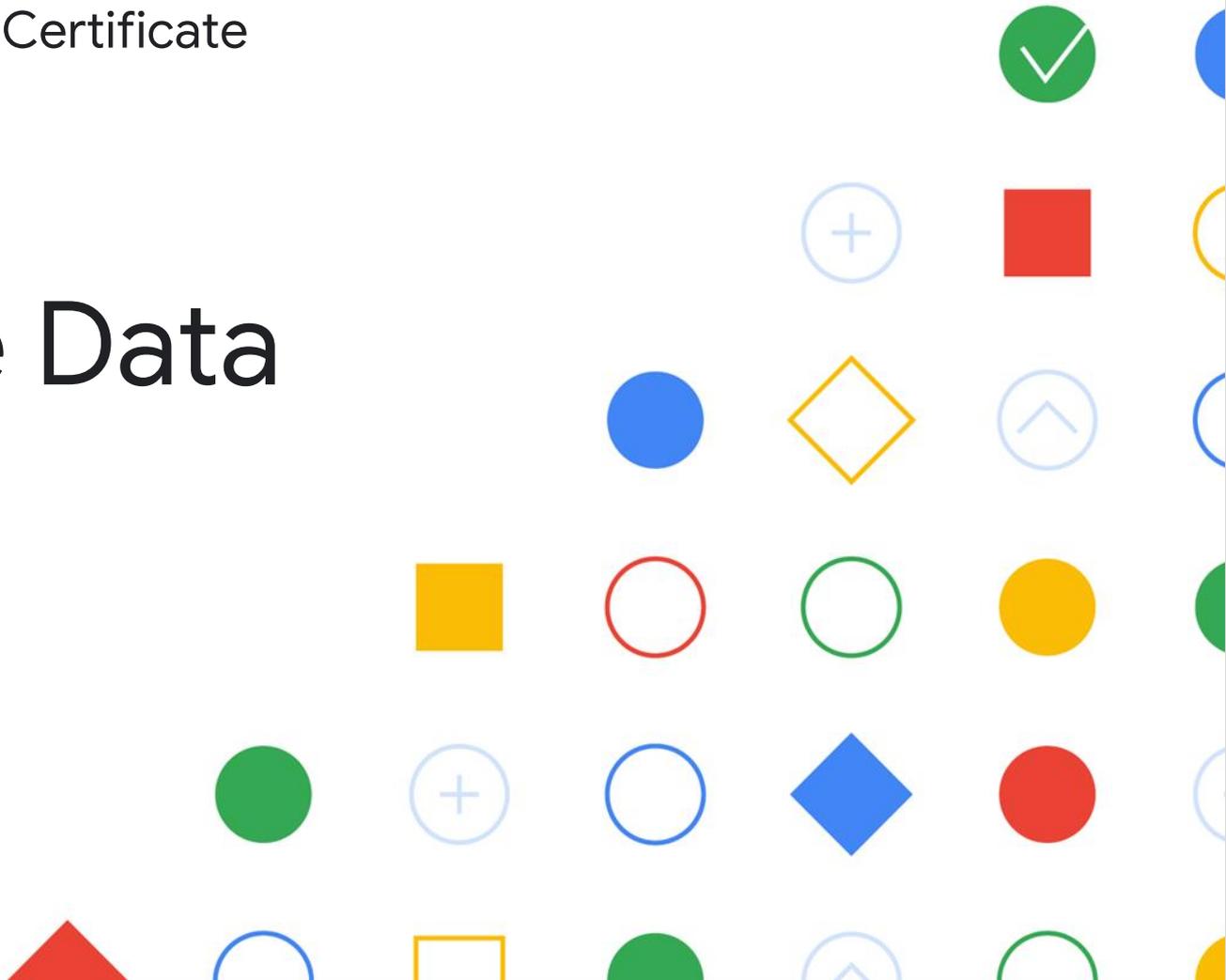


# 3. Prepare Data

*for Exploration*



# Overview:

01

Data Types and Structures

02

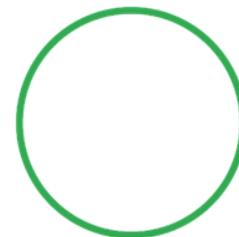
Bias, Credibility, Privacy, Ethics and Access

03

Databases: Where data lives

04

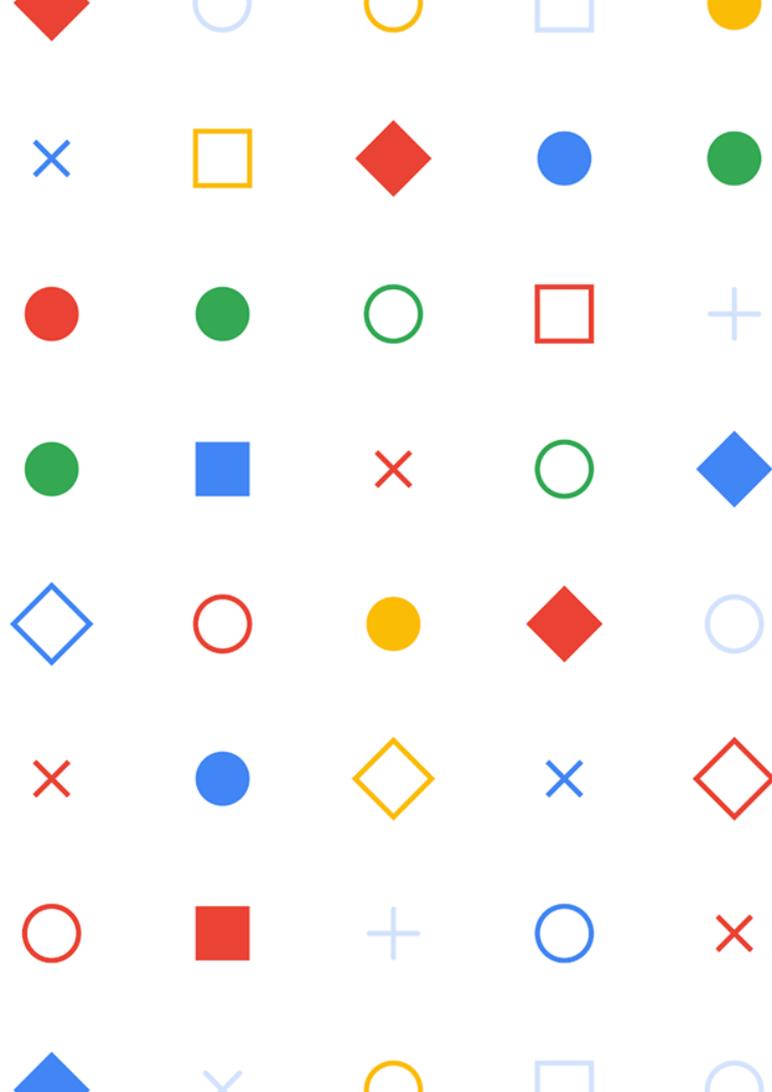
Organizing and Protecting Your Data



# Data Types and Structures

## — *Part 1:*

1. Data Collection Process
2. Data Formats
3. Data Modeling
4. Boolean Logic (True/False)
5. Transforming Data: Long vs. Wide



