

Data Analysis Tools with Pandas 2 - SF Salaries Exercise

แบบฝึกหัดนี้เป็นแบบฝึกหัดทดสอบทักษะการใช้งาน library pandas ด้วย Salaries.csv โดยให้ทำการคำสั่งต่อไปนี้

6521603795 ปฏิพิธ เอี่ยมรัมย์ 700

Import pandas as pd.

In [1]: `import pandas as pd`

ให้นำเข้าข้อมูลจากไฟล์ Salaries.csv มาในรูปของ dataframe โดยตั้งชื่อตัวแปรว่า sal

In [61]: `sal = pd.read_csv('Salaries.csv')`

Check the head of the DataFrame.

In [62]: `sal.head()`

			Id	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalP
0	1		NATHANIEL FORD	GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY	167411.18	0.00	400184.25	NaN	567595.43		
1	2		GARY JIMENEZ	CAPTAIN III (POLICE DEPARTMENT)	155966.02	245131.88	137811.38	NaN	538909.28		
2	3		ALBERT PARDINI	CAPTAIN III (POLICE DEPARTMENT)	212739.13	106088.18	16452.60	NaN	335279.91		
3	4		CHRISTOPHER CHONG	WIRE ROPE CABLE MAINTENANCE MECHANIC	77916.00	56120.71	198306.90	NaN	332343.61		
4	5		PATRICK GARDNER	DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT)	134401.60	9737.00	182234.59	NaN	326373.19		

ใช้คำสั่ง `.info()` method to ในการดูภาพรวมของข้อมูลทั้งหมด

In [63]: `sal.info()`

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 148654 entries, 0 to 148653
Data columns (total 13 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Id               148654 non-null   int64  
 1   EmployeeName     148654 non-null   object  
 2   JobTitle         148654 non-null   object  
 3   BasePay          148045 non-null   float64 
 4   OvertimePay      148650 non-null   float64 
 5   OtherPay         148650 non-null   float64 
 6   Benefits          112491 non-null   float64 
 7   TotalPay         148654 non-null   float64 
 8   TotalPayBenefits 148654 non-null   float64 
 9   Year              148654 non-null   int64  
 10  Notes             0 non-null       float64 
 11  Agency            148654 non-null   object  
 12  Status            0 non-null       float64 

dtypes: float64(8), int64(2), object(3)
memory usage: 14.7+ MB

```

ลบคอลัมน์ Notes และ Status ออกร

```
In [64]: sal = sal.drop(columns=['Notes', 'Status'])
```

```
In [65]: sal.head()
```

	Id	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalP
0	1	NATHANIEL FORD	GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY	167411.18	0.00	400184.25	NaN	567595.43	
1	2	GARY JIMENEZ	CAPTAIN III (POLICE DEPARTMENT)	155966.02	245131.88	137811.38	NaN	538909.28	
2	3	ALBERT PARDINI	CAPTAIN III (POLICE DEPARTMENT)	212739.13	106088.18	16452.60	NaN	335279.91	
3	4	CHRISTOPHER CHONG	WIRE ROPE CABLE MAINTENANCE MECHANIC	77916.00	56120.71	198306.90	NaN	332343.61	
4	5	PATRICK GARDNER	DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT)	134401.60	9737.00	182234.59	NaN	326373.19	

หาค่าเฉลี่ยของ Benefits ใน sal

```
In [66]: print(sal['Benefits'].mean())
```

25007.893150829852

ใน sal แทน Benefits ที่เป็น null ด้วย 0

```
In [67]: sal['Benefits'].fillna(value=0)
```

```
Out[67]: 0      0.0
1      0.0
2      0.0
3      0.0
4      0.0
...
148649  0.0
148650  0.0
148651  0.0
148652  0.0
148653  0.0
Name: Benefits, Length: 148654, dtype: float64
```

In [68]: `sal.head()`

	Id	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalP
0	1	NATHANIEL FORD	GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY	167411.18	0.00	400184.25	NaN	567595.43	
1	2	GARY JIMENEZ	CAPTAIN III (POLICE DEPARTMENT)	155966.02	245131.88	137811.38	NaN	538909.28	
2	3	ALBERT PARDINI	CAPTAIN III (POLICE DEPARTMENT)	212739.13	106088.18	16452.60	NaN	335279.91	
3	4	CHRISTOPHER CHONG	WIRE ROPE CABLE MAINTENANCE MECHANIC	77916.00	56120.71	198306.90	NaN	332343.61	
4	5	PATRICK GARDNER	DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT)	134401.60	9737.00	182234.59	NaN	326373.19	

หาค่าเฉลี่ยนของ Benefits ใน sal อีกครั้ง

In [69]: `print(sal['Benefits'].mean())`

25007.893150829852

จงเพิ่ม colum ชื่อ Year(TH) ใน sal ให้เป็นเลขปี พศ

In [70]: `sal['Year(TH)'] = sal['Year'] + 543`

In [71]: `sal.head()`

Out[71]:

			Id	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalP
0	1		NATHANIEL FORD		GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY	167411.18	0.00	400184.25	NaN	567595.43	
1	2		GARY JIMENEZ		CAPTAIN III (POLICE DEPARTMENT)	155966.02	245131.88	137811.38	NaN	538909.28	
2	3		ALBERT PARDINI		CAPTAIN III (POLICE DEPARTMENT)	212739.13	106088.18	16452.60	NaN	335279.91	
3	4		CHRISTOPHER CHONG		WIRE ROPE CABLE MAINTENANCE MECHANIC	77916.00	56120.71	198306.90	NaN	332343.61	
4	5		PATRICK GARDNER		DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT)	134401.60	9737.00	182234.59	NaN	326373.19	

จะเพิ่มคอลัมน์ Level มีค่าเป็น L เมื่อ TotalPayBenefits น้อยกว่า 1 แสน และเป็น H เมื่อมากกว่าเท่ากับ 1 แสน

In [72]:

```
def fn(x):
    if x >= 100000:
        return 'H'
    else:
        return 'L'
```

In [73]:

```
sal['Level'] = sal['TotalPayBenefits'].apply(fn)
```

In [74]:

```
sal.head()
```

				JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalPay
0	1	NATHANIEL FORD		GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY	167411.18	0.00	400184.25	NaN	567595.43	
1	2	GARY JIMENEZ		CAPTAIN III (POLICE DEPARTMENT)	155966.02	245131.88	137811.38	NaN	538909.28	
2	3	ALBERT PARDINI		CAPTAIN III (POLICE DEPARTMENT)	212739.13	106088.18	16452.60	NaN	335279.91	
3	4	CHRISTOPHER CHONG		WIRE ROPE CABLE MAINTENANCE MECHANIC	77916.00	56120.71	198306.90	NaN	332343.61	
4	5	PATRICK GARDNER		DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT)	134401.60	9737.00	182234.59	NaN	326373.19	

เข้าด้วย Id ให้เป็น index

In [75]: `sal.set_index('Id', inplace = True)`

In [76]: `sal.head()`

			EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalPay
Id										
1		NATHANIEL FORD		GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY	167411.18	0.00	400184.25	NaN	567595.43	56
2		GARY JIMENEZ		CAPTAIN III (POLICE DEPARTMENT)	155966.02	245131.88	137811.38	NaN	538909.28	53
3		ALBERT PARDINI		CAPTAIN III (POLICE DEPARTMENT)	212739.13	106088.18	16452.60	NaN	335279.91	33
4		CHRISTOPHER CHONG		WIRE ROPE CABLE MAINTENANCE MECHANIC	77916.00	56120.71	198306.90	NaN	332343.61	33
5		PATRICK GARDNER		DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT)	134401.60	9737.00	182234.59	NaN	326373.19	32

เปลี่ยนชื่อคอลัมน์ Year เป็น Year(Eng)

```
In [77]: sal.rename( columns= {'Year': 'Year(Eng)'}, inplace=True)
```

```
In [78]: sal.head()
```

Out[78]:

	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalPayE
	Id							
1	NATHANIEL FORD	GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY	167411.18	0.00	400184.25	NaN	567595.43	56
2	GARY JIMENEZ	CAPTAIN III (POLICE DEPARTMENT)	155966.02	245131.88	137811.38	NaN	538909.28	53
3	ALBERT PARDINI	CAPTAIN III (POLICE DEPARTMENT)	212739.13	106088.18	16452.60	NaN	335279.91	33
4	CHRISTOPHER CHONG	WIRE ROPE CABLE MAINTENANCE MECHANIC	77916.00	56120.71	198306.90	NaN	332343.61	33
5	PATRICK GARDNER	DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT)	134401.60	9737.00	182234.59	NaN	326373.19	32

เพิ่มคนชื่อ David Copperfield ทำงานเป็น Magician คอลัมน์น่าเป็น null

```
In [97]: sr = pd.Series(['David Copperfield', 'Magician'], index = ['EmployeeName', 'JobTitle'])  
sr
```

```
Out[97]: EmployeeName    David Copperfield  
JobTitle          Magician  
dtype: object
```

```
In [98]: #เข้าถึงแกร  
sal.loc[len(sal.index)+1] = sr
```

```
In [99]: sal.tail()
```

Out[99]:

	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalPayBen
Id								
148651	Not provided	Not provided	NaN	NaN	NaN	NaN	0.00	
148652	Not provided	Not provided	NaN	NaN	NaN	NaN	0.00	
148653	Not provided	Not provided	NaN	NaN	NaN	NaN	0.00	
148654	Joe Lopez	Counselor, Log Cabin Ranch	0.0	0.0	-618.13	0.0	-618.13	-61
148655	David Copperfield	Magician	NaN	NaN	NaN	NaN	NaN	

สร้าง Dataframe ที่ EmployeeName มีราย A , B และ C ซึ่งมี BasePay เป็น 10000 แล้วนำไปรวมกับ sal

In [103...]

```
df = pd.DataFrame([['A', 10000], ['B', 10000], ['C', 10000]], columns=['EmployeeName', 'BasePay'])
df
```

Out[103...]

	EmployeeName	BasePay
148656	A	10000
148657	B	10000
148658	C	10000

In [105...]

```
sal = pd.concat([sal, df])
```

In [106...]

```
sal.tail()
```

Out[106...]

	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalPayBen
Id								
148654	Joe Lopez	Counselor, Log Cabin Ranch	0.0	0.0	-618.13	0.0	-618.13	-61
148655	David Copperfield	Magician	NaN	NaN	NaN	NaN	NaN	
148656	A	NaN	10000.0	NaN	NaN	NaN	NaN	
148657	B	NaN	10000.0	NaN	NaN	NaN	NaN	
148658	C	NaN	10000.0	NaN	NaN	NaN	NaN	

สร้างตาราง sal_not_B ซึ่งเก็บเฉพาะของคนที่ไม่มี BasePay

In [110...]

```
sal_not_B = sal[sal['BasePay'].isnull()]
```

In [111...]

```
sal_not_B.head()
```

Out[111...]

	EmployeeName	JobTitle	BasePay	OvertimePay	OtherPay	Benefits	TotalPay	TotalPayPer
81392	Kevin P Cashman	Deputy Chief 3	NaN	0.0	149934.11	0.00	149934.11	149934.11
84507	Demetrya Mullens	Licensed Vocational Nurse	NaN	0.0	110485.41	20779.00	110485.41	131274.41
84961	Michael M Horan	Park Patrol Officer	NaN	0.0	120000.00	8841.48	120000.00	128841.48
90526	Thomas Tang	Police Officer 3	NaN	0.0	106079.31	0.00	106079.31	106079.31
90787	Michael C Hill	Deputy Sheriff	NaN	0.0	81299.02	23877.53	81299.02	105176.53

In [112...]

```
sal_not_B.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 610 entries, 81392 to 148655
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   EmployeeName    610 non-null    object  
 1   JobTitle         610 non-null    object  
 2   BasePay          0 non-null     float64
 3   OvertimePay     605 non-null    float64
 4   OtherPay         605 non-null    float64
 5   Benefits         605 non-null    float64
 6   TotalPay         609 non-null    float64
 7   TotalPayBenefits 609 non-null    float64
 8   Year(Eng)        609 non-null    float64
 9   Agency            609 non-null    object  
 10  Year(TH)         609 non-null    float64
 11  Level             609 non-null    object  
dtypes: float64(8), object(4)
memory usage: 62.0+ KB
```

In []: