**Slide 1**

Class period 7

บทที่ 5 การเตรียมข้อมูลสำหรับการแสดงผล 1
Pandas 101 part2

1

**Slide 2**

เฉลย Homework class period 6

- คำนวณ อายุเฉลี่ย ของผู้หญิง และผู้ชาย ของข้อมูลทั้งหมด
- this_data = data_covid[['sex','age','province_of_onset']]
- this_data
- เลือกข้อมูลเฉพาะคอลัมน์ที่ต้องการใช้งานและเก็บไว้ในตัวแปร this data
- female = this_data[this_data['sex']=='หญิง]
- เลือกเฉพาะข้อมูลที่ข้อมูลในคอลัมน์ sex เท่ากับ หญิง เก็บไว้ในตัวแปร female
- female['age']
- เลือกแสดงข้อมูลในตัวแปร female เฉพาะคอลัมน์ age ก็จะได้ข้อมูลอายุของเพศหญิงทั้งหมด

1

**Slide 3**

เฉลย Homework class period 6

- จากนั้นวนลูปเพื่อหาอายุเฉลี่ย
- sum = 0
- N = 0
- for a in female['age']:
-     if a > 0:
-         sum += a # sum = sum + a
-         N += 1
- print(f'อายุเฉลี่ยของผู้หญิง{sum/N}')
- กำหนดตัวแปร sum=0 และ N=0 เพื่อใช้ในการเก็บค่าจากการบวกในการวนลูปและจะรอบจนถึงรอบสุดท้าย โดย sum จะใช้เก็บค่าอายุ และ N ใช้เก็บค่าจำนวนผู้หญิง

2

**Slide 4**

เฉลย Homework class period 6

- sum = 0
- N = 0
- กำหนดตัวแปร sum=0 และ N=0 เพื่อใช้ในการเก็บค่าจากการบวกในการวนลูปและจะรอบจนถึงรอบสุดท้าย โดย sum จะใช้เก็บค่าอายุ และ N ใช้เก็บค่าจำนวนผู้หญิง
- for a in female['age']:
- วนลูปอ่านค่าอายุของผู้หญิงที่ละคนเก็บไว้ในตัวแปร a
- if a > 0:
- ตั้งเงื่อนไขในตาอายุมากกว่า 0 ถึงจะนำค่าอายุมาบวกหาค่าเฉลี่ย เพื่อหลีกเลี่ยงค่า missing(บางคนไม่มีข้อมูลอายุ)

3

**Slide 5**

เฉลย Homework class period 6

- sum += a # sum = sum + a
- N += 1
- นำตัวแปร sum มานวกค่าอายุของผู้หญิงทีละคน จนสุป 1 รอบก็จะเอาผลลัพธ์จากการบวกรอบที่แล้วมานวกต่อไปเรื่อยๆ เพื่อหาค่าอายุรวม
- นำตัวแปร N มานวก 1 เพื่อนับนับจำนวนผู้หญิง
- print(f'อายุเฉลี่ยของผู้หญิง{sum/N}')
- นำตัวแปร sum และ N มาหารกันเพื่อหาค่าเฉลี่ย ผลลัพธ์จะได้

4

**Slide 6**

การจัดการ Missing Value

- มีทั้งหมด 3 แบบ
- 1. ลบ record ที่เป็น missing
- 2. แทนที่ ค่า missing ด้วยค่าที่เหมาะสม mean, default, category-unknown
- 3. ใช้ ค่าจาก columns อื่นๆ ช่วยประมาณค่า ค่าใน column ที่หายไป (regression, deep learning, etc.)

5

**Slide 7**

ลบ record (dropna)

- missing = None, NA(not autorized), NaN (not a number)
- .dropna() เป็นคำสั่งที่ใช้ในการลบข้อมูลแถวที่ไม่มีค่าหรือไม่มีข้อมูล ตัวอย่างเช่น
- data_covid.shape ผลลัพธ์จะได้ขนาดของข้อมูล data_covid
- (839771, 11)
- data_covid.dropna().shape ผลลัพธ์จะได้ขนาดของข้อมูล data_covid ที่ลบแถวข้อมูลที่มีค่าเป็น None
- (599988, 11)

6

**Slide 8**

การใช้งาน .dropna()

- สามารถเลือกลบข้อมูลที่เป็น None เฉพาะในคอลัมน์ที่ต้องการใช้งาน แทนที่จะเลือกลบจากข้อมูลทั้งหมด เช่น
- this_data = data_covid[['sex','age','province_of_onset']]
- this_data.shape ผลลัพธ์จะได้
- (839771, 3)
- this_data.dropna().shape ผลลัพธ์จะได้
- (674906, 3)
- จะเห็นว่าเมื่อเทียบกับ data_covid.dropna().shape ที่เป็นข้อมูลทั้งหมด (599988, 11)
- ข้อมูลที่เลือกเฉพาะคอลัมน์ที่ต้องการใช้งานจะมีจำนวนข้อมูลมากกว่า

7

**Slide 9**

การใส่ตัวแปรเพื่อรับค่า

- this_data.dropna()
- print(this_data.shape) ผลลัพธ์ได้
- (839771, 3)
- ซึ่งไม่ใช่ผลลัพธ์ที่ได้จากการใช้ .dropna() เพื่อลบข้อมูลแถวที่มีค่าเป็น None เนื่องจากไม่ได้มีตัวแปรเข้ารับค่า เช่น
- This_data_dn = this_data.dropna()
- print(This_data_dn.shape ผลลัพธ์ได้
- (674906, 3)

8

## Parameter: inplace ของ .dropna()

- inplace จะเป็นการอัพเดทค่าในตารางเลย โดยที่ไม่จำเป็นต้องมีตัวแปรมารับค่า เช่น

- this_data.dropna(inplace=True)
- print(this_data.shape)
- (674906, 3)

*(slide page 9)*

**10**

## Parameter: subset ของ .dropna()

- subset จะเป็นการเลือกเฉพาะคอลัมน์ที่ต้องการลบแถวข้อมูลที่เป็น None เฉพาะคอลัมน์ที่เลือก เช่น
- this_data = data_covid[['sex','age','province_of_onset']]
- this_data.shape ขนาดของข้อมูลตาราง
- (839771, 3)
- this_data.dropna().shape ขนาดของข้อมูลตารางที่มีการลบข้อมูลแถวที่เป็น None แบบปกติ
- (674906, 3)
- this_data.dropna(subset=['age']).shape
- ขนาดของข้อมูลตารางที่มีการกำหนด subset ลบข้อมูลแถวที่เป็น None เฉพาะ subset ที่กำหนด
- (763606, 3)

**11**

## แทน missing ด้วยค่าที่เหมาะสม .fillna()

- .fillna() เป็นคำสั่งที่ใช้ในการแทนค่าที่ missing หรือมีค่า None ด้วยค่าที่กำหนด เช่น
- this_data = data_covid[['sex','age','province_of_onset']]
- this_data_updated = this_data.fillna(value={'sex':'ไม่รู้','age':-1, 'province_of_onset':'ไม่รู้'})
- หมายความว่า ให้แทนที่ข้อมูลที่เป็น None
- ในคอลัมน์ sex แทนด้วย 'ไม่รู้'
- ในคอลัมน์ age แทนด้วย -1
- และในคอลัมน์ province_of_onset ให้แทนที่ค่า None ด้วย 'ไม่รู้'
- ผลลัพธ์จะได้ข้อมูลตาราง this_data_updated ที่ไม่มีค่า None

**12**

## แทน missing ด้วยค่าที่เหมาะสม .fillna()

- สามารถตรวจสอบว่ามีค่า None ถูกแทนที่ด้วยค่าที่กำหนดไว้แล้วจริง

- this_data_updated[this_data_updated['sex']=='ไม่รู้']
- ผลลัพธ์จะแสดงว่ามีข้อมูลในคอลัมน์ sex ที่ถูกแทนที่ด้วย 'ไม่รู้'

*(slide page 12)*

**13**

## การใช้ logical expression จากข้อมูลตารางอื่น

- data_covid[this_data_updated['province_of_onset']=='ไม่รู้']

- จะเห็นได้ว่าในสวนที่กำหนดเงื่อนไขของ logical expression มาจากตัวแปร this_data_updated ซึ่งนำมาใช้ในข้อมูลตารางของตัวแปร data_covid

- สาเหตุที่สามารถนำมาใช้ข้ามกันได้เนื่องจากข้อมูลตารางทั้งสองมีจำนวนและลำดับของข้อมูลเหมือนกันก็เท่ากับแถว

- ถ้าหาก 2 ตัวแปรมีจำนวนแถวตารางกันในแต่ละแถวและมีข้อมูลเหมือนกันก็สามารถใช้งานได้ แต่ถอดซัพท์ที่ดึงเวลงมาไม่ถูกต้อง เพราะไม่ใช่ข้อมูลเดียวกัน

**13**

## การวนลูป record ในตาราง .iterrows()

- .iterrows()เป็นคำสั่งที่ใช้ในการสามารถวนลูปนำข้อมูลในตาราง
- this_data = data_covid[['sex','age','province_of_onset']]
- for each_row in this_data.iterrows():
- if (each_row[1]['age'] == 20) and (each_row[1]['province_of_onset'] == 'พนมพ'):
- print(each_row)
- หมายความว่า
- ให้วนลูปนำค่าในข้อมูลตารางตัวแปร this_data ทีละแถวและจะเก็บไว้ในตัวแปร each_row
- สวนค่าในแถวกำหนดลูปกำหนดเงื่อนไขค่าที่กำหนดไว้เลือก print ข้อมูลเฉพาะแถวที่มีข้อมูลในตาราง age=20 และ province_of_onset=จะของแผน

**14**