

Class period 8

บทที่ 5 การวิเคราะห์ข้อมูลสำหรับการตลาด 1 (ต่อ)
Pandas 101 part3

1

Quiz สั้นๆ (5 นาที)

- ให้ดูโค้ดที่ทำงานเกี่ยวกับ loc ซึ่งรับ numpy array หรือ matrix ให้ผลลัพธ์เป็น pandas (iloc)

```

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'] == 'samutprakan':
        print(each_row)

```

2

การวนลูปอ่านข้อมูลแบบตาราง (pandas dataframe) เป็น numpy array หรือ matrix (iloc)

```

this_data = data_covid[['sex','age','province_of_onset']]

for each_row in range(this_data.shape[0]):
    if (this_data.iloc[each_row,1] == 20) and (this_data.iloc[each_row,2] == 'samutprakan'):
        print(each_row)
        print(this_data.iloc[each_row,:])

```

3

การวนลูป

```

for each_row in range(this_data.shape[0]):
    # การวนลูปค่าตัวแปรในการเขียนแต่ละครั้งให้ตัวแปร each_row โดยจำนวนการวนลูปในการเขียนจะเท่ากับจำนวนแถว
    this_data.shape
    (839771, 3)
    this_data.shape[0]
    839771
    range(this_data.shape[0]) คือ สร้าง list ด้วยขนาดจำนวนแถวทั้งหมดที่จะใช้วนลูปทั้งหมด
    [0, 1, 2, 3, ..., 839771]

```

4

การใช้ .iloc[] ในการชี้ข้อมูลในตาราง

```

if (this_data.iloc[each_row,1] == 20) and (this_data.iloc[each_row,2] == 'samutprakan'):
    # สร้างเงื่อนไขสำหรับเลือกเฉพาะข้อมูลที่ต้องการคือ age=20 และ province_of_onset=สมุทรปราการ โดยใช้ตัวแปร iloc มาช่วย
    each_row คือ ตัวเลขจากฐานข้อมูลคือตัวแปรเป็นตัวชี้ของข้อมูลที่ต้องการ คือเป็น
    this_data.iloc[each_row,1] == 20 คือ ใช้ข้อมูลแถวตาม index ของตัวชี้ที่ 1 (เป็นค่า 0) ตรวจสอบว่าเท่ากับ 20 หรือไม่ ถ้าใช่ แสดง เป็นที่ต้องการข้อมูล
    this_data.iloc[each_row,2] == 'samutprakan' คือ ใช้ข้อมูลแถวตาม index ของ province_of_onset ตัวชี้ที่ 2) ตรวจสอบว่าเท่ากับ 'สมุทรปราการ' หรือไม่ ถ้าใช่ แสดง เป็นที่ต้องการข้อมูล

```

5

print() แสดงข้อมูล

```

print(each_row)
print(this_data.iloc[each_row,:])

# แสดงเป็นตัวเลข
print(each_row) คือ ตัวเลขตัวเดียว และ
print(this_data.iloc[each_row,:]) คือ ข้อมูลในแถวนั้นๆ ทั้งหมด

```

6

Quiz ในห้อง (15 นาที)

- คัดตารางเฉพาะปี 2021 announce_date ในปี 2021
- Hint
- ระบุค่า index ของปี 2021
- คัดตารางเฉพาะ ปี 2021

7

เฉลย

```

TF=list()
for each_row in data_covid.iterrows():
    if each_row[1]['announce_date'].split('/')[2] == '2021':
        TF.append(True)
    else:
        TF.append(False)

data_covid[TF].head()

```

8

เตรียม list() วาง วนลูปและสร้างเงื่อนไข

```

TF=list()
สำหรับ วนลูปให้ตัวแปร TF (เลือกเก็บเป็นผลลัพธ์ True/False) ซึ่งใช้สำหรับการวนลูป

for each_row in data_covid.iterrows():
    # ระบุค่า index ของตารางเฉพาะ data_covid ที่ต้องการจะดึงค่าตัวแปร each_row
    if each_row[1]['announce_date'].split('/')[2] == '2021':
        # ถ้าเป็นปีที่เราต้องการข้อมูลก็เก็บค่าในคอลัมน์ announce_date เท่ากับ 2021 โดยการใช้ split('/')/2) เพื่อแยกค่าตามวันที่เพื่อ
        # แยกปีมาใช้ในการเก็บเป็นผลลัพธ์
        TF.append(True)
    # ถ้าไม่เป็นปีที่เราต้องการก็เก็บค่าในคอลัมน์ announce_date ไม่เท่ากับ 3/1/2021 คือ ไม่เก็บค่า
    # ดังนั้นผลลัพธ์ที่เก็บมา / ปีที่เราต้องการที่ 2 ปีจาก 0

```

9

เขียน True False เข้าไปใน List

```

if each_row[1]["announce_data"].split('/')[2] == '2021':
    TF.append(True)
else:
    TF.append(False)

```

- ถ้าตามเงื่อนไขให้เขียน True เข้าไปใน list ที่เตรียมไว้ ถ้าไม่ตามให้เขียน False
- หมายเหตุว่า ถ้าข้อมูลเป็นคอมเมนต์ ลากออก, ใส่วันที่ 2021 ให้เขียน True เข้าไปใน list ถ้าไม่ ให้เขียน False
- ผลลัพธ์จะได้ list True False ตามจำนวนคอมเมนต์ข้อมูลข่าว data_covid ที่เลือก True เท่ากับ 2021

10

ผลลัพธ์

- data_covid[TF].head()
- นำข้อมูลเฉพาะ list True False ขึ้นมาทำการดูข้อมูลเฉพาะปี 2021
- ผลลัพธ์จะได้รายชื่อข้อมูลที่เป็นค่าจริงในคอมเมนต์ announce_ ใส่วันที่ 2021

11

Function ตัวช่วยใน pandas

- describe() จำนวนค่าทางสถิติของข้อมูลที่เป็นตัวเลข
- mean() ค่ารวมค่าเฉลี่ยของข้อมูลโดยไม่นับค่า missing
- isnull() ตรวจสอบค่าข้อมูลที่ missing (none)

12

.describe()

- ใช้สำหรับคำนวณค่าทางสถิติของข้อมูลที่เป็นตัวเลข
- นำข้อมูลที่เป็น missing หรือ none มาคำนวณ
- (ค่าของข้อมูลที่ไม่มีค่าเป็น none ไม่ได้อิงผล)

13

.mean()

- ใช้สำหรับคำนวณค่าเฉลี่ยของข้อมูลที่เป็นตัวเลข (ค่าของข้อมูลที่ไม่มีค่าเป็น none ไม่ได้อิงผล)

```

data_covid[data_covid["sex"]=="female"].mean()

```

- ผลลัพธ์จะได้
- 34.962927022130938

14

.isnull()

- ใช้ตรวจสอบค่า missing ในข้อมูลตาราง
- True คือ missing (ค่าว่าง)
- False คือไม่ขาดหาย

15