## Class period 6

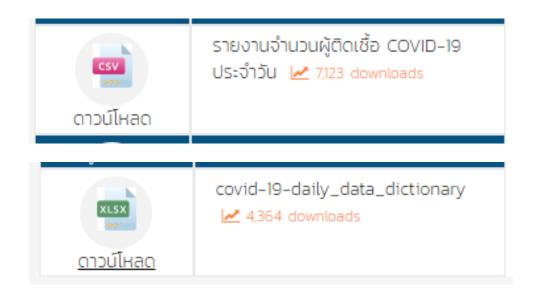
Pandas 101

#### **Pandas**

- Pandas เป็นหนึ่งใน package ที่สำคัญของ python ใช้สำหรับจัดการข้อมูลรูปแบบตาราง .CSV
- import pandas as pd

# Download ข้อมูลรายงาน COVID-19 ประจำวัน ข้อมูล ประจำประเทศไทย

https://data.go.th/dataset/covid-19-daily

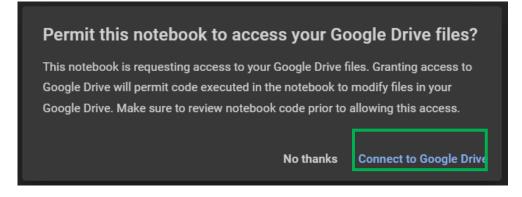


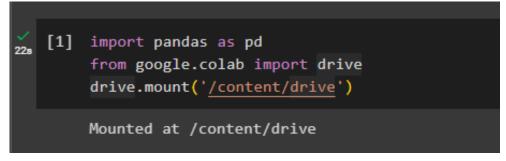
CSV = Comma Separated Values ในการจะใส่ค่าแต่ละค่า จะใช้ comma ในการแยก

Meta data = Data that description data ใช้อธิบายช้อมูล

## การนำข้อมูลเข้า

- 1. สร้าง folder ใน google drive และนำข้อมูล .csv ที่ดาวน์โหลดเข้าไปเก็บไว้ใน folder ที่สร้าง
- 2. นำเข้า package pandas และ package ของ google.colab ที่ชื่อ drive เพื่อเชื่อมต่อ google drive กับ google colab
- import pandas as pd
- from google.colab import drive drive.mount('/content/drive')
- 3. กด Connect Google Drive และเลือก Account
- 4. กด select all และกด continue





#### Import os

- นำเข้า package os เพื่อใช้ในการทำงานต่างๆที่เกี่ยวกับไฟล์ เช่น การชื้ไฟล์ การลบไฟล์ การสร้าง โฟลเดอร์ เป็นต้น
- โดยในกรณี google drive จะใช้ os เพื่อชี้ไฟล์ ว่าไฟล์ที่ต้องการใช้งานอยู่ path ใหนใน google drive ที่เชื่อม
- path คือเส้นทางที่อยู่ไฟล์ จะทำงานเหมือนกับ path ใน window เช่น
- E:\WORKSPACE\Basic Programming\confirmed-cases.csv
- หมายความว่า ไฟล์ confirmed-cases.csv อยู่ใน drive E โฟลเดอร์ WORKSPACE ในโฟลเดอร์ Basic Programming

#### การ set path

- path = '/content/drive/MyDrive/dataviz\_2024\_data'
- Set 'path' ที่ชี้ไปยังโฟลเดอร์ที่เก็บไฟล์ .csv ไว้ใน google drive และเก็บ string ไว้ใน ตัวแปร path
- โดย path หรือเส้นทางที่ชี้ไปยังโฟลเดอร์และไฟล์ต่างๆ ใน os ของ window, mac หรือ linux จะใช้สัญลักษณ์ใน path แตกต่างกัน
- Package os จะช่วยให้สามารถเชื่อม path โดยไม่ต้องคำนึงถึงสัญลักษณ์ เพราะ os จะใส่ สัญลักษณ์เชื่อมให้เองตาม platform ที่ใช้งานอยู่ เช่น
- ถ้าใช้ os ของ window ก็จะเชื่อมด้วย \
- ถ้าใช้ os ของ mac หรือ linux จะเชื่อมด้วย /

## คำสั่ง os.path.join

- เป็นคำสั่งที่ใช้สำหรับเชื่อม path เข้าด้วยกัน
- import os
- covid\_file\_path = os.path.join(path, confirmed-cases.csv')
- หมายความว่า เชื่อมตัวแปร path ที่ set ไว้ก่อนหน้านี้เข้ากับชื่อไฟล์ confirmed-cases.csv เก็บไว้ในตัวแปร covid\_file\_path
- print(covid\_file\_path)
- ผลลัพธ์จะได้เส้นทางไปยังไฟล์ที่ต้องการอยู่ในตัวแปร covid\_file\_path
- /content/drive/My Drive/dataviz\_2024\_data/confirmed-cases.csv

### load data to memory (pd.read\_csv)

- pd.read\_csv เป็นคำสั่งที่ใช้สำหรับโหลดข้อมูล
- data\_covid = pd.read\_csv(covid\_file\_path)
- โหลดข้อมูลไฟล์ confirmed-cases.csv ตามเส้นทาง covid\_file\_path
- data covid พิมพ์ชื่อตัวแปรที่เก็บข้อมูล
- ผลลัพธ์จะได้หน้าไฟล์ CSV

	No.	announce_date	notification_date	sex	age	nationality	province_of_isolation	province_of_onset	district_of_onset	
0	1	1/12/2020	NaN	หญิง	61.0	China	กรุงเทพมหานคร	กรุงเทพมหานคร	NaN	คนต่า เดิน <i>ข</i> จา ปร
1	2	1/17/2020	NaN	หญิง	74.0	China	กรุงเทพมหานคร	กรุงเทพมหานคร	NaN	คนต่า เดิน <i>v</i> จา ปร
2	3	1/22/2020	NaN	หญิง	73.0	Thailand	นครปฐม	นครปฐม	เมือง	คนต่า เดินท จา ปร
3	4	1/22/2020	NaN	ชาย	68.0	China	กรุงเทพมหานคร	กรุงเทพมหานคร	NaN	คนต่า เดินท จา ปร
4	5	1/24/2020	NaN	หญิง	66.0	China	นนทบุรี	กรุงเทพมหานคร	NaN	คนต่า เดินท จา ปร
12648	12649	1/20/2021	1/19/2021	หญิง	44.0	Thailand	ชลบุรี	ชลบุรี	บางละมุง	Quara

## คำสั่ง .head()

• data\_covid.shape ชื่อตัวแปรที่เก็บข้อมูลตามด้วย .head() ใช้เพื่อให้แสดงชื่อ คอลัมน์และข้อมูลในตารางเฉพาะ 5 แถวแรก int, default=5

- สามารถกำหนดจำนวนคอลัมน์ที่ต้องการให้แสดงได้ เช่น
- data covid.head(10) จะแสดงชื่อคอลัมน์และข้อมูลในตาราง 10 แถว

## คำสั่ง .shape

- data\_covid.shape ชื่อตัวแปรที่เก็บข้อมูลตามด้วย .shape ใช้ตรวจสอบขนาดของ ข้อมูล ผลลัพธ์จะได้
- (839771, 11)
- หมายความว่า มีข้อมูลทั้งหมด 839, 771 แถว มีคอลัมน์ 11 คอลัมน์

## การชี้ค่าในข้อมูลตารางแบบ basic

- ใช้ชื่อคอลัมน์ในการดึงข้อมูลในคอลัมน์ที่ต้องการ
- data\_covid['province\_of\_onset']

```
data_covid['province_of_onset']
          กรุงเทพมหานคร
          กรุงเทพมหานคร
                  นครปฐม
          กรุงเทพมหานคร
          กรุงเทพมหานคร
839766
              กาญจนบรี
839767
              กาญจนบุรี
839768
              กาญจนบุรี
839769
              กาญจนบุรี
839770
              กาญจนบุรี
Name: province_of_onset, Length: 839771, dtype: object
```

