

# การจัดการข้อมูลสูญหาย

การรวบรวมข้อมูลมาวิเคราะห์นั้น บางครั้งอาจจะมีข้อมูลที่ได้มา ไม่ครบข้าง ตกหล่นหรือขาดหายไป บ้างเรียกส่วนนี้ว่า Missing Data หรือ Missing Value ในหัวข้อนี้จะมาตรวจสอบข้อมูลและจัดการข้อมูลสูญหาย (Clean Data)

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

df = pd.read_csv("datasets/Employee.csv")
df
```

```
Out[1]:
```

|    | Name | Job        | Age  | Salary   | Bonus | Address |
|----|------|------------|------|----------|-------|---------|
| 0  | A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| 1  | B    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| 2  | C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| 3  | D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| 4  | E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| 5  | F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| 6  | NaN  | NaN        | NaN  | NaN      | NaN   | NaN     |
| 7  | H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| 8  | I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| 9  | H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| 10 | E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

```
In [2]: import pandas as pd

df = pd.read_csv("datasets/Employee.csv", index_col="Name")
df
```

Out [2]:

|             | Job        | Age  | Salary   | Bonus | Address |
|-------------|------------|------|----------|-------|---------|
| <b>Name</b> |            |      |          |       |         |
| <b>A</b>    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| <b>B</b>    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| <b>C</b>    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| <b>D</b>    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| <b>F</b>    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| <b>NaN</b>  | NaN        | NaN  | NaN      | NaN   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>I</b>    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

In [3]: `df.shape`

Out[3]: (11, 5)

In [4]: `df.info()`

```
<class 'pandas.core.frame.DataFrame'>
Index: 11 entries, A to E
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Job         8 non-null      object
1   Age         8 non-null      float64
2   Salary      9 non-null      float64
3   Bonus       10 non-null     object
4   Address     1 non-null      float64
dtypes: float64(3), object(2)
memory usage: 528.0+ bytes
```

In [5]: `#การตรวจสอบข้อมูลสูญหายด้วย isnull()`  
`df.isnull()`

Out [5]:

|  | Job | Age | Salary | Bonus | Address |
|--|-----|-----|--------|-------|---------|
|--|-----|-----|--------|-------|---------|

Name

|            |       |       |       |       |       |
|------------|-------|-------|-------|-------|-------|
| <b>A</b>   | False | False | False | False | False |
| <b>B</b>   | False | False | True  | False | True  |
| <b>C</b>   | False | True  | False | False | True  |
| <b>D</b>   | True  | False | False | False | True  |
| <b>E</b>   | False | False | False | False | True  |
| <b>F</b>   | False | True  | False | False | True  |
| <b>NaN</b> | True  | True  | True  | True  | True  |
| <b>H</b>   | False | False | False | False | True  |
| <b>I</b>   | True  | False | False | False | True  |
| <b>H</b>   | False | False | False | False | True  |
| <b>E</b>   | False | False | False | False | True  |

In [6]: *#ตรวจสอบว่ามีคอลัมน์ใดบ้างที่ไม่มีข้อมูล*  
`df.isnull().any()`

Out [6]:

|         |      |
|---------|------|
| Job     | True |
| Age     | True |
| Salary  | True |
| Bonus   | True |
| Address | True |
| dtype:  | bool |

In [7]: *#นับจำนวนคอลัมน์ที่ไม่มีข้อมูล*  
`df.isnull().sum()`

Out [7]:

|         |       |
|---------|-------|
| Job     | 3     |
| Age     | 3     |
| Salary  | 2     |
| Bonus   | 1     |
| Address | 10    |
| dtype:  | int64 |

In [8]: *#การตรวจสอบข้อมูลครบถ้วนด้วย notnull()*  
`df.notnull()`

Out [8]:

|  | Job | Age | Salary | Bonus | Address |
|--|-----|-----|--------|-------|---------|
|--|-----|-----|--------|-------|---------|

Name

|            |       |       |       |       |       |
|------------|-------|-------|-------|-------|-------|
| <b>A</b>   | True  | True  | True  | True  | True  |
| <b>B</b>   | True  | True  | False | True  | False |
| <b>C</b>   | True  | False | True  | True  | False |
| <b>D</b>   | False | True  | True  | True  | False |
| <b>E</b>   | True  | True  | True  | True  | False |
| <b>F</b>   | True  | False | True  | True  | False |
| <b>NaN</b> | False | False | False | False | False |
| <b>H</b>   | True  | True  | True  | True  | False |
| <b>I</b>   | False | True  | True  | True  | False |
| <b>H</b>   | True  | True  | True  | True  | False |
| <b>E</b>   | True  | True  | True  | True  | False |

In [9]: *#ตรวจสอบว่ามีคอลัมน์ใดบ้างที่มีข้อมูล*  
`df.notnull().any()`

Out [9]:

|         |      |
|---------|------|
| Job     | True |
| Age     | True |
| Salary  | True |
| Bonus   | True |
| Address | True |
| dtype:  | bool |

In [10]: *#นับจำนวนคอลัมน์ที่มีข้อมูล*  
`df.isnull().sum()`

Out [10]:

|         |       |
|---------|-------|
| Job     | 3     |
| Age     | 3     |
| Salary  | 2     |
| Bonus   | 1     |
| Address | 10    |
| dtype:  | int64 |

## การจัดการข้อมูลสูญหาย

- แทนที่ด้วยค่าเฉลี่ยข้อมูลทั้งหมด
- แทนที่ด้วยค่าตรงๆที่กำหนดขึ้นมา
- แทนที่ด้วยค่าก่อนหน้า
- แทนที่ด้วยค่าถัดไป
- ลบข้อมูล

## แทนที่ด้วยค่าเฉลี่ยข้อมูลทั้งหมด

In [11]: `df.describe()`

Out[11]:

|              | Age       | Salary        | Address |
|--------------|-----------|---------------|---------|
| <b>count</b> | 8.000000  | 9.000000      | 1.0     |
| <b>mean</b>  | 26.625000 | 53000.000000  | 123.0   |
| <b>std</b>   | 5.998512  | 23097.618925  | NaN     |
| <b>min</b>   | 18.000000 | 30000.000000  | 123.0   |
| <b>25%</b>   | 22.250000 | 40000.000000  | 123.0   |
| <b>50%</b>   | 27.500000 | 40000.000000  | 123.0   |
| <b>75%</b>   | 30.250000 | 60000.000000  | 123.0   |
| <b>max</b>   | 34.000000 | 100000.000000 | 123.0   |

In [12]: `#นำเข้า DataFrame ใหม่`  
`df = pd.read_csv("datasets/Employee.csv", index_col="Name")`  
`df`

Out[12]:

|             | Job        | Age  | Salary   | Bonus | Address |
|-------------|------------|------|----------|-------|---------|
| <b>Name</b> |            |      |          |       |         |
| <b>A</b>    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| <b>B</b>    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| <b>C</b>    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| <b>D</b>    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| <b>F</b>    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| <b>NaN</b>  | NaN        | NaN  | NaN      | NaN   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>I</b>    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

In [13]: `#แทนที่ด้วยค่าเฉลี่ยข้อมูลทั้งหมด`  
`df['Salary'] = df['Salary'].fillna(df['Salary'].mean())`  
`df`

Out [13]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | 53000.0  | 10%   | NaN     |
| C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| NaN  | NaN        | NaN  | 53000.0  | NaN   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

## แทนที่ด้วยค่าตรงๆที่กำหนดขึ้นมา

In [14]: `#นำเข้า DataFrame ใหม่`  
`df = pd.read_csv("datasets/Employee.csv", index_col="Name")`  
`df`

Out [14]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| NaN  | NaN        | NaN  | NaN      | NaN   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

```
In [15]: #แทนที่ด้วยค่าตรงๆที่กำหนดขึ้นมา
df['Salary'] = df['Salary'].fillna(22000)
df
```

