

1 Pandas คืออะไร?

2 เทคนิคการใช้ Pandas

- 2.1 วิธีเช็ค Version Pandas
- 2.2 วิธีการโหลดไฟล์ CSV (Import)
- 2.3 วิธีเช็คจำนวนแถว และจำนวนคอลัมน์
- 2.4 วิธีสุ่มข้อมูลสำหรับเช็ค (Sample)
- 2.5 วิธีเช็คข้อมูลหาความผิดปกติใน DataFrame เบื้องต้น
- 2.6 วิธีแปลงประเภทข้อมูล (Data Type) ใน Data Frame
- 2.7 วิธีเช็ค Summary ของแต่ละคอลัมน์ (count, min, max, mean)
- 2.8 วิธีเช็ค Summary (count, min, max, mean) แบบแยกกลุ่ม
- 2.9 วิธีสร้าง DataFrame ใหม่
- 2.10 วิธีเลือกหลายคอลัมน์จาก DataFrame
- 2.11 วิธีเลือกคอลัมน์ตามเงื่อนไขที่ต้องการ
- 2.12 วิธีเลือกแถวตามเงื่อนไขที่ต้องการ
- 2.13 วิธีเพิ่มคอลัมน์ใหม่
- 2.14 การสลับ Row <-> Column (Transpose)
- 2.15 การต่อ DataFrame
- 2.16 การต่อ DataFrame แบบ Join
- 2.17 การหาค่า Mean, Sum, Max (Aggregate) แบบทั้ง DataFrame
- 2.18 การ Aggregate แบบตามกลุ่มที่ต้องการ
- 2.19 การรัน Function เดียวกันทุกแถว หรือทุกคอลัมน์
- 2.20 รันคำสั่งที่เขียนเองกับทุกแถวใน 1 คอลัมน์
- 2.21 รันคำสั่งที่เขียนเองกับทุกค่า
- 2.22 คำนวณ Correlation & Covariance
- 2.23 คำนวณ Cross Tabulation
- 2.24 วิธีหาค่า Unique ในแต่ละคอลัมน์
- 2.25 วิธีเช็คว่ามีแถวไหนข้อมูลซ้ำมัย (Duplicated)
- 2.26 วิธีการนับจำนวน Duplicate
- 2.27 วิธีการลบ Duplicate
- 2.28 วิธีการลบแถว และลบคอลัมน์
- 2.29 วิธีการลบแถวที่มี Missing Value
- 2.30 วิธีแทนค่า Missing Value ด้วยค่าเฉลี่ย (Mean Imputation)
- 2.31 การลบข้อมูลแต่ละคอลัมน์ และแต่ละแถว
- 2.32 วิธีเปลี่ยน DataFrame จากแบบ Wide เป็น Long (Melt)
- 2.33 วิธีการเปลี่ยนชื่อคอลัมน์ (Rename)
- 2.34 วิธีการใส่คำนำหน้าคอลัมน์ (Prefix)
- 2.35 วิธีการแทนค่าใน DataFrame
- 2.36 วิธีการ Export DataFrame เป็นไฟล์ CSV

3 สรุปการใช้งาน Pandas

Pandas คืออะไร?

Pandas เป็น Library ใน Python ที่ทำให้เราเล่นกับข้อมูลได้ง่ายขึ้น เหมาะมากสำหรับทำ [Data Cleaning / Wrangling](#) ครับผม

วิธีการใช้งาน Pandas คือ โหลดไฟล์ข้อมูล เช่น CSV เข้าไป แล้วเราจะได้ข้อมูลในรูปแบบตาราง (DataFrame) ที่แบ่งข้อมูลตามแถวและคอลัมน์ หรือเหมือน Excel ที่เราใช้กันนั่นเอง

	date	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	sqft_above	sqft_basement	yr_built	yr_renovated
0	20140502T000000	313k	3	1.50	1340	7912	1.5	0	0	3	1340	0	1965	NaN
1	20140502T000000	2.384m	5	2.50	3650	8090	2.0	0	4	5	3370	280	1921	NaN
2	20140502T000000	342k	3	2.00	1930	11947	1.0	0	0	4	1930	0	1966	NaN
3	20140502T000000	420k	3	2.25	2000	8030	1.0	0	0	4	1000	1000	1963	NaN
4	20140502T000000	550k	4	2.50	1940	10500	1.0	0	0	4	1140	800	1976	NaN

ตัวอย่าง DataFrame ของ Pandas เป็นตารางเหมือน Excel เลยครับ

ป.ล. Pandas ไม่เกี่ยวกับหมีแพนด้านะฮะ จริง ๆ แล้วมาจากคำว่า PANel DATa ซึ่งหมายถึงข้อมูลที่มีหลายมิตินั่นเอง

เทคนิคการใช้ Pandas

อย่างที่**แอตมินเคยเล่า** ว่าการทำ Data Wrangling เป็นงานที่ค่อนข้างถึกครับ วันนี้เลยรวบรวมโค้ดต่าง ๆ ในการใช้ Pandas มาให้ ซึ่งน่าจะครอบคลุมการใช้งานประมาณ 80 – 90% แล้วครับผม

โค้ดบางส่วนมาจากคลาส Data Wrangling ที่แอตมินเรียน และจากเว็บไซต์ [MyCheatSheet](#) ครับ

วิธีเช็ค Version Pandas

โค้ดนี้เหมือนไม่สำคัญ แต่จริง ๆ แล้วสำคัญมากเวลาเราอ่าน Documentation ครับ เพราะถ้าเกิดมีอะไร พัง เราจะเทียบได้ว่า Pandas ของเราเป็นเวอร์ชันตามใน Documentation มั้ย

```
1. print ("Pandas version",pandas.__version__)
```

วิธีการโหลดไฟล์ CSV (Import)

จุดเริ่มต้นของการทำ Data Exploration & Analysis ใน Pandas ก็คือการโหลดไฟล์ข้อมูลแบบ CSV มาใช้งานนั่นเองครับ

เราสามารถใช้คำสั่ง `.head` หรือ `.tail` เพื่อดูข้อมูลแถวบนสุด หรือแถวล่างสุดได้

```
1. # Read DF
2. csvdf = pd.read_csv('data.csv')
3. # Sometimes reading CSV for Excel need encoding
4. csvdf = pd.read_csv('data.csv',encoding = "ISO-8859-1")
5. # Print head and tail
6. csvdf.head()
7. csvdf.tail()
```

วิธีเช็คจำนวนแถว และจำนวนคอลัมน์

ใน Pandas มีฟังก์ชันสำหรับนับจำนวนแถว และจำนวนคอลัมน์แบบง่าย ๆ โดยใช้

```
1. csvdf.shape
```

วิธีสุ่มข้อมูลสำหรับเช็ค (Sample)

ปกติเราเช็คข้อมูลว่าถูกต้องมั้ยด้วย `head` กับ `tail` ซึ่งเป็นการเช็คจากด้านบนหรือด้านล่าง อีกวิธีที่น่าสนใจ คือ เช็คแบบสุ่มข้อมูลขึ้นมาดูนั่นเองครับ ทำได้ง่าย ๆ โดยใช้

```
1. csvdf.sample()
```

วิธีเช็คข้อมูลหาความผิดปกติใน DataFrame เบื้องต้น

หลังจากโหลดข้อมูลมาแล้ว เราอยากรู้ว่าข้อมูลมีกี่แถว, Missing value เท่าไหร่, แต่ละคอลัมน์เป็น Data Type อะไรบ้าง ก็รันคำสั่งนี้ได้เลย มีประโยชน์มากครับ

```
1. df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4600 entries, 0 to 4599
Data columns (total 20 columns):
date                4600 non-null object
price               4600 non-null int64
bedrooms            4600 non-null int64
bathrooms           4600 non-null float64
sqft_living          4600 non-null int64
sqft_lot             4600 non-null int64
floors              4600 non-null float64
waterfront          4600 non-null int64
view                4600 non-null int64
condition            4600 non-null int64
sqft_above           4600 non-null int64
sqft_basement        4600 non-null int64
yr_built             4600 non-null int64
yr_renovated         224 non-null datetime64[ns]
street              4600 non-null object
city                4600 non-null object
statezip            4600 non-null object
country             4600 non-null object
state               4600 non-null object
zipcode             4600 non-null object
dtypes: datetime64[ns](1), float64(2), int64(10), object(7)
memory usage: 718.8+ KB

```

df.info() จะแสดงสรุปข้อมูลมาให้

นอกจากนั้นยังมีคำสั่ง df.dtypes (ไม่มีวงเล็บ) สำหรับดู Data Type แต่ละคอลัมน์อย่างเดียว

วิธีแปลงประเภทข้อมูล (Data Type) ใน Data Frame

บางครั้งประเภทข้อมูลของคอลัมน์เป็น String แต่เราต้องการ Integer หรือเราต้องการ Date เราสามารถแปลงข้อมูลได้ง่าย ๆ ดังนี้เลยครับ

```

1. df['hour'] = pd.to_numeric(df['hour']) # แปลงเป็น Numeric
2. df['hour'] = df['hour'].astype('int') # อีกวิธีในการแปลงค่า สามารถใช้วิธีนี้แปลงเป็น float ได้

```

วิธีเช็ค Summary ของแต่ละคอลัมน์ (count, min, max, mean)

ถ้าเราอยากรู้ Distribution คร่าว ๆ ของแต่ละคอลัมน์ว่าเป็นอย่างไร สามารถใช้คำสั่ง describe() ได้

```

1. df.describe()

```

วิธีเช็ค Summary (count, min, max, mean) แบบแยกกลุ่ม

บางครั้งเราไม่ได้ต้องการรู้ Summary ของทั้งคอลัมน์ แต่อยากให้แยกตามแต่ละค่าในคอลัมน์นั้น ๆ ครับ ซึ่งมีประโยชน์มากเวลาเราทำ Data Analysis แล้วอยากรู้ว่าบางกลุ่มมีอะไรผิดปกติหรือเปล่า

```
1. test = df.groupby(['Gender'])
2. test.describe()
```

วิธีสร้าง DataFrame ใหม่

วิธีสร้างแบบง่ายที่สุด ถ้าต้องการข้อมูลหลายรูปแบบ เราสามารถใช้ Dictionary แบบนี้เลยครับ

```
1. dataframe = pandas.DataFrame({
2.     'C1': pandas.date_range('20170101', periods=4),
3.     'C2' : [10,20,30,40],
4.     'C3': pandas.Categorical(['A','B','C','D']),
5.     'C4': 1})
```

แต่ถ้าเราต้องการแค่เป็นแบบตัวเลขทั่วไป ใช้ Numpy แบบนี้ได้เลย

```
1. array = numpy.array([(1,2,3), (4,5,6), (7,8,9)])
2. dataframe = pandas.DataFrame(array, columns=['C1', 'C2', 'C3'])
```

วิธีเลือกหลายคอลัมน์จาก DataFrame

ปกติถ้าเราต้องการเลือกแค่ 1 Column ก็เขียนแบบนี้ได้เลย

```
1. df['C1']
```

แต่ถ้าต้องการเลือกหลายคอลัมน์ ให้ทำแบบนี้

```
1. df[['C1', 'C2']]
```

วิธีเลือกคอลัมน์ตามเงื่อนไขที่ต้องการ

บางทีเราอยาก Filter เฉพาะคอลัมน์ที่มีค่าตามที่เรต้องการโดยใช้ .loc ได้ โดยสามารถเลือก Filter แบบ .all() (ทุกค่าในคอลัมน์ต้องตรงตามเงื่อนไข) หรือ .any() (บางค่าในคอลัมน์ต้องตรงตามเงื่อนไข)

```
1. dataframe2 = dataframe.loc[:,(dataframe>50).any()]
2. dataframe3 = dataframe.loc[:,(dataframe>50).all()]
```

เราสามารถใช้หาคอลัมน์ที่มี Missing Values หรือหาคอลัมน์ที่ไม่มี Missing Values เลยก็ได้

```
1. dataframe2 = dataframe.loc[:,dataframe.isnull().any()]
2. dataframe3 = dataframe.loc[:,dataframe.notnull().all()]
```

วิธีเลือกแถวตามเงื่อนไขที่ต้องการ

```
1. dataframe[dataframe['C1']>50] # เงื่อนไขแบบง่าย ๆ
2. dataframe2 = dataframe.loc[dataframe.C1.isin([1,2,3])] # เงื่อนไขแบบซับซ้อน
```

ถ้ามีหลายเงื่อนไขเราสามารถใส่ & (and) หรือ | (or) ได้

```
1. dataframe[(dataframe['C1']>50) & ((dataframe['C2']<25) |
(dataframe['C2']>75))]
```

หรือใช้ Query เป็นเงื่อนไขได้ด้วย มีประโยชน์มากเวลาเรามีเงื่อนไขแปลก ๆ ไม่ต้องเขียนลูปขึ้นมาเองเลยครับ

```
1. dataframe2 = dataframe.query('C1 > C2')
```

วิธีเพิ่มคอลัมน์ใหม่

สามารถเพิ่มคอลัมน์ใหม่ได้ 2 แบบ คือ

1. เพิ่มโดยอิงจากคอลัมน์เดิม (เช่น เอาคอลัมน์เดิม + 10 หรือ เอาคอลัมน์ A - คอลัมน์ B มีประโยชน์มากตอนทำ Feature Engineering)
2. เพิ่มคอลัมน์โดยตั้งค่า Fix ไปเลยสำหรับทุกแถว ส่วนใหญ่จะใช้วิธีนี้เวลาเราอยากได้ค่าอะไรแปลก ๆ ที่ต้องเขียนลูปเพื่อใส่ค่า ก็สร้างคอลัมน์แบบ Fix ค่าก่อน แล้วต่อด้วยลูป

```
1. df['new'] = dataframe['old'] + 10 # use old values
2. df['new2'] = 5 # apply the same value
```

การสลับ Row <-> Column (Transpose)

ถ้าเราต้องการ Transpose (อาร์เรย์เหมือน Vector) เราสามารถใช้คำสั่งนี้ได้เลย

```
1. dataframe.T
```

การต่อ DataFrame

การต่อ Data Frame คือการเอา Data Set 2 ชุดมาต่อกันในแถวตั้งหรือแนวนอน สำหรับการต่อแบบปะติดไปเลย

มี 2 คำสั่งที่เหมือนกัน คือ concat กับ append แต่ให้ใช้ concat ไปเลย เพราะ append เป็นคำสั่งที่ไม่ Memory Efficient

```
1. pd.concat([df1,df2], axis=1) # รวมกัน 2 คอลัมน์ (axis = 0 คือแถว, axis = 1 คือ คอลัมน์)
2. pd.concat([df1,df2,df3]) # รวมมากกว่า 2 คอลัมน์ก็ได้
3. pd.concat(..., ignore_index=True) # รวมเสร็จแล้ว reset index ให้อยู่ในลำดับใหม่
4. pd.concat(..., join='inner') # รวมเฉพาะคอลัมน์ที่ df1 กับ df2 มีทั้งคู่
5. pd.concat(..., keys=['source1', 'source2']) # เพิ่มคอลัมน์เข้าไปด้วยเพื่อระบุว่า Row แต่ละอันมาจาก Data Frame อันไหน
6. pd.concat(..., join_axes=[df2.index]) # เลือกเฉพาะ row index ที่เรากำหนดได้
```

การต่อ DataFrame แบบ Join

ถ้าต้องการต่อ DataFrame แบบ Advance หน่อย เราก็สามารถ Join DataFrame ได้เหมือน Join Table ใครเขียน SQL มาก่อนน่าจะถนัดเลย

```
1. pd.merge(df1, df2, left_on="col1", right_on="col2", how="inner")
```

เราสามารถเปลี่ยนตรง how="inner" เป็น "outer", "left", "right" เพื่อเปลี่ยนเป็น Outer Join, Left Join, Right Join ได้อีกด้วย

การหาค่า Mean, Sum, Max (Aggregate) แบบทั้ง DataFrame

Pandas สามารถสั่ง Aggregate เพื่อหาค่า Mean, Sum, และ Max ได้เลย เหมาะมากเวลาเราต้องการรวบรวมข้อมูลก่อนเอาไป Visualize หรือต้องการทำ Feature Engineering ก็ได้

```
1. newdf = df.agg(['sum', 'max', 'mean'])
```

การ Aggregate แบบตามกลุ่มที่ต้องการ

บางทีเราอยาก Aggregate ข้อมูลตามการจัดกลุ่มในคอลัมน์อื่น เช่น เราอยากได้รายจ่ายทั้งหมดของแต่ละคน (ต้อง aggregate sum ของคอลัมน์รายจ่าย โดยแบ่งกลุ่มตามคอลัมน์ User ID) ใช้แบบนี้

```
1. aggregate = dataframe.groupby('C1').sum()
```

การรัน Function เดียวกันทุกแถว หรือทุกคอลัมน์

เวลาเราอยากรันคำสั่งอะไรสักอย่างสำหรับทุกแถว หรือทุกคอลัมน์ เราสามารถเขียนได้แบบนี้

```
1. # sum for columns
2. sum_columns = dataframe[['C1', 'C2']].apply(sum, axis=0)
3. # sum for rows
4. sum_rows = dataframe[['C1', 'C2']].apply(sum, axis=1)
```

เหมือนกับฟังก์ชัน apply() ใน R นั่นเอง

รันคำสั่งที่เขียนเองกับทุกแถวใน 1 คอลัมน์

ถ้าต้องการรันคำสั่ง (Function) ที่เขียนเอง สำหรับทุกแถวในคอลัมน์อันใดอันหนึ่ง ใช้แบบนี้ได้

```
1. dataframe['C1'] = dataframe['C1'].map(lambda x: x-100)
```


รันคำสั่งที่เขียนเองกับทุกค่า

ถ้าต้องการรันคำสั่งที่เขียนเองกับทุกค่าใน DataFrame ใช้โค้ดนี้

```
1. function_result = dataframe.applymap(lambda x: x*10)
```

หรือใช้ transform ก็ได้

```
1. new_dataframe = dataframe.transform(lambda x: x*100)
```

คำนวณ Correlation & Covariance

เวลาเราอยากรู้ว่าค่าต่าง ๆ ใน Data Set เรา Correlate กันมัย

```
1. dataframe.corr() # Correlation
2. dataframe.cov() # Covariance
```

แต่ค่าที่ออกมาเป็นตัวเลขอาจจะดูยากนิดนึง เราสามารถพลอตสวย ๆ ด้วย Seaborn ได้ครับ สามารถใช้โค้ดด้านล่างนี้ได้เลย

```
1. import seaborn as sns
2.
3. corr = modeldf.corr()
4.
5. # Set up the matplotlib figure
6. f, ax = plt.subplots(figsize=(15, 8))
7.
8. # Generate a custom diverging colormap
9. cmap = sns.diverging_palette(10, 10, as_cmap=True)
10.
11. # Draw the heatmap with the mask and correct aspect ratio
12. sns.heatmap(corr, annot=True)
```



คำนวณ Cross Tabulation

Cross Tabulation มีประโยชน์มากเวลาเราอยากรู้ว่ามี Data ที่ตรงกับกรุป A ของคอลัมน์ 1 และกรุป B ของคอลัมน์ 2 เท่าไหร่ เช่น มีนักเรียนผู้ชาย (คอลัมน์ gender) ที่คนในมัธยมปลาย (คอลัมน์ education) แบบนี้เป็นต้น

หรือถ้าใครใช้ PivotTable ใน Excel มาก่อน ก็เหมือนกันเลยครับ

```
1. aggregate = pandas.crosstab(dataframe.C1, dataframe.C2)
```

วิธีหาค่า Unique ในแต่ละคอลัมน์

คำสั่งนี้มีประโยชน์มาก เอาไว้ใช้เช็คค่าแต่ละคอลัมน์มีค่าแปลก ๆ มั้ย

ตัวอย่างการใช้งานก็คือ เราอยากรู้ว่า มีบ้านไหนที่มีจำนวนห้องนอนแปลก ๆ มั้ย (เช่น 50 ห้องนอน หรือ -5 ห้องนอน) ก็หาค่า unique จากคอลัมน์ "bedrooms"

```
1. dataframe['C1'].unique()
```

วิธีเช็คว่ามีแถวไหนข้อมูลซ้ำมัย (Duplicated)

อันนี้มีประโยชน์มาก เอาไว้ใช้เช็คว่ามีข้อมูลแปลก ๆ มั้ย เช่น ทุกคอลัมน์ซ้ำกันหมด (อันนี้มีโอกาสว่าเป็นข้อมูลซ้ำ อาจจะต้องลบออก) หรือซ้ำกันบางคอลัมน์ (อันนี้ต้องเช็คอีกทีว่าคืออะไร)

```
1. dataframe.duplicated() # หาอันที่เหมือนกันทุกคอลัมน์
2. dataframe.duplicated('C1') # หาอันที่ซ้ำกันเฉพาะคอลัมน์ C1
3. dataframe.duplicated(['C1', 'C2']) # หาอันที่ซ้ำกันเฉพาะคอลัมน์ C1 และ C2
```

ปกติแล้วถ้ามีไอเทมซ้ำ คำสั่งนี้จะไม่แสดงไอเทมแรกในกลุ่มที่ซ้ำ (เช่น ถ้า C1=5 มี 2 แถว มันจะแสดงเฉพาะแถวที่ 2) เราสามารถใส่ Argument **keep=False** เข้าไปเพื่อบังคับให้มันแสดงทุกแถวได้

นอกจากนั้นเรายังสามารถนับจำนวนแถวที่ Duplicate และลบทิ้งได้ด้วย

วิธีการนับจำนวน Duplicate

```
1. len(df[ df.duplicated(['A', 'B', 'C'], keep = False) ])
```

วิธีการลบ Duplicate

เอาไว้ใช้ตอนเราเจอว่าทุกคอลัมน์ซ้ำกันหมดเลย ซึ่งเป็นเคสที่บอกว่าข้อมูลน่าจะซ้ำ ลบออกได้ (ขึ้นอยู่กับข้อมูลด้วยนะครับ บางข้อมูลอาจจะไม่ได้แปลว่าซ้ำแล้วลบได้):

```
1. # Remove the duplicates
2. df.drop_duplicates(['A', 'B', 'C'], inplace=True)
3.
4. # Reset dataframe index after drop_duplicates.
5. df.reset_index(drop=True, inplace=True)
6.
7. len(df)
```

สำหรับโค้ดข้างบน จะเห็นว่าเราต้อง reset index หลังลบ duplicate ด้วยนะครับ

วิธีการลบแถว และลบคอลัมน์

ลบคอลัมน์สามารถทำได้แบบนี้

```
1. dataframe = dataframe.drop('C1', axis=1)
2. df.drop(['C1'], axis=1, inplace=True) # แบบนี้ก็ใช้ได้
3. df.drop(['C1', 'C2', 'C3'], 1, inplace=True) # ลบทีละหลายคอลัมน์ก็ได้
```

ส่วนการลบแถวจะลำบากหน่อย เพราะต้องใส่ Row Index (เลขที่อยู่ซ้ายสุดเวลาเราพิมพ์ DataFrame)

```
1. dataframe = dataframe.drop(5, axis=0)
2. dataframe.reset_index(drop=True) # Reset index
```

ลบแถวแล้วอย่าลืมเช็คด้วยว่าที่ลบไปถูกต้องมั้ย และหลังจากลบแถวต้อง Reset Index ด้วย

วิธีการลบแถวที่มี Missing Value

ข้อควรระวัง: การที่อยู่ ๆ เราลบแถวที่มี Missing Value ทิ้งไปเลยอาจจะไม่ใช่วิธีที่ดีที่สุดในการทำ Data Analysis เสมอไปนะครับ บางเคสการ Impute (คำนวณหาค่าไปใส่) จะดีกว่าครับ

```
1. dataframe2 = dataframe.dropna(axis=0)
```

วิธีแทนค่า Missing Value ด้วยค่าเฉลี่ย (Mean Imputation)

วิธีหนึ่งในการแทนค่าที่หายไป คือการทำสิ่งที่เรียกว่า Mean Imputation หรือหาค่าเฉลี่ยของคอลัมน์นั้น แล้วเอามาแทนค่าที่หายไปนั่นเองครับ

ข้อดีของการทำ Mean Imputation คือ สามารถทำได้ง่าย แต่ก็ต้องระวังเรื่องข้อเสีย เช่น ทำแบบนี้จะเป็นการไม่สนใจความสัมพันธ์ระหว่างตัวแปร ทำให้เกิด Bias สูง ควรใช้เฉพาะเวลา Missing Value ไม่เยอะเท่านั้นครับ

สามารถรันโค้ดด้านล่างเพื่อทำ Mean Imputation ได้ง่าย ๆ เลย

```
1. import numpy as np
2. meanAge = np.mean(df.Age) # Get mean value
3. df.Age = df.Age.fillna(meanAge) # Fill missing values with mean
```

การลูบข้อมูลแต่ละคอลัมน์ และแต่ละแถว

การลูบมีประโยชน์มากถ้าเราต้องการเขียนฟังก์ชันแปลก ๆ ใช้เองที่ Pandas ไม่รองรับ (หรืออาจจะรองรับแต่เราหาไม่เจอ เขียนเองง่ายกว่า) สามารถลูบได้ทั้งแต่ละคอลัมน์ และแต่ละแถว

```
1. for col_idx, data in dataframe.iteritems():
2.     print ("column:", col_idx)
3.     print ("column data:")
4.     print (data, "\n")
```

การลูบข้อมูลแต่ละแถว

```
1. for col_idx, data in dataframe.iterrows():
2.     print ("row:", col_idx)
3.     print ("row data:")
4.     print (data, "\n")
```

วิธีเปลี่ยน DataFrame จากแบบ Wide เป็น Long (Melt)

การ Melt Data มีประโยชน์มากเวลาเราต้องการเอาข้อมูลไปพลอต Data Visualization หรือเราต้องการ Aggregate ครับ

```
1. dataframe2 = dataframe.melt()
```

วิธีการเปลี่ยนชื่อคอลัมน์ (Rename)

บางที่เราต้องการเปลี่ยนชื่อเพื่อให้สั้นลง ให้พิมพ์สะดวกขึ้น สามารถทำได้ดังนี้

```
1. dataframe.rename(columns={'old': 'new'}, inplace=True)
```

วิธีการใส่คำนำหน้าคอลัมน์ (Prefix)

อันนี้มีประโยชน์มากตอนเรามีข้อมูลหลาย ๆ ชุด และต้องการ Merge โดยอยากให้ชื่อคอลัมน์ไม่ซ้ำกัน

```
1. thisdata = thisdata.add_prefix('data_')
```

วิธีการแทนค่าใน DataFrame

เหมาะมากเวลาต้องการแก้ Typo Error เช่น เราอยากได้ค่า Bangkok แต่เรา发现有คนเขียนเป็น BKK อะไรแบบนี้ (รันคำสั่ง .unique เพื่อดูก่อน)

เราสามารถ Replace ทั้ง DataFrame ได้เลยแบบนี้

```
1. dataframe2 = dataframe.replace(1, -100)
```

เราสามารถ Replace หลายค่าพร้อมกันได้ด้วยครับ และสามารถกำหนด Column ที่ต้องการให้แทนค่าได้ด้วย

```
1. df['city'].replace({
2.     'BKK': 'Bangkok',
3.     'BNK': 'Bangkok'
4. }, inplace=True)
```

วิธีการ Export DataFrame เป็นไฟล์ CSV

หลังจากที่เราจัดการ Data เรียบร้อยแล้ว ก็สามารถ Export เป็น CSV เอาไปใช้ต่อกับโปรแกรมอื่น หรืองานส่วนอื่น ๆ ได้ (แอดทำบ่อยเพราะบางทีต้องสลับ Python <-> R รว้ ๆ)

```
1. dataframe.to_csv('dataframe.csv')
```

สรุปการใช้งาน Pandas

จากเทคนิคต่าง ๆ ด้านบน จะเห็นว่า Pandas มีฟังก์ชันให้เราทำงานได้ง่ายขึ้นมาก ๆ ครับ หวังว่าบทความนี้จะมีประโยชน์กับท่านที่กำลังหัดใช้งาน Pandas กันอยู่นะครับ :)