## การชี้ค่าในข้อมูลตารางแบบ basic

- การใช้ชื่อคอลัมน์และลำดับแถวในการดึงข้อมูลในแถวและคอลัมน์ที่ต้องการ
- data\_covid['province\_of\_onset'][4]
- ผลลัพธ์จะได้ข้อมูลแถวที่ 4 นับจาก 0 ในคอลัมน์ province\_of\_onset
- 'กรุงเทพมหานคร'

## การชี้ค่าในข้อมูลตารางแบบ .iloc

- โดยการมองมุมมองข้อมูลตารางในรูปแบบ numpy array หรือ matrix จะใช้ตำแหน่งเพื่อชี้ข้อมูล ที่ต้องการ เช่น
- data\_covid.iloc[4,9]
- ผลลัพธ์จะได้ข้อมูลแถวที่ 4 คอลัมน์ที่ 9 (ในมุมมอง matrix คือหลักที่ 9) นับจาก 0 คือคอลัมน์ province of onset
- 'กรุงเทพมหานคร'

## Table slicing การเลือกข้อมูลเฉพาะคอลัมน์ที่ต้องการ

- การเลือกข้อมูลเฉพาะคอลัมน์ที่ต้องการมาเก็บไว้ในตัวแปรเพื่อนำไปใช้งาน
- smaller\_table = data\_covid[['announce\_date','province\_of\_onset','risk']]
- หมายความว่า เลือกข้อมูลคอลัมน์ชื่อ announce\_date, province\_of\_onset, risk ในข้อมูลที่เก็บอยู่ในตัวแปร data\_covid และเก็บข้อมูลเฉพาะคอลัมน์ที่เลือกไว้ในตัวแปร

smaller table

• ผลลัพธ์จะได้

-	smaller_ smaller_	_	ovid[['announce_date	e','province_of_onset','risk']	1
		announce_date	province_of_onset	risk	E
	0	12/1/2020	กรุงเทพมหานคร	คนต่างชาติเดินทางมาจากต่างประเทศ	11.
	1	17/1/2020	กรุงเทพมหานคร	คนต่างชาติเดินทางมาจากต่างประเทศ	
	2	22/1/2020	นครปฐม	คนต่างชาติเดินทางมาจากต่างประเทศ	
	3	22/1/2020	กรุงเทพมหานคร	คนต่างชาติเดินทางมาจากต่างประเทศ	
	4	24/1/2020	กรุงเทพมหานคร	คนต่างชาติเดินทางมาจากต่างประเทศ	
	839766	12/8/2021	กาญจนบุรี	ทัณฑสถาน/เรือนจำ	
	839767	12/8/2021	กาญจนบุรี	ทัณฑสถาน/เรือนจำ	
	839768	12/8/2021	กาญจนบุรี	ทัณฑสถาน/เรือนจำ	
	839769	12/8/2021	กาญจนบุรี	ทัณฑสถาน/เรือนจำ	
	839770	12/8/2021	กาญจนบุรี	ทัณฑสถาน/เรือนจำ	
	839771 rc	ows × 3 columns			

## Table slicing การเลือกข้อมูลเฉพาะแถวที่ต้องการแบบง่าย

- การเลือกข้อมูลเฉพาะแถวที่ต้องการแบบง่าย คือการมองมุมมองข้อมูลตารางในรูปแบบ array แต่การ นำไปใช้งาน ใช้งานอะไรไม่ค่อยได้
- data covid.iloc[1:5,:]
- หมายความว่า
- 1:5 คือเลือกข้อมูลที่อยู่ในแถวที่ 1 ไปจนถึงแถวที่ 4
- , : คือเลือกทุกคอลัมน์ ดังนั้น
- data\_covid.iloc[1:5,:] คือเลือกข้อมูลในตัวแปร data\_covid ที่อยู่ ในแถวที่ 1 ไปจนถึงแถวที่ 4 และเลือกทุกคอลัมน์

# Table slicing การเลือกข้อมูลเฉพาะแถวที่ต้องการแบบ advanced

- การเลือกข้อมูลเฉพาะแถวที่ต้องการแบบ advanced คือการใช้ logic query ในการเลือกข้อมูล
- data\_covid[data\_covid['province\_of\_onset'] == 'ขอนแก่น']
- หมายความว่า เลือกข้อมูลที่อยู่ในตัวแปร data\_covid โดยกำหนดชื่อคอลัมน์ที่ต้องการคือ province\_of\_onset และต้องการข้อมูลทุกแถวที่มีข้อมูลในคอลัมน์ province of onset เป็นจังหวัดขอนแก่น

## วิธีการเลือกข้อมูลเฉพาะแถวที่ต้องการแบบ advanced

- การทำงานของการเลือกข้อมูลเฉพาะแถวที่ต้องการแบบ advanced
- คือ การใส่แถวที่ต้องการ และใส่ list True/False ที่มีขนาดเท่ากับจำนวนแถว เพื่อเปรียบเทียบ ข้อมูลในแถวนั้นๆ ด้วยเงื่อนไข logical expression (True/False)) เช่น
- สร้างตารางใช้สำหรับยกตัวอย่าง
- eight\_rows\_covid = data\_covid.iloc[:8,:]
   eight\_rows\_covid
- หมายความว่า เลือกข้อมูลในตัวแปร data\_covid แถวที่ 0 ถึงแถวที่ 7 ทุกคอลัมน์เก็บไว้ในตัว แปร eight rows covid

## การทำงานของการเลือกข้อมูลเฉพาะแถวที่ต้องการแบบ advanced

- ใส่แถวที่ต้องการด้วยการกำหนดค่า True (แถวที่ต้องการ) / False (แถวที่ไม่ต้องการ)
- eight\_rows\_covid[[True, True, False, True, True, True, True, False]]
- ผลลัพธ์จะได้ข้อมูลตารางตามค่า True/False ที่เลือกใน list คือแถวที่ 0, 1, 3, 4, 5, 6
- เช่นเดียวกันกับการสร้าง list ของ logical expression แต่แทนที่จะเลือกเอง โดยการใส่ list True/False ให้กำหนดเงื่อนไขและข้อมูลที่ต้องการ เพื่อเปรียบเทียบและ เลือกข้อมูลที่ตรงตามเงื่อนไข โดยถ้าตรงตามเงื่อนไขคือ True ไม่ตรงตามเงื่อนไขคือ False

### การสร้าง list ของ logical expression

```
• eight rows covid['province of onset'] == 'กรุงเทพมหานคร'
• ผลลัพธ์จะได้
        True
• 1
        True
• 2
       False
• 3
        True
• 4
       True
• 5
        True
• 6
        True
• 7
       False
• Name: province of onset, dtype: bool
```

## นำ list ของ logical expression ที่สร้างมาใช้งาน

- ซึ่งเมื่อนำมาใช้งานเลือกข้อมูลเฉพาะแถวที่ต้องการแบบ advanced คือการใช้ logic query ในการ เลือกข้อมูล นั้นเอง
- eight\_rows\_covid[eight\_rows\_covid['province\_of\_onset'] == 'กรุงเทพมหานคร']
- ผลลัพธ์จะได้ข้อมูลทุกแถวที่มีข้อมูลในคอลัมน์ **province\_of\_onset** เป็น กรุงเทพมหานคร คือแถว ที่ 0, 1, 3, 4, 5, 6

### Homework class period 6

- (ให้ใช้เฉพาะที่อาจารย์สอนไปแล้วในวิชานี้)
- คำนวณ อายุเฉลี่ย ของผู้หญิง และผู้ชาย ของข้อมูลทั้งหมด
- คำนวณ อายุเฉลี่ย ของผู้หญิง และผู้ชาย ของผู้ป่วยในจังหวัดขอนแก่น
- หาจำนวนผู้ป่วยที่เป็นคน "คนต่างชาติเดินทางมาจากต่างประเทศ"