```
Out[15]:
```

|             | Job        | Age  | Salary   | Bonus | Address |
|-------------|------------|------|----------|-------|---------|
| <b>Name</b> |            |      |          |       |         |
| <b>A</b>    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| <b>B</b>    | Programmer | 18.0 | 22000.0  | 10%   | NaN     |
| <b>C</b>    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| <b>D</b>    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| <b>F</b>    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| <b>NaN</b>  | NaN        | NaN  | 22000.0  | NaN   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>I</b>    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

## แทนที่ด้วยค่าก่อนหน้า

```
In [16]: #นำเข้า DataFrame ใหม่
df = pd.read_csv("datasets/Employee.csv", index_col="Name")
df
```

Out [16]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| NaN  | NaN        | NaN  | NaN      | NaN   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

In [17]:

```
#แทนที่ด้วยค่าก่อนหน้า
df.fillna(method='pad')
```

```
/var/folders/83/3fg00w11r7bf7rcsz4nznlh0000gn/T/ipykernel_21816/3352054385.
py:2: FutureWarning: DataFrame.fillna with 'method' is deprecated and will r
aise in a future version. Use obj.ffill() or obj.bfill() instead.
df.fillna(method='pad')
```

Out [17]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | 30000.0  | 10%   | 123.0   |
| C    | Developer  | 18.0 | 32000.0  | 10%   | 123.0   |
| D    | Developer  | 23.0 | 40000.0  | 10%   | 123.0   |
| E    | Sale       | 29.0 | 40000.0  | 10%   | 123.0   |
| F    | Manager    | 29.0 | 75000.0  | 10%   | 123.0   |
| NaN  | Manager    | 29.0 | 75000.0  | 10%   | 123.0   |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | 123.0   |
| I    | Maketing   | 26.0 | 100000.0 | 10%   | 123.0   |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | 123.0   |
| E    | Sale       | 29.0 | 40000.0  | 10%   | 123.0   |



```
In [18]: #นำเข้า DataFrame ใหม่
df = pd.read_csv("datasets/Employee.csv", index_col="Name")
df
```

```
Out[18]:
```

|             | Job        | Age  | Salary   | Bonus | Address |
|-------------|------------|------|----------|-------|---------|
| <b>Name</b> |            |      |          |       |         |
| <b>A</b>    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| <b>B</b>    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| <b>C</b>    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| <b>D</b>    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| <b>F</b>    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| <b>NaN</b>  | NaN        | NaN  | NaN      | NaN   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>I</b>    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

```
In [19]: #แทนที่ด้วยค่าถัดไป
df.fillna(method='bfill')
```

```
/var/folders/83/3fg00w111r7bf7rcsz4nznlh0000gn/T/ipykernel_21816/2964409054.
py:2: FutureWarning: DataFrame.fillna with 'method' is deprecated and will r
aise in a future version. Use obj.ffill() or obj.bfill() instead.
df.fillna(method='bfill')
```

Out [19]:

|             | Job        | Age  | Salary   | Bonus | Address |
|-------------|------------|------|----------|-------|---------|
| <b>Name</b> |            |      |          |       |         |
| <b>A</b>    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| <b>B</b>    | Programmer | 18.0 | 32000.0  | 10%   | NaN     |
| <b>C</b>    | Developer  | 23.0 | 32000.0  | 10%   | NaN     |
| <b>D</b>    | Sale       | 23.0 | 40000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| <b>F</b>    | Manager    | 34.0 | 75000.0  | 10%   | NaN     |
| <b>NaN</b>  | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>I</b>    | Maketing   | 26.0 | 100000.0 | 10%   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

## ลบข้อมูล

- ลบทั้งทั้งหมด
- ลบแถวบางส่วน
- ลบคอลัมน์บางส่วน
- ลบค่าซ้ำ

## ลบทั้งทั้งหมด

```
In [20]: #นำเข้า DataFrame ใหม่
df = pd.read_csv("datasets/Employee.csv", index_col="Name")
df
```

Out [20]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| NaN  | NaN        | NaN  | NaN      | NaN   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

In [21]: `#ลบทิ้งทั้งหมด`  
`df.dropna()`

Out [21]:

|      | Job        | Age  | Salary  | Bonus | Address |
|------|------------|------|---------|-------|---------|
| Name |            |      |         |       |         |
| A    | Programmer | 20.0 | 30000.0 | 10%   | 123.0   |

ลบแถวบางส่วนที่มีค่าว่าง

In [22]: `#นำเข้า DataFrame ใหม่`  
`df = pd.read_csv("datasets/Employee.csv", index_col="Name")`  
`df`

Out [22]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| NaN  | NaN        | NaN  | NaN      | NaN   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

In [23]:

```
#ลบแถวบางส่วนที่มีค่าว่าง
df.dropna(subset=['Age', 'Job'])
```

Out [23]:

|      | Job        | Age  | Salary  | Bonus | Address |
|------|------------|------|---------|-------|---------|
| Name |            |      |         |       |         |
| A    | Programmer | 20.0 | 30000.0 | 10%   | 123.0   |
| B    | Programmer | 18.0 | NaN     | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0 | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0 | 10%   | NaN     |

ลบคอลัมน์บางส่วนที่มีค่าว่าง

In [24]:

```
#นำเข้า DataFrame ใหม่
df = pd.read_csv("datasets/Employee.csv", index_col="Name")
df
```

Out [24]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| NaN  | NaN        | NaN  | NaN      | NaN   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

In [25]: `df.dropna(axis='columns')`

Out [25]:

| Name |
|------|
| A    |
| B    |
| C    |
| D    |
| E    |
| F    |
| NaN  |
| H    |
| I    |
| H    |
| E    |

ลบค่าซ้ำ

In [26]: `#นำเข้า DataFrame ใหม่`  
`df = pd.read_csv("datasets/Employee.csv", index_col="Name")`  
`df`

Loading [MathJax]/extensions/Safe.js

Out [26]:

|             | Job        | Age  | Salary   | Bonus | Address |
|-------------|------------|------|----------|-------|---------|
| <b>Name</b> |            |      |          |       |         |
| <b>A</b>    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| <b>B</b>    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| <b>C</b>    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| <b>D</b>    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| <b>F</b>    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| <b>NaN</b>  | NaN        | NaN  | NaN      | NaN   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>I</b>    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |
| <b>H</b>    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| <b>E</b>    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |

In [27]:

```
#เช็คค่าซ้ำ
df[df.duplicated]
```

Out [27]:

|             | Job      | Age  | Salary  | Bonus | Address |
|-------------|----------|------|---------|-------|---------|
| <b>Name</b> |          |      |         |       |         |
| <b>H</b>    | Maketing | 34.0 | 60000.0 | 10%   | NaN     |
| <b>E</b>    | Sale     | 29.0 | 40000.0 | 10%   | NaN     |

In [28]:

```
#ลบค่าซ้ำ
df.drop_duplicates()
```

Out [28]:

|      | Job        | Age  | Salary   | Bonus | Address |
|------|------------|------|----------|-------|---------|
| Name |            |      |          |       |         |
| A    | Programmer | 20.0 | 30000.0  | 10%   | 123.0   |
| B    | Programmer | 18.0 | NaN      | 10%   | NaN     |
| C    | Developer  | NaN  | 32000.0  | 10%   | NaN     |
| D    | NaN        | 23.0 | 40000.0  | 10%   | NaN     |
| E    | Sale       | 29.0 | 40000.0  | 10%   | NaN     |
| F    | Manager    | NaN  | 75000.0  | 10%   | NaN     |
| NaN  | NaN        | NaN  | NaN      | NaN   | NaN     |
| H    | Maketing   | 34.0 | 60000.0  | 10%   | NaN     |
| I    | NaN        | 26.0 | 100000.0 | 10%   | NaN     |

In [ ]: