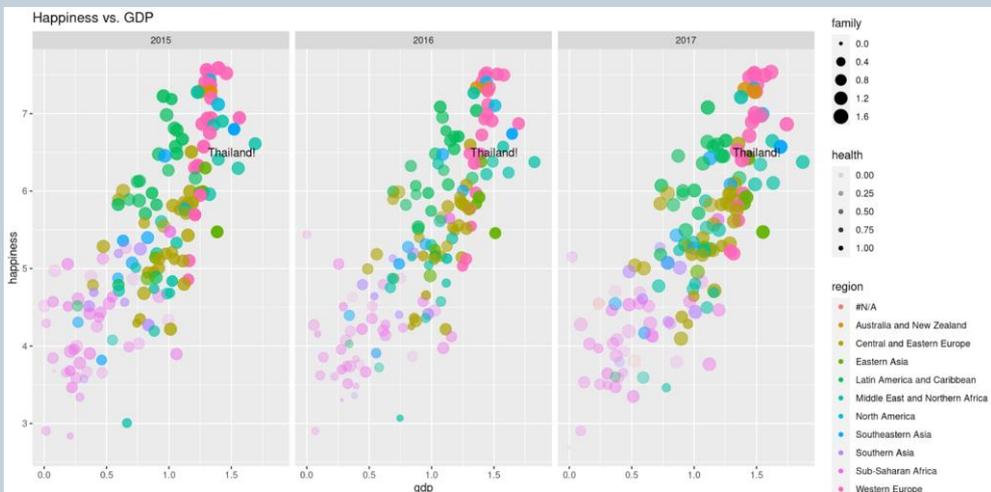
```
ggplot(data=my_table) +  
  geom_xxx(mapping=aes(x=col1,y=col2,...)) +  
  facet_wrap(~col3) +  
  labs(title='Plot Name') +  
  annotate('text',x=2,y=3,label='Google!')
```

Notes:

- geom = geometry คือ ชนิดเรขาคณิตของพล็อต เช่น geom\_point, geom\_line, geom\_bar, geom\_boxplot
- aes = aesthetics คือ "หน้าตา" ของพล็อต เช่น ตำแหน่ง (ค่า x และ y) สี (color) ความทึบ (alpha) ขนาด (size) รูปร่าง (shape)
- facet คือ แบ่งข้อมูลเป็นพล็อตย่อย (subplots) มีแบบ facet\_wrap กับ facet\_grid

Ex. สร้าง Scatter Plot หาความสัมพันธ์ระหว่าง GDP และ Happiness Score ของประเทศทั่วโลก แยกปี โดยให้สีของจุดแทนภูมิภาคของประเทศ ขนาดจุดแทนขนาดครอบครัว และความทึบของสีแทนดัชนีสุขภาพ

```
1 ggplot(data=happy_table) +  
2   geom_point(mapping=aes(x=gdp,y=happiness,color=region,alpha=health,size=family)) +  
3   facet_wrap(~year) +  
4   labs(title='Happiness vs. GDP') +  
5   annotate("text",x=1.5,y=6.5,label="Thailand!")
```



# Example: Basic Scatter Plots

สร้างด้วย `geom_point()`

Syntax:

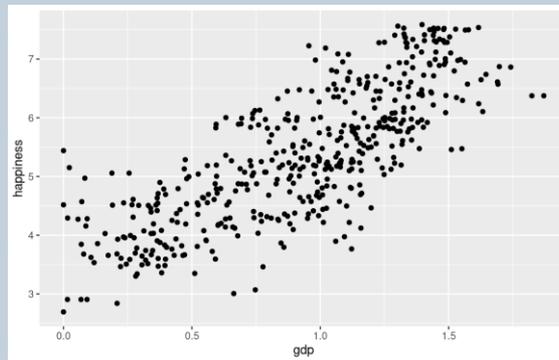
```
ggplot(data=my_table) +  
  geom_point(mapping=aes(x=col1,y=col2,...))
```

Notes:

- การเพิ่ม "Layer" ของพล็อต เราใช้เครื่องหมาย + ไม่ใช่ %>%
- ถ้าจุดมันทับกัน อาจใช้ `geom_jitter()` แทนเพื่อให้โปรแกรม "แอบเนียน" สุ่มขยับจุดให้เล็กน้อย เพื่อให้จุดไม่ทับกัน

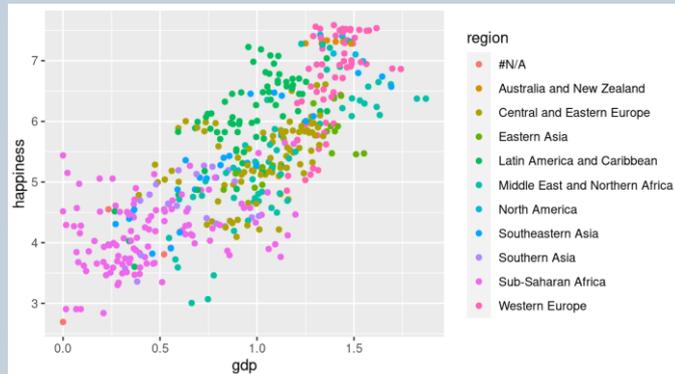
Ex. สร้าง Scatter Plot ทหาความสัมพันธ์ระหว่าง GDP และ Happiness Score

```
1 ggplot(data=happy_table) +  
2   geom_point(mapping=aes(x=gdp,y=happiness))
```



Ex. สร้าง Scatter Plot ทหาความสัมพันธ์ระหว่าง GDP และ Happiness Score ใช้สีแทนภูมิภาค และใช้ `geom_jitter()` แทนเพื่อให้จุดไม่ทับกันละ

```
1 ggplot(data=happy_table) +  
2   geom_jitter(mapping=aes(x=gdp,y=happiness,color=region))
```



# Example: Scatter Plots with Trend Line

สร้างด้วย `geom_point()` และ `geom_smooth()`

Syntax:

```
ggplot(data=table) +  
  geom_point(mapping=aes(x=c1, y=c2, ...)) +  
  geom_smooth(mapping=aes(x=c1, y=c2, ...))
```

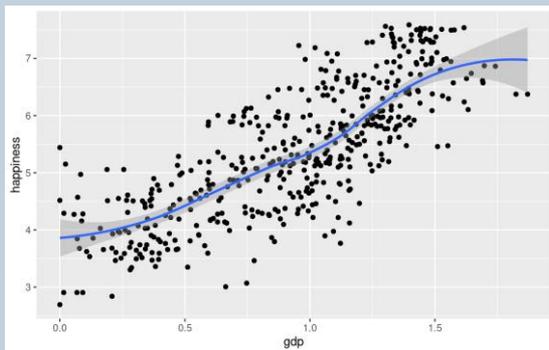
```
ggplot(data=table, mapping=aes(x=c1, y=c2)) +  
  geom_point(mapping=aes(color=...)) +  
  geom_smooth(method='lm')
```

Notes:

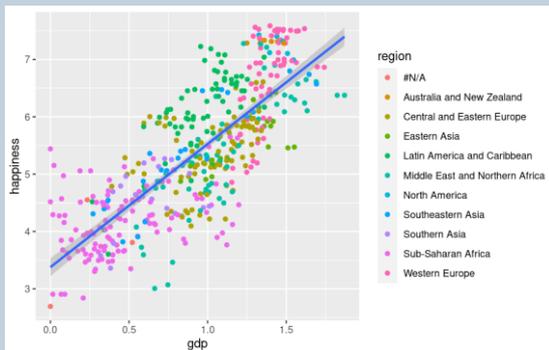
- หากต้องการเส้นตรง ต้องระบุ `method='lm'` (Linear Model)
- เอา `mapping=aes(...)` ไปใส่ใน `ggplot()` ก็ได้ถ้าใช้ค่าร่วมกันทุก Layer

Ex. สร้าง Scatter Plot และเส้น Trend Line แสดงความสัมพันธ์ระหว่าง GDP และ Happiness Score

```
1 ggplot(data=happy_table) +  
2   geom_point(mapping=aes(x=gdp, y=happiness)) +  
3   geom_smooth(mapping=aes(x=gdp, y=happiness))
```



```
1 ggplot(data=happy_table, mapping=aes(x=gdp, y=happiness)) +  
2   geom_point(mapping=aes(color=region)) +  
3   geom_smooth(method='lm')
```



# Example: Basic Bar Charts

สร้างแผนภูมิแท่งความถี่ด้วย `geom_bar()`

Syntax:

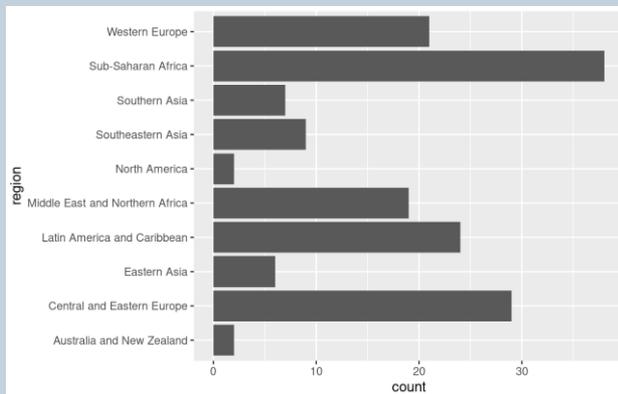
```
ggplot(data=my_table) +  
  geom_bar(mapping=aes(x=col))
```

Notes:

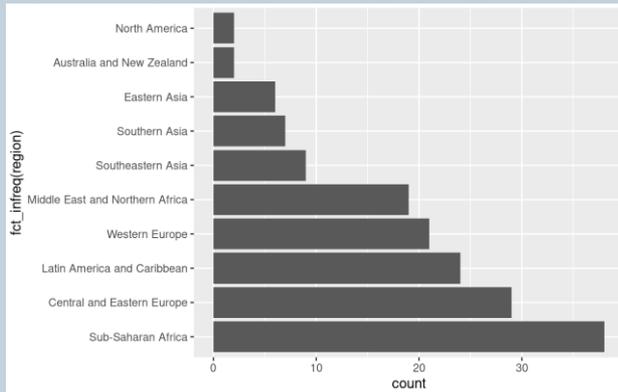
- ถ้าต้องการแท่งแนวนอน เปลี่ยนเป็น `aes(y=col)` แทน
- [Bonus] ถ้าต้องการเรียงค่าให้ใช้ฟังก์ชัน `fct_infreq()` ช่วย ซึ่งเป็นของ library ชื่อ `forcats` (มีอยู่ใน `tidyverse` แล้ว)

Ex. สร้าง Bar chart แสดงจำนวนประเทศที่เข้าร่วมในปี 2016 แบ่งตามภูมิภาค

```
1 happy_table %>%  
2 filter(year==2016) %>%  
3 ggplot() +  
4 geom_bar(mapping=aes(y=region))
```



```
1 happy_table %>%  
2 filter(year==2016) %>%  
3 ggplot() +  
4 geom_bar(mapping=aes(y=fct_infreq(region)))
```



# Example: Facet Functions

สร้างกราฟย่อย (Subplots) ด้วย `facet_wrap()` หรือ `facet_grid()`

Syntax:

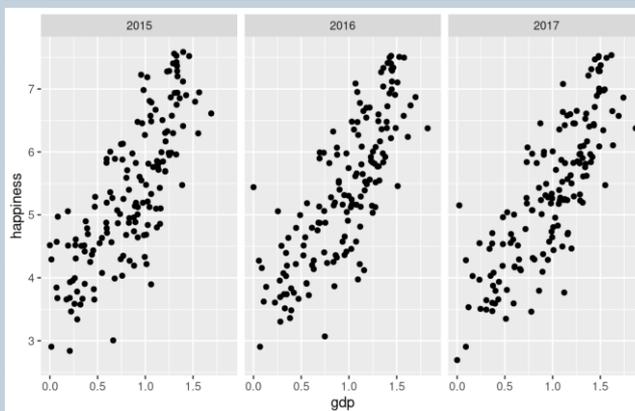
```
ggplot(data=table) +  
  geom_xxx(mapping=aes(x=c1,y=c2,...)) +  
  facet_wrap(~c3)
```

```
ggplot(data=table) +  
  geom_xxx(mapping=aes(x=c1,y=c2,...)) +  
  facet_grid(c3~c4...)
```

Note: หากต้องการใช้ `facet_wrap` แต่แบ่งย่อยมากกว่า 1 ตัวแปร ก็เติม `~col` เพิ่มได้ เช่น `facet_wrap(~year~region)`

Ex. สร้าง Scatter Plot แสดงความสัมพันธ์ระหว่าง GDP และ Happiness Score แยกตามปี

```
1 ggplot(data=happy_table) +  
2   geom_point(mapping=aes(x=gdp,y=happiness)) +  
3   facet_wrap(~year)
```



# Example: Facet Functions

สร้างกราฟย่อย (Subplots) ด้วย `facet_wrap()` หรือ `facet_grid()`

Syntax:

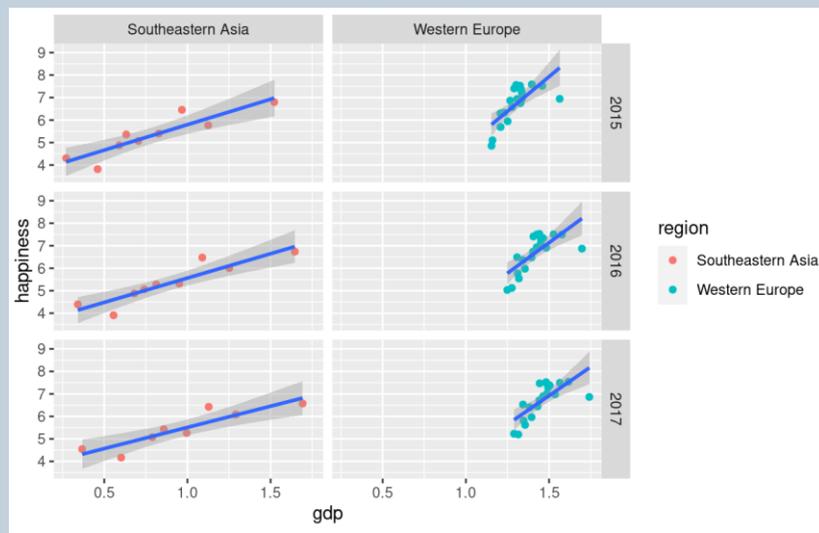
```
ggplot(data=table) +  
  geom_xxx(mapping=aes(x=c1,y=c2,...)) +  
  facet_wrap(~c3)
```

```
ggplot(data=table) +  
  geom_xxx(mapping=aes(x=c1,y=c2,...)) +  
  facet_grid(c3~c4...)
```

Note: หากต้องการใช้ `facet_wrap` แต่แบ่งย่อยมากกว่า 1 ตัวแปร ก็เติม `~col` เพิ่มได้ เช่น `facet_wrap(~year~region)`

Ex. สร้าง Scatter Plot แสดงความสัมพันธ์ระหว่าง GDP และ Happiness Score แยกตามปี

```
1 happy_table %>%  
2   filter(region %in% c('Southeastern Asia', 'Western Europe')) %>%  
3   ggplot(mapping=aes(x=gdp,y=happiness)) +  
4     geom_point(mapping=aes(color=region)) +  
5     geom_smooth(method='lm') +  
6     facet_grid(year~region)
```



# Example: Labels and Annotations

เราใส่ Labels และ Annotations ได้ด้วยฟังก์ชัน labs() และ annotate() ตามลำดับ

Syntax:

```
ggplot(data=table) +  
  geom_xxx(mapping=aes(x=c1,y=c2,...)) +  
  labs(title='aa', subtitle='bb', caption='cc') +  
  annotate('text', x=2, y=3, label='Google!')
```

Notes:

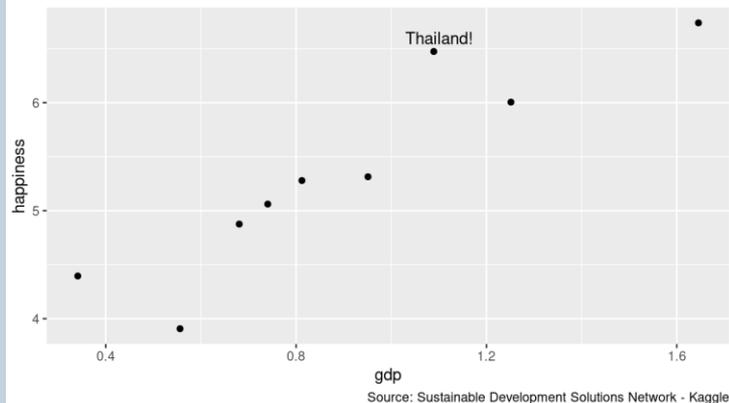
- เราสามารถปรับแต่งขนาด (size) สี (color) และรูปแบบฟอนต์ (font style) ของตัวหนังสือได้ด้วยการใส่ arguments เพิ่มเติม ลองพิมพ์ ?annotate ใน console ของ R Studio แล้วกด Enter ดูได้ครับ
- หากใครสนใจวิธีสำหรับ Annotate ชื่อประเทศของแต่ละจุดแบบครบทุกประเทศ ลองศึกษาฟังก์ชัน geom\_text() ดูได้ครับ (พิมพ์ ?geom\_text ใน console)

Ex. สร้าง Scatter Plot แสดงความสัมพันธ์ระหว่าง GDP และ Happiness Score เฉพาะประเทศใน Southeastern Asia ในปี 2016 และแสดงว่าจุดไหนคือประเทศไทย

```
# A tibble: 9 × 11  
  country    region    year happiness_rank happiness    gdp  
  <chr>      <chr>    <dbl>      <dbl>      <dbl> <dbl>  
1 Singapore Southeastern Asia 2016         22    6.74 1.65  
2 Thailand  Southeastern Asia 2016         33    6.47 1.09  
3 Malaysia Southeastern Asia 2016         47    6.00 1.25  
4 Indonesia Southeastern Asia 2016         79    5.31 0.951  
5 Philippines Southeastern Asia 2016         82    5.28 0.812  
6 Vietnam  Southeastern Asia 2016         96    5.06 0.740  
7 Laos     Southeastern Asia 2016        102    4.88 0.680  
8 Myanmar Southeastern Asia 2016        119    4.39 0.341  
9 Cambodia Southeastern Asia 2016        140    3.91 0.556
```

```
1 happy_table %>%  
2   filter(year==2016 & region=='Southeastern Asia') %>%  
3   ggplot(mapping=aes(x=gdp,y=happiness)) +  
4   geom_point() +  
5   labs(title='Happiness Score vs. GDP', subtitle='Southeast Asia 2016',  
6        caption='Source: Sustainable Development Solutions Network - Kaggle') +  
7   annotate('text', x=1.10, y=6.6, label='Thailand!')
```

Happiness Score vs. GDP  
Southeast Asia 2016



# Saving Your Visualizations

บันทึก Plots ได้ 2 วิธีคือ

1. คลิก Export -> Save as Image
2. ใช้ฟังก์ชัน ggsave()

Syntax:

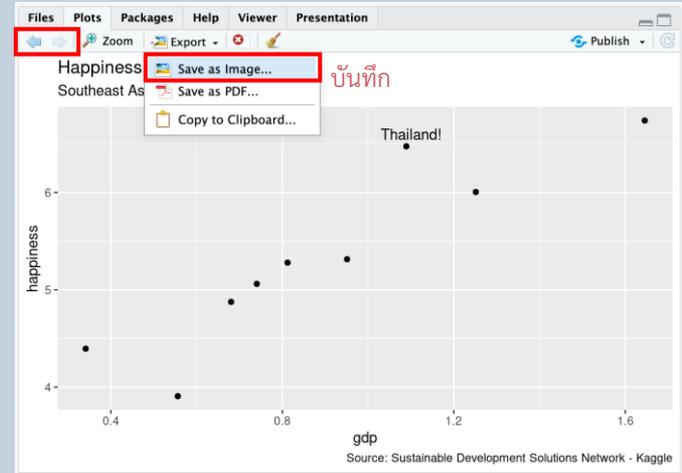
```
ggsave('filename.png', ...)
```

Notes:

- หากต้องการบันทึก Plot ในอดีตด้วย ggsave() ต้องบันทึก Plot ในอดีตไว้ตัวแปรก่อน เช่น `p <- ggplot() + ...` แล้วบันทึกด้วย `ggsave('filename.png', plot=p)`
- หากต้องการระบุสมบัติอื่นๆ ของรูป เช่น ความกว้าง (width) ความสูง (height) ความละเอียด (dpi) สามารถใส่ arguments เพิ่มเติมได้เลยครับ ศึกษาเพิ่มเติม พิมพ์ `?ggsave` ใน console ของ R Studio แล้วกด Enter

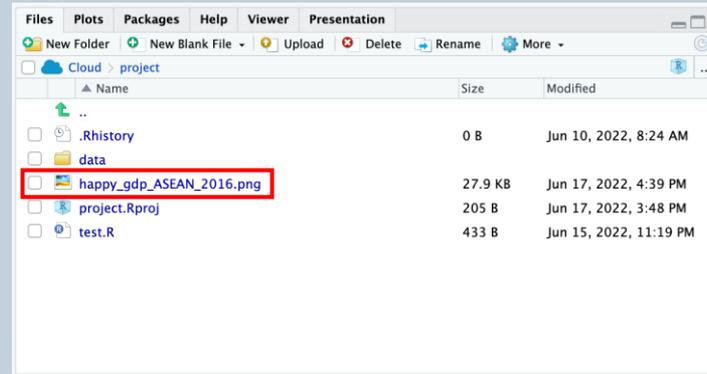
Ex. บันทึก Plot ล่าสุด ด้วยวิธี Export -> Save as Image

เลือก Plot  
ที่ต้องการ



Ex. บันทึก Plot ล่าสุด ด้วยวิธี ggsave()

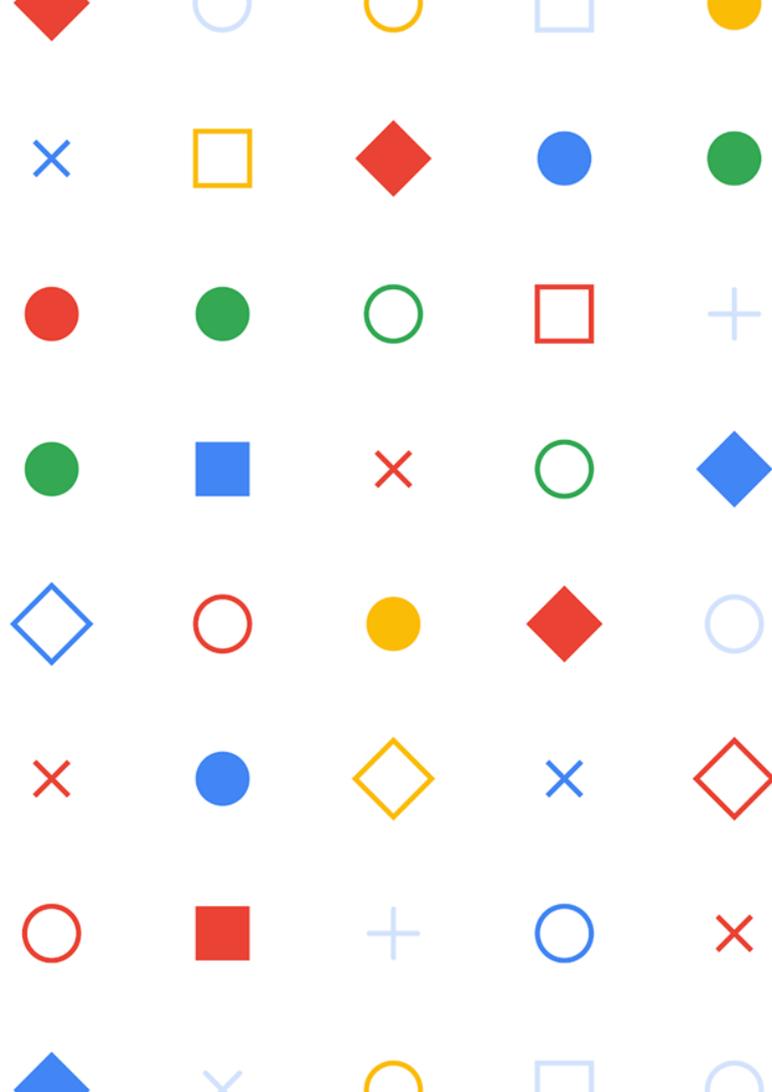
```
1 ggsave('happy_gdp_ASEAN_2016.png')
```



# Documentation and Reports

—— *Part5:*

1. Intro to R Markdown (.Rmd file)
2. R Markdown: Structure
3. R Markdown: Syntax
4. R Markdown: Code Chunks



# Intro to R Markdown

Markdown คือไวยากรณ์ (Syntax) ที่ใช้ในการปรับแต่งไฟล์เอกสารข้อความ

```
1 ---
2 title: "Thailand is Happy"
3 author: "Great"
4 date: '2022-06-17'
5 output: html_document
6 ---
7
8 ## Import Libraries
9 We need tidyverse and janitor libraries for this project.
10 ```{r Import Libraries, message=FALSE}
11 library(tidyverse)
12 library(janitor)
13 ```
14
15 ## Import Data
16
17 ```{r Import World Happiness Dataset, message=FALSE}
18 happy_table <- read_csv('./data/World_Happiness_2015_2017_ccol.csv')
19 ```
20
21 ```{r Head of Table}
22 head(happy_table)
23 ```
24
25 ## Clean Data
26 We clean the column names so they are easier to work with.
27 ```{r Clean Data}
28 happy_table <- clean_names(happy_table)
29 colnames(happy_table)
30 ```
31
32 ## Create Plot
33 Finally, we create a plot to see how happy Thailand is compared to other ASEAN countries.
34 We use the 2016 data for simplicity.
35 ```{r, echo=FALSE}
36 happy_table %>%
37   filter(year==2016 & region=='Southeastern Asia') %>%
38   ggplot(mapping=aes(x=gdp,y=happiness)) +
39   geom_point() +
40   labs(title='Happiness Score vs. GDP', subtitle='Southeast Asia 2016',
41        caption='Source: Sustainable Development Solutions Network - Kaggle') +
42   annotate('text',x=1.10,y=6.6,label='Thailand!')
43 ```
```



## Thailand is Happy

Great  
2022-06-17

### Import Libraries

We need tidyverse and janitor libraries for this project.

```
library(tidyverse)
library(janitor)
```

### Import Data

```
happy_table <- read_csv('./data/World_Happiness_2015_2017_ccol.csv')
```

```
head(happy_table)
```

```
## # A tibble: 4 x 11
##   Country Region Year Happiness Rank GDP Family Health Freedom
##   <chr> <chr> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 Switzerland Westc. 2015 7.59 1.40 1.39 0.941 0.466
## 2 Finland Westc. 2015 7.56 1.30 1.49 0.948 0.459
## 3 Denmark Westc. 2015 7.53 1.33 1.38 0.970 0.449
## 4 Norway Westc. 2015 7.50 1.46 1.29 0.995 0.430
## 5 Canada Northc. 2015 7.43 1.33 1.32 0.906 0.433
## 6 Finland Westc. 2015 7.41 1.29 1.33 0.889 0.442
## # ... with 3 more variables: Trust <dbl>, Generosity <dbl>
```

### Clean Data

We clean the column names so they are easier to work with.

```
happy_table <- clean_names(happy_table)
colnames(happy_table)
```

```
## [1] "country" "region" "year" "happiness_rank"
## [2] "happiness" "gdp" "family" "health"
## [3] "freedom" "trust" "generosity"
```

### Create Plot

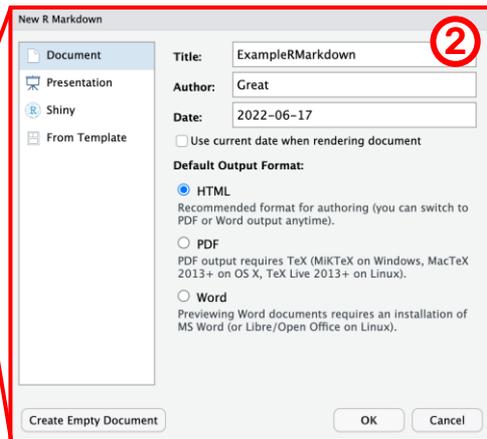
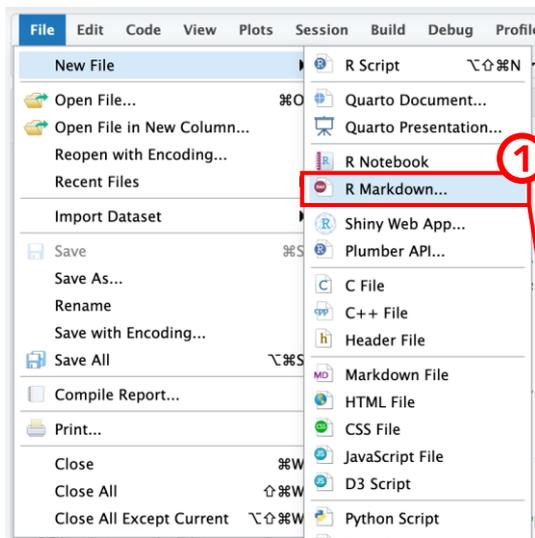
Finally, we create a plot to see how happy Thailand is compared to other ASEAN countries. We use the 2016 data for simplicity.

Happiness Score vs. GDP  
Southeast Asia 2016

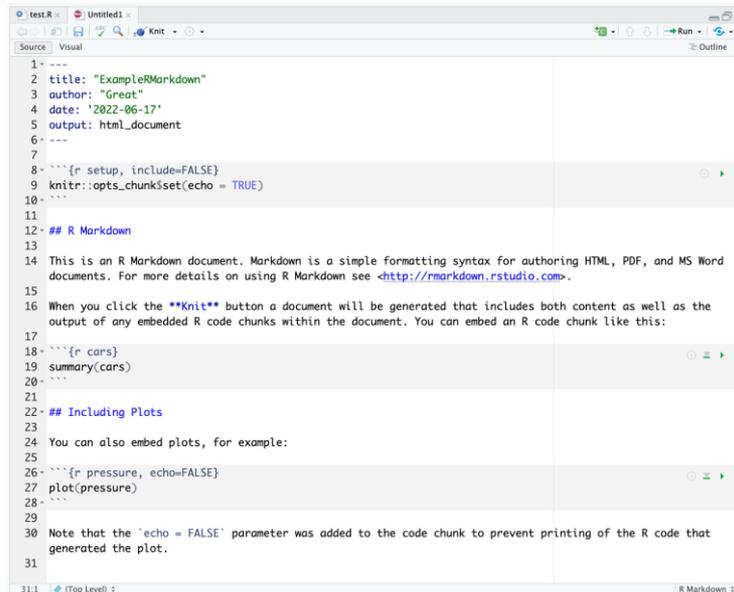
Source: Sustainable Development Solutions Network - Kaggle

# Intro to R Markdown

ใน R Studio เราสร้างไฟล์ R Markdown (.Rmd) ได้ง่าย ๆ ดังนี้

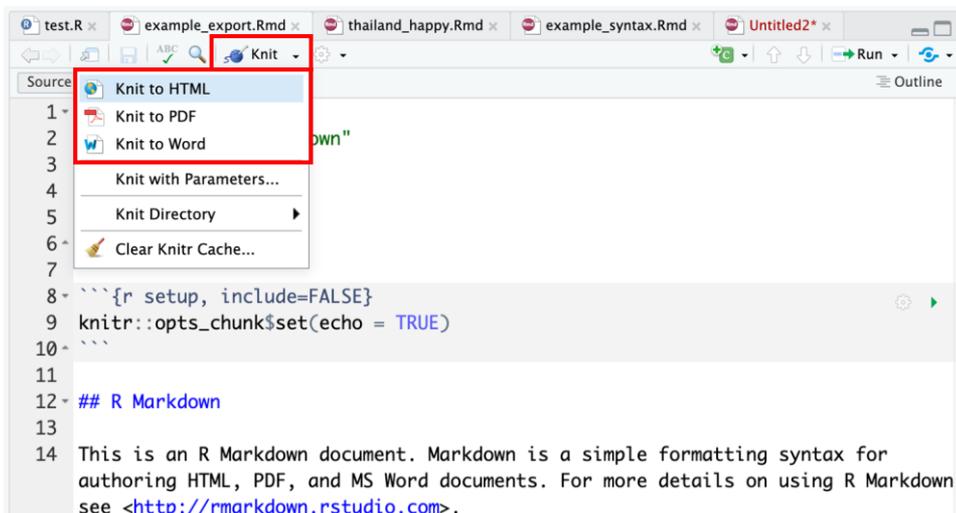


## Result



# Intro to R Markdown

สามารถ Export เป็นไฟล์สกุลต่างๆ เช่น .html .pdf .docx เป็นต้นด้วยปุ่ม Knit



```
1- 2- 3- 4- 5- 6- 7- 8- {r setup, include=FALSE}
9- knitr::opts_chunk$set(echo = TRUE)
10- 11- 12- ## R Markdown
13- 14- This is an R Markdown document. Markdown is a simple formatting syntax for
authoring HTML, PDF, and MS Word documents. For more details on using R Markdown
see <http://rmarkdown.rstudio.com>.
```



## Result

### ExampleRMarkdown

Great  
2022-06-17

#### R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

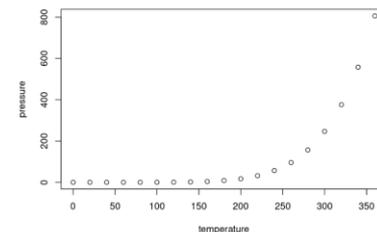
When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(mtcars)

##      speed      dist
##  Min.   4.0   Min.   1.000
## 1st Qu. 12.0   1st Qu. 16.000
##  Median 15.0   Median 19.000
##   Mean 15.4   Mean   19.298
## 3rd Qu. 19.0   3rd Qu. 26.000
##   Max. 21.0   Max.   31.000
```

#### Including Plots

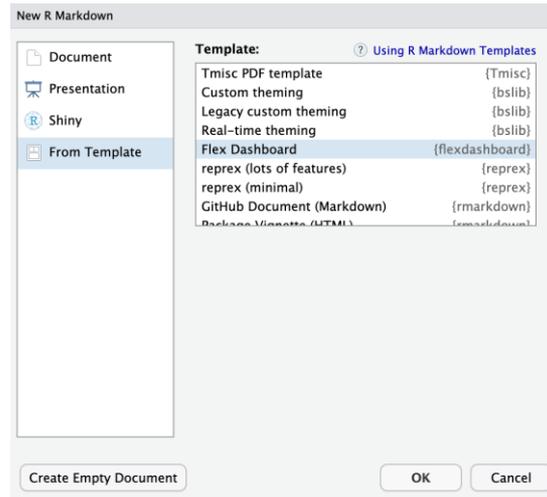
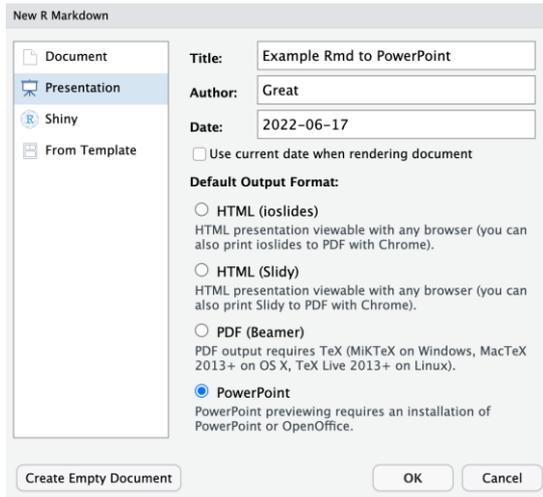
You can also embed plots, for example:



Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.

# Intro to R Markdown

เรายังสามารถ Export ไฟล์ .Rmd เป็นไฟล์ Presentations เช่น PowerPoint หรือแม้แต่เป็น Dashboard ได้ด้วย!

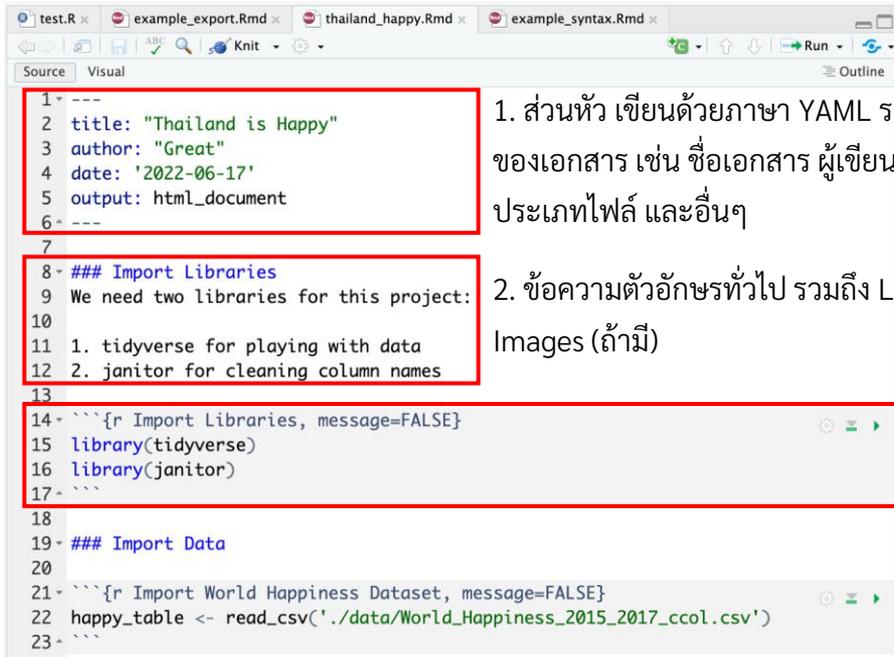


สำหรับ Flex Dashboard ต้องลง library ชื่อ "flexdashboard" ก่อนถึงจะมีให้เลือกใช้ครับ ดูตัวอย่างผลลัพธ์ได้ที่ <https://rmarkdown.rstudio.com/gallery.html>

**Trick:** เราสามารถใช้ Template หรือโครงร่างของเอกสารที่มีอยู่ก่อนแล้วได้ เพื่อประหยัดเวลาและเพิ่มความคล่องตัวในการทำงาน (Streamline the work process) เช่น Templates สำหรับรายงานยอดขายรายเดือน ใครสนใจสร้าง Template เองลองศึกษาเพิ่มเติมได้ที่ <https://bookdown.org/yihui/rmarkdown/document-templates.html> ครับ

# R Markdown: Structure

โครงสร้างเบื้องต้นของไฟล์ R Markdown มีดังนี้



```
1 ---
2 title: "Thailand is Happy"
3 author: "Great"
4 date: '2022-06-17'
5 output: html_document
6 ---
7
8 ### Import Libraries
9 We need two libraries for this project:
10
11 1. tidyverse for playing with data
12 2. janitor for cleaning column names
13
14 ```{r Import Libraries, message=FALSE}
15 library(tidyverse)
16 library(janitor)
17 ```
18
19 ### Import Data
20
21 ```{r Import World Happiness Dataset, message=FALSE}
22 happy_table <- read_csv('./data/World_Happiness_2015_2017_ccol.csv')
23 ```
```

1. ส่วนหัว เขียนด้วยภาษา YAML ระบุ Metadata ของเอกสาร เช่น ชื่อเอกสาร ผู้เขียน วันที่ผลิต ประเภทไฟล์ และอื่นๆ

2. ข้อความตัวอักษรทั่วไป รวมถึง Link และ Images (ถ้ามี)

3. ส่วนของโค้ดที่เราต้องการ Run ใส่เป็นช่วง ๆ ในเอกสาร เรียกว่า Inline Code Chunks (Inline ในที่นี้แปลว่า ใส่อยู่ท่ามกลางส่วนอื่นๆ ของเอกสารได้เลย)

# R Markdown: Syntax

เราสร้างชื่อหัวข้อ (Headers) ขนาดต่างๆ ได้ ด้วย  
เครื่องหมาย Hashtags #

Syntax:

```
1 - ---
2  title: "Example Syntax"
3  author: "Great"
4  date: '2022-06-17'
5  output: html_document
6 ^ ---
7
8 - # One Hashtag Header
9 - ## Two Hashtags Header
10 - ### Three Hashtags Header
11 - #### Four Hashtags Header
12 - ##### Five Hashtags Header
```

Result

## Example Syntax

Great

2022-06-17

### One Hashtag Header

### Two Hashtags Header

### Three Hashtags Header

Four Hashtags Header

Five Hashtags Header

Note: สังเกตว่า ยิ่งจำนวน # มาก ตัวหนังสือยิ่งเล็กลง

# R Markdown: Syntax

เราปรับแต่งรูปแบบตัวหนังสือให้เป็นตัวหนา ตัวเอียง ได้ด้วยวิธีดังนี้

Syntax:

```
16 - #### Example Styles
17
18 Normal Text
19
20 **Bold Text**
21
22 *Italic Text* or _Italic Text_
```

Result

**Example Styles**

Normal Text

**Bold Text**

*Italic Text* or Italic Text

# R Markdown: Syntax

เราสร้าง Lists แบบ Bullet Points หรือแบบเป็นตัวเลขได้ดังนี้

Notes:

- เราเรียกเครื่องหมายดอกจัน \* ภาษาอังกฤษว่า Asterisk
- สังเกตว่า ต้องมีการเว้นบรรทัด ก่อนจะขึ้นส่วนที่เป็น Bullet Points หรือ Numbered List (ดูการเว้นบรรทัดที่ 27 และ 32) ซึ่งถ้าไม่เว้น ผลลัพธ์จะไม่ออกมาเป็นแบบที่ต้องการ

Syntax:

```
24 - #### Example Lists
25
26 Bullet Points
27
28 * Point 1
29 * Point 2
30
31 Numbered List
32
33 1. Item 1
34 2. Item 2
```

Result

## Example Lists

Bullet Points

- Point 1
- Point 2

Numbered List

1. Item 1
2. Item 2

# R Markdown: Syntax

เราสร้าง Link ได้ 2 แบบ ดังนี้

Syntax:

```
36 - ### Links
37
38 <https://www.skoldio.com/>
39
40 [Skoldio is Awesome!](https://www.skoldio.com/)
```

## Result

### Links

<https://www.skoldio.com/>

[Skoldio is Awesome!](https://www.skoldio.com/)

# R Markdown: Syntax

เราแสดงรูปภาพจาก Link ไปยังรูปภาพได้ดังนี้

Syntax:

```
42 - #### Image  
43  
44 ![Skooldio Logo](https://www.skooldio.com/static/images/Skooldio_Logo.svg)
```

Result



# R Markdown: Code Chunks

เราสร้าง Code Chunks ได้ดังนี้

Notes:

- Code Chunk 1 ก้อน จะเริ่มต้นและจบลงด้วยเครื่องหมาย backtick 3 อัน ```
  - เครื่องหมายลักษณะนี้เรียกว่า "Delimiter" หรือ "ตัวคั่น" สำหรับกั้นข้อมูล (Data items) เป็นส่วน ๆ
  - เครื่องหมายนี้อยู่ตรงมุมซ้ายบนของคีย์บอร์ดด้านล่างปุ่ม Esc
- ตัว r ด้านใน { ... } บอกโปรแกรมว่าภาษาที่ใช้คือภาษา R
- My Chunk No. 1 ด้านใน { ... } เอาไว้ใช้บอกชื่อของ Chunk เฉย ๆ เพื่อความสะดวกในการค้นหา โดยจะไม่โผล่ในผลลัพธ์ที่ Export

Syntax:

```
46 - ### Inline Code Chunk
47
48 - ```{r My Chunk No.1}
49   summary(cars)
50 - ```
51
```

Result

Inline Code Chunk

```
summary(cars)
```

```
##      speed      dist
## Min.   : 4.0    Min.   : 2.00
## 1st Qu.:12.0    1st Qu.: 26.00
## Median :15.0    Median : 36.00
## Mean   :15.4    Mean   : 42.98
## 3rd Qu.:19.0    3rd Qu.: 56.00
## Max.   :25.0    Max.   :120.00
```

# R Markdown: Code Chunks

เราสามารถ Run โค้ดได้โดยตรงระหว่างกำลัง  
เขียนไฟล์ .Rmd ได้ โดยการกดปุ่ม Play

Syntax:

```
46- #### Inline Code Chunk
47-
48- ```{r My Chunk No.1}
49- summary(cars)
50- ```
51-
```

click ตรงนี้ 

Result

```
46- #### Inline Code Chunk
47-
48- ```{r My Chunk No.1}
49- summary(cars)
50- ```
```

speed	dist
Min. : 4.0	Min. : 2.00
1st Qu.:12.0	1st Qu.: 26.00
Median :15.0	Median : 36.00
Mean :15.4	Mean : 42.98
3rd Qu.:19.0	3rd Qu.: 56.00
Max. :25.0	Max. :120.00

# R Markdown: Code Chunks

เราสามารถซ่อน Code แล้วโชว์แค่ผลลัพธ์ได้ ด้วยการใช้พารามิเตอร์ `echo=FALSE`

Note: หากกดปุ่มเกียร์ตรงมุมขวาบน จะมี Options การแสดงผลให้เลือกอีกมากมาย

click ตรงนี้ 



Chunk Name: My Another Awesome Chur  
Output: Show output only  
 Show warnings  
 Show messages  
 Use paged tables  
 Use custom figure size  
[? Chunk options](#)

Syntax:

```
52 ### Plots  
53 {r My Another Awesome Chunk, echo=FALSE}  
54 plot(pressure)  
55 }  
56
```

Result

