# Data Analysis Tools with Pandas 2 - SF Salaries Exercise

แบบฝึกหัดนี้เป็นแบบฝึกหัดทดสอบทักษะการใช้งาน library pandas ด้วย SF Salaries Dataset (https://www.kaggle.com/kaggle/sf-salaries) จากเว็ปไซต์ Kaggle ให้ทำตามคำสั่ง ต่อไปนี้

---

**Import pandas as pd.**

In [1]:
```
1  import pandas as pd
```

**ให้นำเข้าข้อมูลจากไฟล์ Salaries.csv มาในรูปของ dataframe โดยตั้งชื่อตัวแปรว่า sal**

In [2]:
```
1  sal = pd.read_csv("Salaries.csv")
```

**Check the head of the DataFrame.**

In [3]:
```
1  sal.head()
```

Out[3]:

|   | Id | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits | 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 |

**ใช้คำสั่ง .info() method to ในการดูภาพรวมของข้อมูลทั้งหมด**

In [4]: 

```
1 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 ออก**

In [5]: 

```
1 sal.drop("Notes", axis = 1, inplace = True)
2 sal.drop("Status", axis = 1, inplace = True)
```

In [6]: 

```
1 sal
```

Out[6]:

|   | Id | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits |
|---|----|--------------|----------|---------|-------------|----------|----------|
| 0 | 1 | NATHANIEL FORD | GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00 | 400184.25 | NaN |
| 1 | 2 | GARY JIMENEZ | CAPTAIN III (POLICE DEPARTMENT) | 155966.02 | 245131.88 | 137811.38 | NaN |
| 2 | 3 | ALBERT PARDINI | CAPTAIN III (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | NaN |
| 3 | 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | NaN |
| | | | DEPUTY CHIEF | | | | |

**หาค่าเฉลี่ยของ Benefits ใน sal**

In [7]:
```
1  sal["Benefits"].mean()
```

Out[7]: 25007.893150829852

### ใน sal แทน Benefits ที่เป็น null ด้วย 0

In [8]:
```
1  sal["Benefits"].fillna(value = 0, inplace = True)
```

In [9]:
```
1  sal.head()
```

Out[9]:

|   | Id | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits | TotalPay |
|---|-----|--------------|----------|---------|-------------|----------|----------|----------|
| 0 | 1 | NATHANIEL FORD | GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00 | 400184.25 | 0.0 | 567595.43 |
| 1 | 2 | GARY JIMENEZ | CAPTAIN III (POLICE DEPARTMENT) | 155966.02 | 245131.88 | 137811.38 | 0.0 | 538909.28 |
| 2 | 3 | ALBERT PARDINI | CAPTAIN III (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | 0.0 | 335279.91 |
| 3 | 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | 0.0 | 332343.61 |
| 4 | 5 | PATRICK GARDNER | DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT) | 134401.60 | 9737.00 | 182234.59 | 0.0 | 326373.19 |

### หาค่าเฉลี่ยนของ Benefits ใน sal อีกครั้ง

In [10]:
```
1  sal["Benefits"].mean()
```

Out[10]: 18924.23283887417

### จงเพิ่มคอลัมน์ Year(TH) ใน sal ให้เป็นเลขปี พศ

In [11]:
```
1  sal["Year(TH)"] = sal["Year"] + 543
```

In [12]:
```
1  sal
```

Out[12]:

| | Id | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits |
|---|---|---|---|---|---|---|---|
| **0** | 1 | NATHANIEL FORD | GENERAL MANAGER- METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00 | 400184.25 | 0.0 |
| **1** | 2 | GARY JIMENEZ | CAPTAIN III (POLICE DEPARTMENT) | 155966.02 | 245131.88 | 137811.38 | 0.0 |
| **2** | 3 | ALBERT PARDINI | CAPTAIN III (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | 0.0 |
| **3** | 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | 0.0 |
| | | | DEPUTY CHIEF | | | | |

**จงเพิ่มคอลัมน์ Level มีค่าเป็น L เมื่อ TotalPayBenefits น้อยกว่า 1 แสน และเป็น H เมื่อมากกว่าเท่ากับ 1 แสน**

In [24]:
```
1  def cal(TPB):
2
3      if (TPB < 100000) :
4          return "L"
5      else:
6          return "H"
```

In [25]:
```
1  sal["Level"] = sal["TotalPayBenefits"].apply(cal)
```

In [26]:  1  sal

| | | | BasePay | OvertimePay | OtherPay | Benefits | TotalP |
|---|---|---|---|---|---|---|---|
| 3 | ALBERT PARDINI | (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | 0.0 | 335279 |
| 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | 0.0 | 332343 |
| 5 | PATRICK GARDNER | DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT) | 134401.60 | 9737.00 | 182234.59 | 0.0 | 326373 |
| ... | ... | ... | ... | ... | ... | ... | |
| 148654 | Joe Lopez | Counselor, Log Cabin Ranch | 0.00 | 0.00 | -618.13 | 0.0 | -618 |
| 148656 | David Copperfield | Magician | NaN | NaN | NaN | NaN | N |
| 0 | A | NaN | 10000.00 | NaN | NaN | NaN | N |

### เซ็ต Id ให้เป็น index

In [13]:  1  sal.set_index("Id", inplace = True)

In [14]:  1  sal

Out[14]:

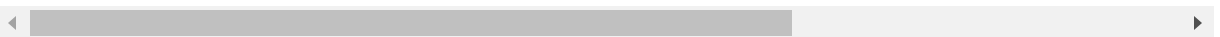| | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits | TotalP |
|---|---|---|---|---|---|---|---|
| **Id** | | | | | | | |
| 1 | NATHANIEL FORD | GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00 | 400184.25 | 0.0 | 567595 |
| 2 | GARY JIMENEZ | CAPTAIN III (POLICE DEPARTMENT) | 155966.02 | 245131.88 | 137811.38 | 0.0 | 538909 |
| 3 | ALBERT PARDINI | CAPTAIN III (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | 0.0 | 335279 |
| 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | 0.0 | 332343 |

### เปลี่ยนชื่อคอลัมน์ Year เป็น Year(Eng)

In [15]:
```
1  sal.rename(columns = {"Year" : "Year(Eng)"})
```

Out[15]:

| Id | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits | TotalPay |
|---|---|---|---|---|---|---|---|
| 1 | NATHANIEL FORD | GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00 | 400184.25 | 0.0 | 567595.43 |
| 2 | GARY JIMENEZ | CAPTAIN III (POLICE DEPARTMENT) | 155966.02 | 245131.88 | 137811.38 | 0.0 | 538909.28 |
| 3 | ALBERT PARDINI | CAPTAIN III (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | 0.0 | 335279.91 |
| 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | 0.0 | 332343.61 |
| 5 | PATRICK GARDNER | DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT) | 134401.60 | 9737.00 | 182234.59 | 0.0 | 326373.19 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 148650 | Roy I Tillery | Custodian | 0.00 | 0.00 | 0.00 | 0.0 | 0.00 |
| 148651 | Not provided | Not provided | NaN | NaN | NaN | 0.0 | 0.00 |
| 148652 | Not provided | Not provided | NaN | NaN | NaN | 0.0 | 0.00 |
| 148653 | Not provided | Not provided | NaN | NaN | NaN | 0.0 | 0.00 |
| 148654 | Joe Lopez | Counselor, Log Cabin Ranch | 0.00 | 0.00 | -618.13 | 0.0 | -618.13 |

148654 rows × 11 columns

### เพิ่มคนชื่อ David Copperfield ทำงานเป็น Magician คอลัมน์อื่นๆเป็น null

In [18]:
```
1  sal.loc["148656", "EmployeeName"] = "David Copperfield"
2  sal.loc["148656", "JobTitle"] = "Magician"
```
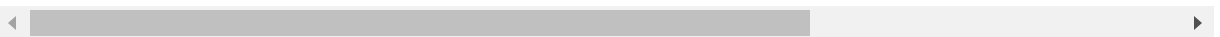
In [19]:
```
1 sal
```

Out[19]:

| Id | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits | TotalPay |
|---|---|---|---|---|---|---|---|
| 1 | NATHANIEL FORD | GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00 | 400184.25 | 0.0 | 567595.43 |
| 2 | GARY JIMENEZ | CAPTAIN III (POLICE DEPARTMENT) | 155966.02 | 245131.88 | 137811.38 | 0.0 | 538909.28 |
| 3 | ALBERT PARDINI | CAPTAIN III (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | 0.0 | 335279.91 |
| 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | 0.0 | 332343.61 |
| 5 | PATRICK GARDNER | DEPUTY CHIEF OF DEPARTMENT, (FIRE DEPARTMENT) | 134401.60 | 9737.00 | 182234.59 | 0.0 | 326373.19 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 148651 | Not provided | Not provided | NaN | NaN | NaN | 0.0 | 0.00 |
| 148652 | Not provided | Not provided | NaN | NaN | NaN | 0.0 | 0.00 |
| 148653 | Not provided | Not provided | NaN | NaN | NaN | 0.0 | 0.00 |
| 148654 | Joe Lopez | Counselor, Log Cabin Ranch | 0.00 | 0.00 | -618.13 | 0.0 | -618.13 |
| 148656 | David Copperfield | Magician | NaN | NaN | NaN | NaN | NaN |

148655 rows × 11 columns

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

In [20]:
```python
df = pd.DataFrame({"EmployeeName" : ["A", "B", "C"],
                   "BasePay" : [10000, 10000, 10000]},
index = [0, 1, 2])
df
```

Out[20]:

|   | EmployeeName | BasePay |
|---|---|---|
| 0 | A | 10000 |
| 1 | B | 10000 |
| 2 | C | 10000 |

In [21]:
```python
sal = pd.concat([sal, df])
sal
```

Out[21]:

| | EmployeeName | JobTitle | BasePay | OvertimePay | OtherPay | Benefits | TotalF |
|---|---|---|---|---|---|---|---|
| 1 | NATHANIEL FORD | GENERAL MANAGER-METROPOLITAN TRANSIT AUTHORITY | 167411.18 | 0.00 | 400184.25 | 0.0 | 567595 |
| 2 | GARY JIMENEZ | CAPTAIN III (POLICE DEPARTMENT) | 155966.02 | 245131.88 | 137811.38 | 0.0 | 538909 |
| 3 | ALBERT PARDINI | CAPTAIN III (POLICE DEPARTMENT) | 212739.13 | 106088.18 | 16452.60 | 0.0 | 335279 |
| 4 | CHRISTOPHER CHONG | WIRE ROPE CABLE MAINTENANCE MECHANIC | 77916.00 | 56120.71 | 198306.90 | 0.0 | 332343 |
| | | DEPUTY CHIEF | | | | | |

**สร้างตาราง salB ซึ่งเก็บเฉพาะของคนที่ไม่มี BasePay**

In [22]:
```python
salB = sal[sal["BasePay"].isnull()]
```

In [23]:  `salB.head()`

Out[23]:

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

In [165]:  `salB.info()`

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 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         609 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
```

## ----- ภาวนามยปัญญา ปัญญาที่เกิดจากการลงมือทำ! -----