# Data Collection Process

## กระบวนการเก็บข้อมูล

1

### Select the right type of data

เลือกประเภทข้อมูลที่ตอบโจทย์ที่ต้องการ

2

### Determine the time frame

กำหนดกรอบเวลาของข้อมูลที่จะใช้และเวลาที่มี เช่นถ้าไม่มีเวลาเก็บข้อมูลเพิ่ม อาจใช้ข้อมูลในอดีต (Historical data)

3

need new data

### Decide how data will be collected and how much

ตัดสินใจว่าจะเก็บข้อมูลด้วยวิธีใดและเก็บปริมาณเท่าไร

- การสัมภาษณ์ (Interviews)
- การสังเกตและบันทึกผล (Observations)
- แบบฟอร์ม (Forms)
- แบบสอบถามหรือการสำรวจ (Questionnaires or Surveys)

use existing data

### Choose data sources and what data to use

เลือกแหล่งข้อมูลและข้อมูลที่จะใช้

- First-party data: เก็บข้อมูลเองด้วยทรัพยากรตนเอง
- Second-party data: เก็บข้อมูลโดยองค์กรอื่น โดยองค์กรนั้นเป็นคนเก็บเอง
- Third-party data: มาจากองค์กรที่รวบรวมข้อมูลซึ่งตัวข้อมูลมาจากแหล่งอื่นอีกที่

# Data Formats

## Formats

เก็บเองหรือไม่

เก็บเอง      **Primary data**  
คนอื่นเก็บ      **Secondary data**

เก็บไว้ที่ไหน

ภายในองค์กร      **Internal data**  
ภายนอกองค์กร      **External data**

ลักษณะค่าข้อมูล

มีค่าต่อเนื่อง      **Continuous data**  
มีค่าเป็นขั้น ๆ      **Discrete data**

เชิงปริมาณหรือคุณภาพ

เชิงปริมาณ      **Quantitative data**  
เชิงคุณภาพ      **Qualitative data**

มีลำดับชั้นหรือไม่

มีลำดับชั้น      **Ordinal data**  
ไม่มีลำดับชั้น      **Nominal data**

โครงสร้างข้อมูล

โครงสร้างชัด      **Structured data**  
โครงสร้างไม่ชัด      **Unstructured data**

## Examples

- First-party data
- Second-party and Third-party data

- ข้อมูลเงินเดือนพนักงาน ข้อมูลยอดขาย
- ข้อมูลเครดิตบูโร ข้อมูลสถิติโควิด

- น้ำหนัก ส่วนสูง อุณหภูมิ ระยะทาง
- **เงิน (เพิ่มลดที่ละ 1 สตางค์!) จำนวนคน จำนวนชั้น จำนวนวัน**

- ยอดขาย สัดส่วนลูกค้าคงเหลือ
- ยี่ห้อที่ชอบ สีที่ใช่ รสชาติที่โปรดปราน

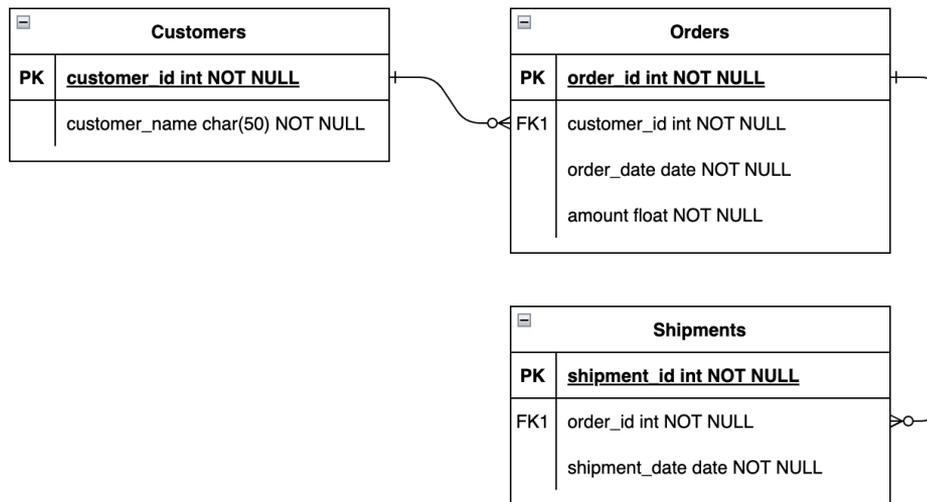
- จำนวนดาวในรีวิว ระดับความเสี่ยงกองทุน ขนาดเสื้อ
- ประเภทลูกค้า (เก่า ใหม่) เพศสภาพ

- ข้อมูลตารางทั่วไป
- รูปภาพ วิดีโอ เพลง อีเมล ข้อความโซเชียล

# Data Modeling

ใช้แผนภาพแสดงโครงสร้างและความสัมพันธ์ของข้อมูลที่มี

Ex. Entity Relationship Diagram (ER Diagram)



การนำเอา Structured data มาทำ Data model จะช่วยให้ นักวิเคราะห์ข้อมูลสามารถทำสิ่งเหล่านี้ได้ง่ายขึ้น

- บันทึก (Store)
- ค้นหา (Search)
- วิเคราะห์ (Analyze)

# Boolean Logic (True/False)

ตรรกะของข้อมูลประเภทจริง/เท็จ (True/False, Yes/No, 1/0)

เราสามารถเอาประโยคจริง/เท็จ มาผสมกันด้วยตัวดำเนินการ "และ" (AND) "หรือ" (OR) กับ "ไม่" (NOT)

	Rule	Examples
 <p>AND</p>	ต้องจริงทั้งคู่ ถึงจะจริง	<pre>CASE WHEN object = 'shirt' AND color = 'red' THEN 'buy' ELSE 'ignore' END</pre> <p>   </p>
 <p>OR</p>	จริงอย่างใดอย่างหนึ่งก็พอ	<pre>CASE WHEN object = 'shirt' OR color = 'red' THEN 'buy' ELSE 'ignore' END</pre> <p>   </p>
 <p>NOT</p>	เปลี่ยนเป็นตรงข้าม	<pre>CASE WHEN NOT object = 'shirt' THEN 'buy' ELSE 'ignore' END</pre> <p>   </p>

# Transforming Data: Long vs. Wide

ปรับโครงสร้างข้อมูลจากแบบ "ยาว" (Long) เป็นแบบ "กว้าง" (Wide)

## Long Data

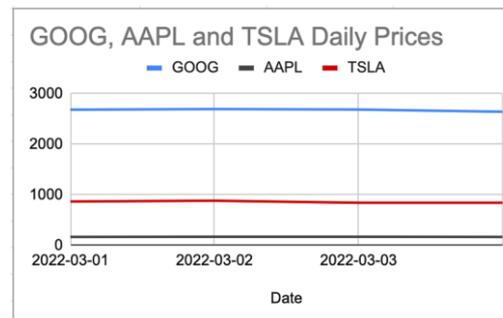
Symbol	Date	Close
GOOG	2022-03-01	2683.36
GOOG	2022-03-02	2695.03
GOOG	2022-03-03	2686.16
GOOG	2022-03-04	2642.44
AAPL	2022-03-01	163.2
AAPL	2022-03-02	166.56
AAPL	2022-03-03	166.23
AAPL	2022-03-04	163.17
TSLA	2022-03-01	864.37
TSLA	2022-03-02	879.89
TSLA	2022-03-03	839.29
TSLA	2022-03-04	838.29

- เก็บข้อมูลปริมาณมากๆ ได้สะดวกกว่า (ลองนึกว่ามีหุ้นเป็นพัน ๆ ตัว หรือมีหุ้นตัวใหม่เข้ามาเพิ่ม)
- นำไปวิเคราะห์ข้อมูลแบบแบ่งกลุ่มได้ง่ายกว่า เช่น ค่าสถิติต่างๆ ของราคาหุ้นแต่ละตัว

## Wide Data

Date	GOOG	AAPL	TSLA
2022-03-01	2683.36	163.2	864.37
2022-03-02	2695.03	166.56	879.89
2022-03-03	2686.16	166.23	839.29
2022-03-04	2642.44	163.17	838.29

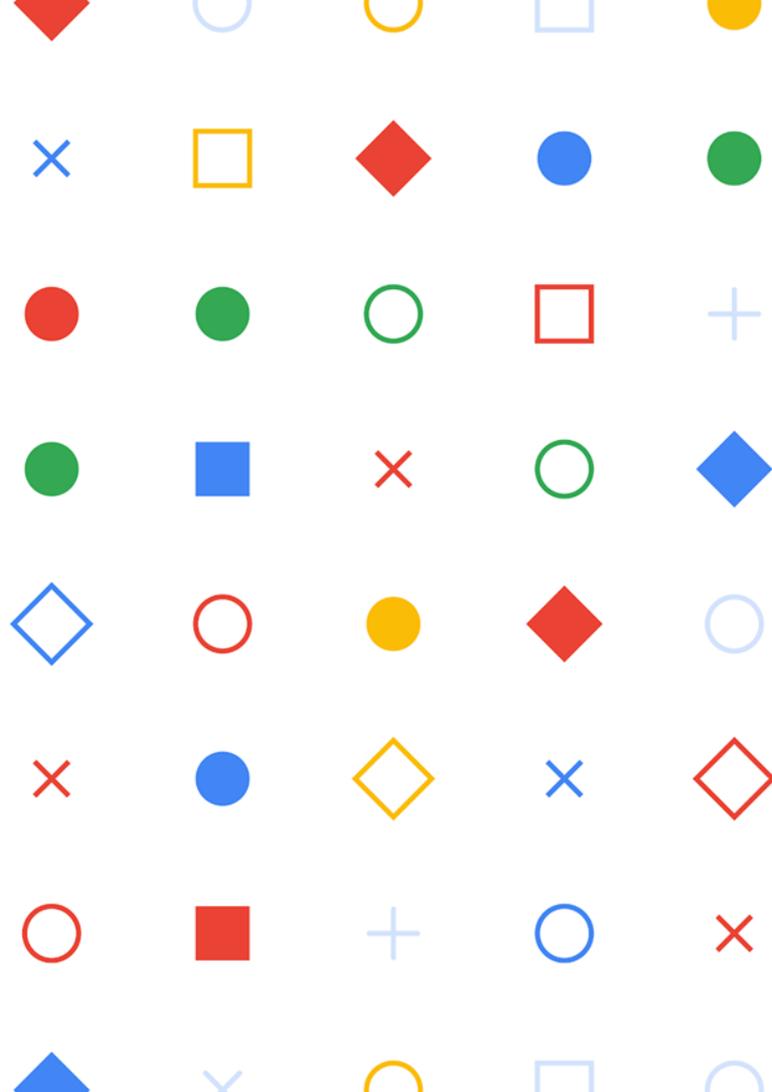
- อ่านและเปรียบเทียบค่าได้ง่ายกว่า
- สร้างกราฟของแต่ละคอลัมน์ได้ง่ายกว่า



# Bias, Credibility, Ethics, Privacy and Access

— *Part2:*

1. Four Types of Data Bias
2. Five Qualities of Good Data Sources
3. Six Aspects of Data Ethics
4. Open Data Standards



# 4 Types of Data Bias

ความลำเอียงเกี่ยวกับข้อมูลมีหลัก ๆ อยู่ 4 แบบคือ



# 5 Qualities of Good Data Sources

## 5 ลักษณะของแหล่งข้อมูลที่ดี

1

### Reliable

มาจากแหล่งข้อมูลที่น่าเชื่อถือ ว่าเก็บข้อมูลอย่างแม่นยำ ครบถ้วน ไม่ลำเอียง

2

### Original

ไม่ได้ไปลอกใครมาอีกที (พยายามหลีกเลี่ยง Third-party data ถ้าทำได้)

3

### Comprehensive

มีข้อมูลสำคัญครบถ้วนในการตอบโจทย์

4

### Current

ไม่ล้าสมัย

5

### Cited

มีการอ้างอิงถึงเจ้าของข้อมูลอย่างชัดเจนและถูกต้อง (ใครสร้างข้อมูลชุดนี้? มาจากองค์กรที่น่าเชื่อถือไหม?)

# 6 Aspects of Data Ethics

## 6 ลักษณะของจริยธรรมในการเก็บ ใช้และเผยแพร่ข้อมูล

1

### Ownership

คนให้ข้อมูลคือเจ้าของข้อมูล โดยเป็นผู้กำหนดว่าจะให้ใช้และแชร์ข้อมูลอย่างไรได้บ้าง

2

### Transaction transparency

ทุกกระบวนการประมวลผลข้อมูลต้องได้รับการอธิบายให้เจ้าของข้อมูลเข้าใจ

3

### Consent

แต่ละคนมีสิทธิที่จะทราบว่าข้อมูลส่วนบุคคลจะถูกนำไปใช้ทำอะไรบ้างและเพราะเหตุใดก่อนตกลงที่จะให้ข้อมูล

4

### Currency

แต่ละคนควรได้รับทราบว่าข้อมูลส่วนบุคคลของเขาจะถูกนำไปใช้หารายได้หรือไม่ และมากน้อยแค่ไหน

5

### Privacy

ปกป้องความเป็นส่วนตัวของข้อมูลส่วนบุคคลในทุกขั้นตอน

6

### Openness

ข้อมูล que ควรเปิดเผยต่อสาธารณะ ควรได้รับการเปิดเผยและเข้าถึงโดยสาธารณชน



## Data Anonymization

ปรับข้อมูล que สู่ถึงตัวบุคคล (Personal Identifiable Information: PII)

เช่น ชื่อ เบอร์โทร รหัสบัตรประชาชน ฯลฯ ให้เป็นนิรนาม

ด้วยเทคนิคต่าง ๆ เช่น การเว้นค่า (Blanking)

การแฮช (Hashing) การปกปิดค่า (Masking)

# Open Data Standards

1

**Accessible to the public as a complete dataset**

คนทั่วไปสามารถเข้าถึงข้อมูลได้ทั้งหมด

2

**Allow to be reused and redistributed**

อนุญาตให้นำไปใช้และเผยแพร่ซ้ำได้

3

**Allow universal participation**

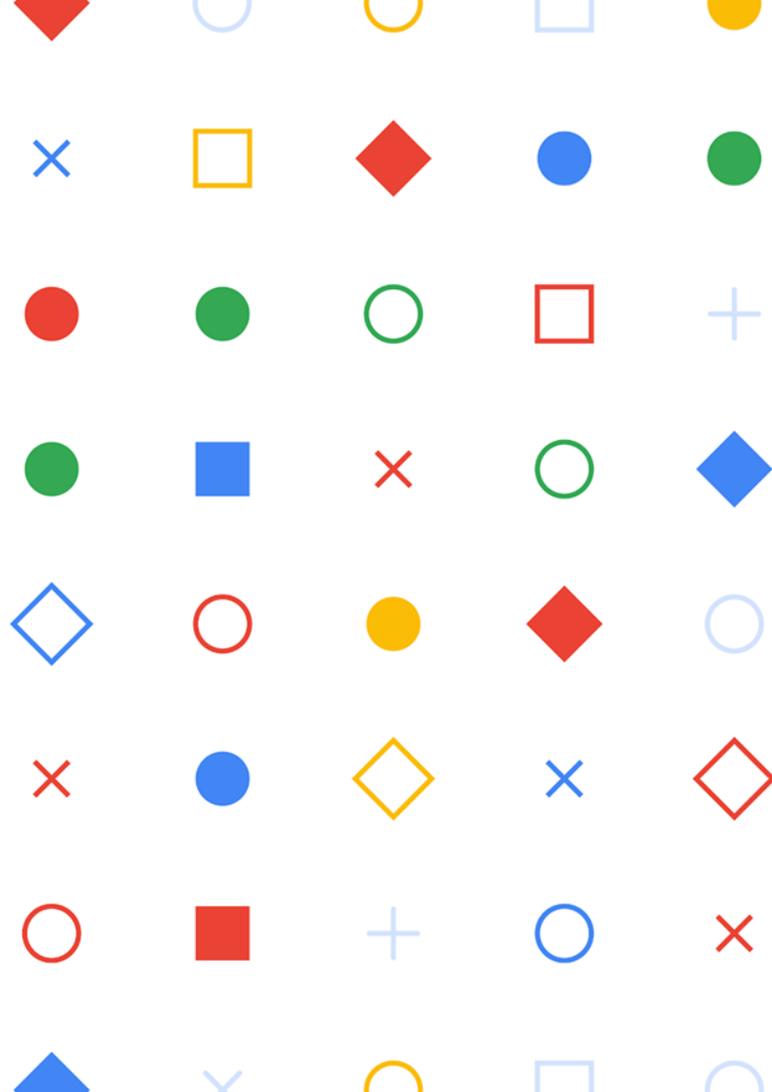
เปิดให้คนทุกกลุ่มเข้าถึงข้อมูลได้ ไม่กีดกัน (Discriminate) คนกลุ่มใดกลุ่มหนึ่ง

**Note:** Open Data ช่วยกระตุ้นความร่วมมือระหว่างองค์กรต่าง ๆ ได้ เช่น วงการวิทยาศาสตร์ วงการสาธารณสุข แต่ความท้าทายอย่างหนึ่งของ Open Data คือถ้ามีข้อมูลจากหลายแหล่ง แหล่งข้อมูลควรจะสามารถเชื่อมโยงและแลกเปลี่ยนข้อมูลได้อย่างราบรื่น (Data Interoperability)

# Databases: Where data lives

— *Part3:*

1. Relational Database
2. Metadata
3. Importing Data
4. Sorting and Filtering Data
5. BigQuery and more SQL knowledge



# Relational Database

ฐานข้อมูลที่ประกอบด้วยตารางที่มีความสัมพันธ์กันและสามารถเชื่อมต่อกันได้ด้วย "Keys" มี 2 ประเภทคือ



## Primary key

คีย์หลัก: ในตารางหนึ่ง มีได้เพียง 1 อันและห้ามมีค่าซ้ำกัน (Unique)



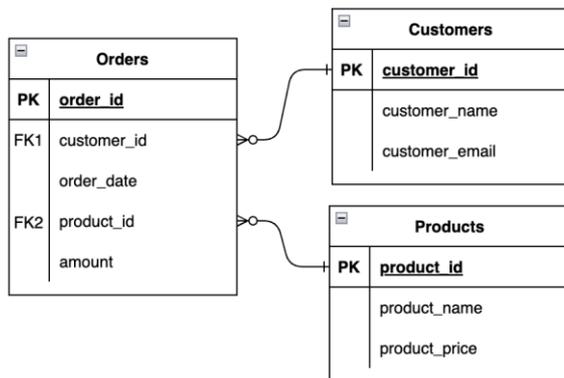
## Foreign key

คีย์นอก: ในตารางหนึ่ง มีได้มากกว่า 1 อันและมีค่าซ้ำกันได้

Ex. ฐานข้อมูลร้านค้าออนไลน์อย่างง่าย

### ตารางรายการสั่งซื้อ

- order\_id ไม่ซ้ำกัน (1 การซื้อ 1 order\_id)
- ลูกค้า 1 คนซื้อได้หลายครั้ง ดังนั้นอาจมีหลายแถวที่มี customer\_id เดียวกัน



### ตารางรายการซื้อลูกค้า

- สังเกตว่า customer\_id เป็น primary key ของตารางนี้ แต่เป็น foreign key ของตารางรายการสั่งซื้อ

# Relational Database

เราอาจต้องทำการ "Normalize" ตารางเพื่อไม่ให้เกิดการเก็บข้อมูลเดียวกันซ้ำซ้อนเปลืองพื้นที่

order_id	customer_id	customer_name	customer_email	order_date	product_id	product_name	product_price	amount
A0001	C1234	Great	ggg@cher.com	2022-05-31	P0069	Playstation 5	23,000	1
A0002	C1234	Great	ggg@cher.com	2022-06-01	P0070	PS5 Joystick	2,000	2
A0003	C5678	Elon	doge@moon.co	2022-06-01	P4234	Lactasoy 125ml	5	10,000
A0004	C5678	Elon	doge@moon.co	2022-06-02	P4234	Lactasoy 125ml	5	20,000
A0005	C1234	Great	ggg@cher.com	2022-06-03	P4234	Lactasoy 125ml	5	30,000



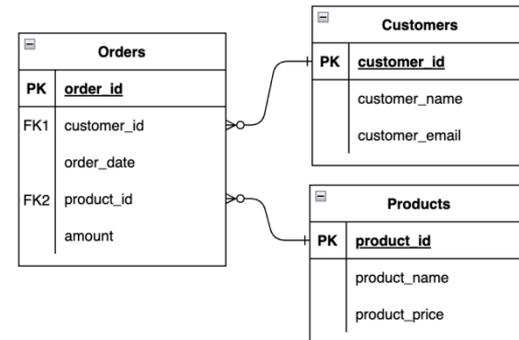
Normalize

order_id	customer_id	order_date	product_id	amount
A0001	C1234	2022-05-31	P0069	1
A0002	C1234	2022-06-01	P0070	2
A0003	C5678	2022-06-01	P4234	10,000
A0004	C5678	2022-06-02	P4234	20,000
A0005	C1234	2022-06-03	P4234	30,000

customer_id	customer_name	customer_email
C1234	Great	ggg@cher.com
C5678	Elon	doge@moon.co

product_id	product_name	product_price
P0069	Playstation 5	23,000
P0070	PS5 Joystick	2,000
P4234	Lactasoy 125ml	5

Orders	
PK	order_id
	customer_id
	customer_name
	customer_email
	order_date
	product_id
	product_name
	product_price
	amount



# Metadata

ข้อมูลสำหรับอธิบายสมบัติต่าง ๆ ของข้อมูลดิบอีกทีหนึ่ง (Data about data) มี 3 ประเภทคือ

## Definition

## Examples

### Descriptive Metadata

ใช้บรรยายข้อมูลทั่วไปของสิ่ง ๆ หนึ่ง  
เพื่อการสืบค้นในอนาคต

- หมายเลข ISBN ของหนังสือ  
(ด้านในมีข้อมูลต่างๆ เช่น ชื่อหนังสือ คนเขียน ฯลฯ)
- รหัสสินค้า (ด้านในมีข้อมูลชื่อสินค้า ชื่อผู้ผลิต ราคาขาย ฯลฯ)

### Structural Metadata

ใช้ระบุว่าข้อมูลถูกจัดระเบียบอย่างไร และ  
เป็นสมาชิกของชุดข้อมูลใดบ้าง

- สารบัญของหนังสือ (หน้าต่างๆ ถูกรวมเป็นบทอย่างไร  
และแต่ละบทเรียงกันเป็นโครงสร้างจากต้นไปท้ายอย่างไร)

### Administrative Metadata

ใช้ระบุข้อมูลทางเทคนิค  
ของสินทรัพย์ดิจิทัล (Digital asset)

- ข้อมูลสิทธิการนำสินทรัพย์ไปใช้ต่อ เช่น Creative Commons License
- ข้อมูลวันที่รูปถูกถ่าย รวมถึงสกุลไฟล์
- ข้อมูลสิทธิและประวัติการแก้ไขไฟล์ Google Doc

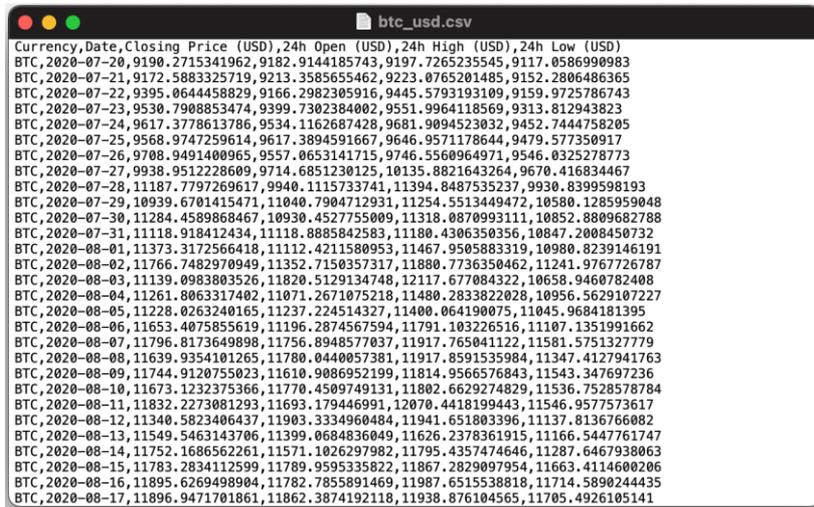
# Metadata

สิ่งที่ควรรู้เพิ่มเติมเกี่ยวกับ Metadata สำหรับ Data Analysts มีดังนี้

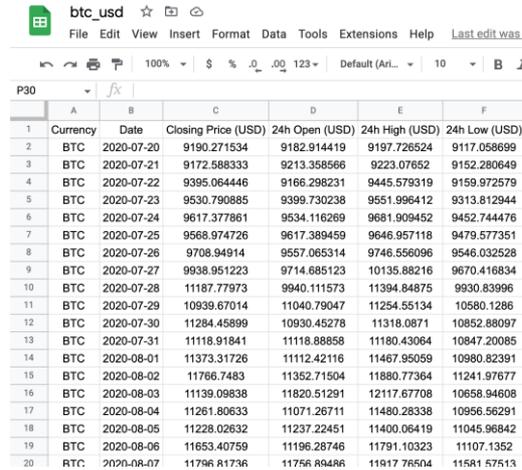
- การศึกษา Metadata ก่อนวิเคราะห์ข้อมูล จะช่วยให้เราประเมินคุณภาพและความน่าเชื่อถือของข้อมูลได้ในเบื้องต้น
- หากข้อมูลมีหลายส่วน เราควรมีระบบจัดการ Metadata แบบรวมศูนย์ ซึ่งอาจอยู่ในรูปฐานข้อมูลที่รวบรวม Metadata ทั้งหมดของชุดข้อมูลต่างๆ ไว้อย่างเป็นระเบียบ ซึ่งเรียกว่า **Metadata Repository**
- เราเรียกกระบวนการในการจัดการข้อมูลของบริษัทอย่างเป็นทางการว่า **Data Governance**

# Importing Data

ก่อนจะเริ่มวิเคราะห์ข้อมูล เราต้องนำเข้าข้อมูล (Import data) ก่อน ซึ่งสำหรับข้อมูลตาราง หนึ่งในสกุลไฟล์ที่นิยมคือ CSV (Comma-Separated Values) ซึ่งก็คือข้อมูลต่างๆ ถูกคั่น (Delineated) ด้วยเครื่องหมายลูกน้ำ ", " (Comma)



```
currency,date,closing price (usd),24h open (usd),24h high (usd),24h low (usd)
BTC,2020-07-20,9190.2715341962,9182.9144185743,9197.7265235545,9117.0586999893
BTC,2020-07-21,9172.5883325719,9213.3585655462,9223.0765201485,9152.2806486365
BTC,2020-07-22,9395.0644458829,9166.2982305916,9445.5793193109,9159.9725786743
BTC,2020-07-23,9530.7908853474,9399.7302384002,9551.9964118569,9313.812943823
BTC,2020-07-24,9617.3778613786,9534.1162687428,9681.9094523032,9452.7444758205
BTC,2020-07-25,9568.9747259614,9617.3894591667,9646.9571178644,9479.577350917
BTC,2020-07-26,9708.9491400695,9557.0653141715,9746.5560964971,9546.8325278773
BTC,2020-07-27,9938.9512228609,9714.6851230125,10135.8821643264,9670.416834467
BTC,2020-07-28,11187.797269617,9940.1115733741,11394.8487535237,9930.8399598193
BTC,2020-07-29,10939.6701415471,11040.7904712931,11254.5513449472,10580.1285959048
BTC,2020-07-30,11284.4589868467,10930.4527755009,11318.0870993111,10852.8809682788
BTC,2020-07-31,11118.918412434,11118.8885842583,11180.4306350356,10847.2000450732
BTC,2020-08-01,11373.3172566418,11112.4211580953,11467.9505883319,10980.8239146191
BTC,2020-08-02,11766.7482970949,11352.7150357317,11800.677084322,11241.9767726787
BTC,2020-08-03,11139.0983803526,11820.5129134748,12117.677084322,10658.9460782408
BTC,2020-08-04,11261.8063317402,11071.2671075218,11480.2833822028,10956.5629107227
BTC,2020-08-05,11228.0263240165,11237.224514327,11400.064190075,11045.9684181395
BTC,2020-08-06,11653.4075855619,11196.2874567594,11791.103226516,11107.1351991662
BTC,2020-08-07,11796.8173649898,11756.8948577037,11917.765041122,11581.5751327779
BTC,2020-08-08,11639.9354101265,11700.0440057381,11917.8591535984,11347.4127941763
BTC,2020-08-09,11744.9120755023,11610.9086952199,11814.9566576843,11543.347697236
BTC,2020-08-10,11673.1232375366,11770.4509749131,11802.6629274829,11536.7528578784
BTC,2020-08-11,11832.2273081293,11693.179446991,12070.4418199443,11546.547761747
BTC,2020-08-12,11340.5823406437,11903.3334960484,11941.651803396,11137.8136766082
BTC,2020-08-13,11549.5463143706,11399.0684836049,11626.2378361915,11166.5447761747
BTC,2020-08-14,11752.1686562261,11571.1026297982,11795.4357474646,11287.6467938063
BTC,2020-08-15,11783.2834112599,11789.9595335822,11867.2829097954,11663.4114600206
BTC,2020-08-16,11895.6269498904,11782.7855091469,11987.6515538818,11714.5890244435
BTC,2020-08-17,11896.9471701861,11862.3874192118,11938.876104565,11705.4926105141
```



	A	B	C	D	E	F
1	Currency	Date	Closing Price (USD)	24h Open (USD)	24h High (USD)	24h Low (USD)
2	BTC	2020-07-20	9190.271534	9182.914419	9197.726524	9117.058699
3	BTC	2020-07-21	9172.588333	9213.358566	9223.07652	9152.280649
4	BTC	2020-07-22	9395.064446	9166.298231	9445.579319	9159.972579
5	BTC	2020-07-23	9530.790885	9399.730238	9551.996412	9313.812944
6	BTC	2020-07-24	9617.377861	9534.116269	9681.909452	9452.744476
7	BTC	2020-07-25	9568.974726	9617.389459	9646.957118	9479.577351
8	BTC	2020-07-26	9708.94914	9557.065314	9746.556096	9546.032528
9	BTC	2020-07-27	9938.951223	9714.685123	10135.88216	9670.416834
10	BTC	2020-07-28	11187.77973	9940.111573	11394.84875	9930.83996
11	BTC	2020-07-29	10939.67014	11040.79047	11254.55134	10580.1286
12	BTC	2020-07-30	11284.45899	10930.45278	11318.0871	10852.88097
13	BTC	2020-07-31	11118.91841	11118.88858	11180.43064	10847.20085
14	BTC	2020-08-01	11373.31726	11112.42116	11467.95059	10980.82391
15	BTC	2020-08-02	11766.7483	11352.71504	11880.77364	11241.97677
16	BTC	2020-08-03	11139.09838	11820.51291	12117.67708	10658.94608
17	BTC	2020-08-04	11261.80633	11071.26711	11480.28338	10956.56291
18	BTC	2020-08-05	11228.02632	11237.22451	11400.06419	11045.98842
19	BTC	2020-08-06	11653.40759	11196.28746	11791.10323	11107.1352
20	BTC	2020-08-07	11796.81736	11756.89486	11917.76504	11581.57513

Note: หากไฟล์ข้อมูลมีขนาดใหญ่มาก ๆ อาจไม่เหมาะกับการเปิดใน Google Sheets หรือ Excel

# Sorting and Filtering

การเรียงลำดับ (Sorting) มี 2 แบบ: จากน้อยไปมาก (Ascending: A-Z) และจากมากไปน้อย (Descending: Z-A)

การกรองค่า (Filtering): กรองเอาแถวที่มีค่าที่สอดคล้องเงื่อนไขคงไว้ในตารางและแถวที่มีค่าที่ไม่ต้องการออกไปจากตาราง

1

2

3

## Results

Sort by Age in ascending order

	A	B	C	D	E
1	Name	Age	Height	Gender	Birthday
2	Bambam	25	178	M	1997-05-02
3	Cherprang	26	160	F	1996-05-02
4	Jennie	26	163	F	1996-01-16
5	Great	31	172	M	1991-01-07
6	Sergey	48	173	M	1973-06-21
7	Larry	49	181	M	1973-03-28
8	Elon	50	187	M	1971-06-28
9	Jisoo		162	F	1995-01-03

Filter by Age > 30  
(Filter out rows with Age <= 30)

	A	B	C	D	E
1	Name	Age	Height	Gender	Birthday
5	Great	31	172	M	1991-01-07
6	Sergey	48	173	M	1973-06-21
7	Larry	49	181	M	1973-03-28
8	Elon	50	187	M	1971-06-28

# Sorting and Filtering

Note: ระวังเรื่องการเรียงลำดับวันที่ จากน้อยไปมาก = เก่าไปใหม่! (ตรงข้ามกับ Age!)

Sort by Birthday in ascending order (Older → Younger)

	A	B	C	D	E
1	Name	Age	Height	Gender	Birthday
2	Elon	50	187	M	1971-06-28
3	Larry	49	181	M	1973-03-26
4	Sergey	48	173	M	1973-06-21
5	Great	31	172	M	1991-01-07
6	Jisoo	27	162	F	1995-01-03
7	Jennie	26	163	F	1996-01-16
8	Cherprang	26	160	F	1996-05-02
9	Bambam	25	178	M	1997-05-02

Z-A

A-Z

Sort by Birthday in descending order (Younger → Older)

	A	B	C	D	E
1	Name	Age	Height	Gender	Birthday
2	Bambam	25	178	M	1997-05-02
3	Cherprang	26	160	F	1996-05-02
4	Jennie	26	163	F	1996-01-16
5	Jisoo	27	162	F	1995-01-03
6	Great	31	172	M	1991-01-07
7	Sergey	48	173	M	1973-06-21
8	Larry	49	181	M	1973-03-26
9	Elon	50	187	M	1971-06-28

A-Z

Z-A

# BigQuery and more SQL knowledge

## Query

การขอข้อมูล (request for data)  
จากฐานข้อมูล (database)

Syntax:

**SELECT** รายชื่อคอลัมน์ที่ต้องการ  
**FROM** ชื่อตาราง  
**WHERE** เงื่อนไข (Optional)  
**LIMIT** จำนวนแถวที่ต้องการ (Optional)

Ex. ดึงคอลัมน์ชื่อ make, engine\_type, price ของตาราง car\_info โดยเอาเฉพาะยี่ห้อ audi และราคาเกิน \$17,000

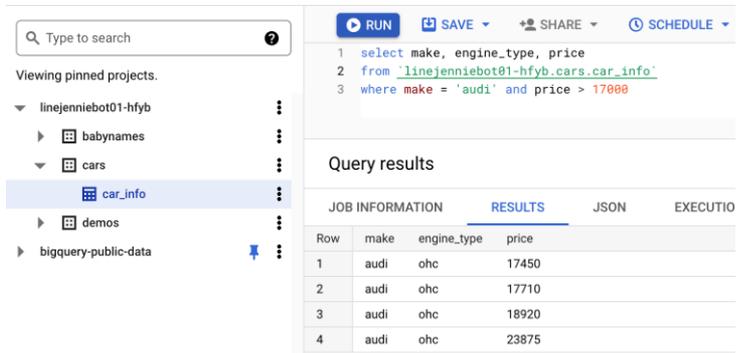
The screenshot shows the BigQuery web interface. On the left, a sidebar displays a project tree with 'linejenniebot01-hfyb' expanded to show 'cars' and 'car\_info' selected. The main area contains a search bar, a 'RUN' button, and a 'SAVE' dropdown. Below the search bar, the text 'Viewing pinned projects.' is visible. The query editor shows the following SQL code:

```
1 SELECT make, engine_type, price
2 FROM linejenniebot01-hfyb.cars.car_info
3 WHERE make = 'audi' AND price > 17000
```

Below the query editor, the 'Query results' section is active, showing a table with 4 rows and 4 columns: 'Row', 'make', 'engine\_type', and 'price'. The results are as follows:

Row	make	engine_type	price
1	audi	ohc	17450
2	audi	ohc	17710
3	audi	ohc	18920
4	audi	ohc	23875

# BigQuery and more SQL knowledge



The screenshot shows the BigQuery console. On the left, a sidebar lists projects, with 'linejennibot01-hfyb' selected and 'car\_info' highlighted. The main area displays a SQL query:

```
1 select make, engine_type, price
2 from `linejennibot01-hfyb.cars.car_info`
3 where make = 'audi' and price > 17000
```

Below the query, the 'Query results' section shows a table with 4 rows and 4 columns: 'Row', 'make', 'engine\_type', and 'price'.

Row	make	engine_type	price
1	audi	ohc	17450
2	audi	ohc	17710
3	audi	ohc	18920
4	audi	ohc	23875

backtick



## Notes:

1. เรามักตั้งชื่อคอลัมน์ด้วยตัวพิมพ์เล็กล้วน (lowercase) และเว้นวรรคด้วยเครื่องหมาย underscore `_` เรียกว่า snake\_case (ยาวเป็นงู!) เช่น `engine_type`
2. พวกคำสั่งเช่น `SELECT`, `FROM`, `WHERE`, `LIMIT`, `AND` ไม่จำเป็นต้องเป็นตัวพิมพ์ใหญ่ล้วนก็ได้ (แต่นิยมใช้ตัวพิมพ์ใหญ่ล้วนเพราะทำให้อ่านง่ายขึ้น)
3. ทืออยู่ของตาราง จะคร่อมด้วยเครื่องหมาย backtick ``projname.car.car_info`` หรือไม่คร่อม `projname.car.car_info` ก็ได้สำหรับ BigQuery
4. การตั้งชื่อตาราง พบเจอได้ทั้งแบบ snake\_case (`car_info`) หรือ CamelCase / camelCase (ขึ้นลงเหมือนหลังอูฐ เช่น `CarInfo` / `carInfo`) แต่ส่วนมากนิยม snake\_case

# BigQuery and more SQL knowledge

Syntax:

**SELECT** เอาค่าในคอลัมน์มาคำนวณ  
**FROM** ชื่อตาราง  
**WHERE** เงื่อนไข (Optional)  
**LIMIT** จำนวนแถวที่ต้องการ (Optional)

Note: เราตั้งชื่อคอลัมน์ของผลลัพธ์ใหม่ได้โดยใช้คำสั่ง **AS**

Ex. คำนวณหาผลรวม ค่าเฉลี่ย ค่าต่ำสุด ค่าสูงสุด และจำนวนแบบของเครื่องยนต์ที่ไม่ซ้ำกันของรถยี่ห้อ jaguar

The screenshot shows the BigQuery interface with a query editor containing the following SQL code:

```
1 SELECT make, engine_type, price
2 FROM linejenniebot01-hfyb.cars.car_info
3 WHERE make = 'jaguar'
```

The query results are displayed in a table:

Row	make	engine_type	price
1	jaguar	dohc	32250
2	jaguar	dohc	35550
3	jaguar	ohcv	36000

The screenshot shows the BigQuery interface with a query editor containing the following SQL code:

```
1 SELECT
2   SUM(price) AS sum_price,
3   AVG(price) AS avg_price,
4   MIN(price) AS min_price,
5   MAX(price) AS max_price,
6   COUNT(*) AS n_rows,
7   COUNT(DISTINCT engine_type) AS n_engine_type
8 FROM linejenniebot01-hfyb.cars.car_info
9 WHERE make = 'jaguar'
```

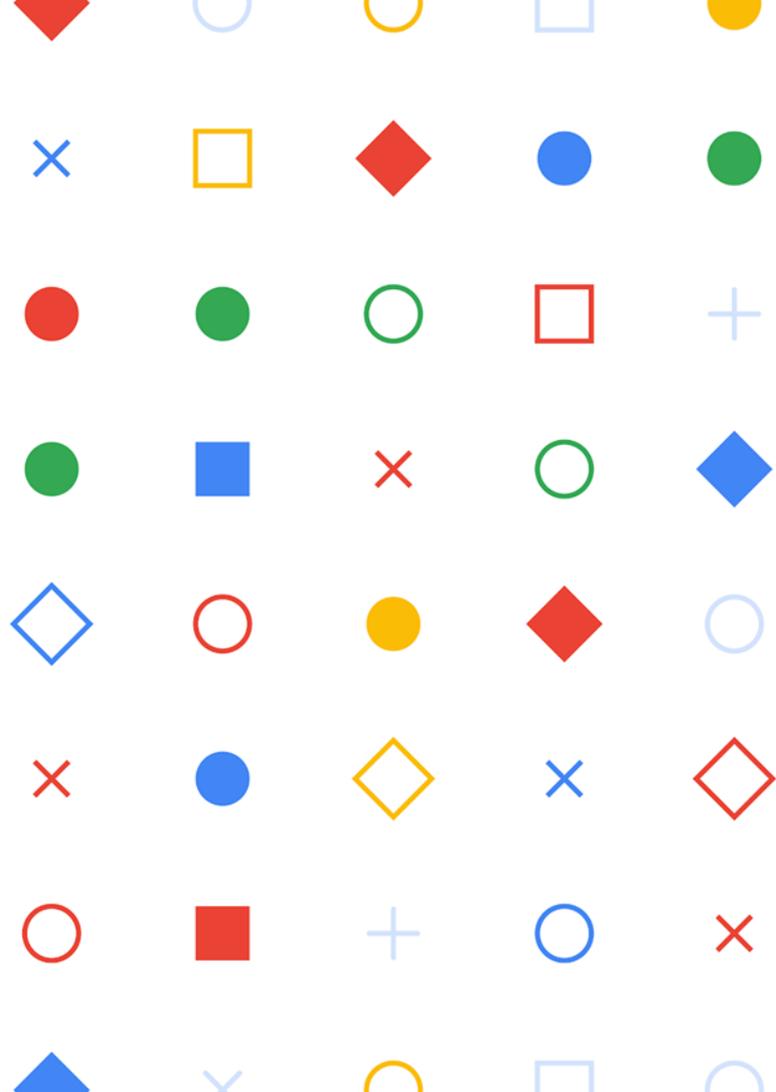
The query results are displayed in a table:

Row	sum_price	avg_price	min_price	max_price	n_rows	n_engine_type
1	103800	34600.0	32250	36000	3	2

# Organizing and Protecting Your Data

— *Part4:*

1. Organizing Files
2. Data Security Measures



# Organizing Files

เราควรจัดระเบียบไฟล์ต่าง ๆ เพื่อให้ทำงานง่าย

## Definition

### File Naming Convention

แนวทาง (Guidelines) ในการตั้งชื่อไฟล์ให้เรียงลำดับได้ และค้นหาได้ง่าย เช่น ใช้วันที่ในแบบ YYYYMMDD แทนที่จะเป็น DDMMYYYY

### Folders and subfolders

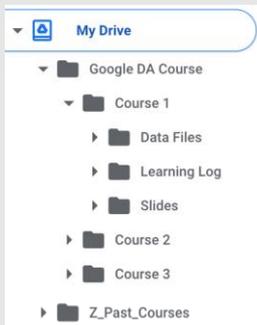
จัดโฟลเดอร์ให้งานที่เกี่ยวข้องกันอยู่ด้วยกัน และจัดลำดับความสำคัญ (Hierarchy) จากภาพกว้าง (Broad) ไปภาพย่อย (Specific)

### Archiving older files

มีพื้นที่ในการเก็บไฟล์งานเก่าที่ทำเสร็จไปแล้ว เพื่อลดความรกรุงรัง (Clutter)

## Examples

- sales\_2021\_12.csv, sales\_2022\_01.csv, sales\_2022\_02.csv
- sales\_report\_20220301\_v01.xlsx
- ExecutiveReport20220301.pdf



# Data Security Measures for Spreadsheets

## มาตรการความปลอดภัยในการรักษาความลับของข้อมูลในไฟล์ Spreadsheets

- ใน Google Sheets เราสามารถกำหนด Sharing Permissions ได้ว่าจะให้ใครสามารถ ดู (View) คอมเมนต์ (Comment) และ/หรือ แก้ไข (Edit) ได้
  - ส่วน Excel เราสามารถเข้ารหัสไฟล์ (Encrypt with Password) ได้
- **WARNING:** เราซ่อน (Hide) Tab ได้ แต่คนอื่นที่เป็น Editor ก็ยังสามารถนำกลับมา (Unhide) ได้!  
(ส่วน Commenter/Viewer ยังเห็นชื่อ Tab นั้นได้อยู่ แม้จะ Unhide กลับมาไม่ได้)

