


```
[12]: df=pd.read_excel('D:\Pokemon.xls')
df
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49	49	65	65	45	1	False
1	2	Ivysaur	Grass	Poison	405	60	62	63	80	80	60	1	False
2	3	Venusaur	Grass	Poison	525	80	82	83	100	100	80	1	False
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100	123	122	120	80	1	False
4	4	Charmander	Fire	NaN	309	39	52	43	60	50	65	1	False
...
805	721	Volcanion	Fire	Water	600	80	110	120	130	90	70	6	1
806	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined
807	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined	undefined
808	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
809	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

810 rows x 13 columns

去除最后的无意义行和空行

```
[13]: df=df[:-4]
df
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49	49	65	65	45	1	False
1	2	Ivysaur	Grass	Poison	405	60	62	63	80	80	60	1	False
2	3	Venusaur	Grass	Poison	525	80	82	83	100	100	80	1	False
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100	123	122	120	80	1	False
4	4	Charmander	Fire	NaN	309	39	52	43	60	50	65	1	False
...
801	719	Diancie	Rock	Fairy	600	50	100	150	100	150	50	6	1
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160	110	160	110	110	6	1
803	720	HoopaHoopa Confined	Psychic	Ghost	600	80	110	60	150	130	70	6	1
804	720	HoopaHoopa Unbound	Psychic	Dark	680	80	160	60	170	130	80	6	1
805	721	Volcanion	Fire	Water	600	80	110	120	130	90	70	6	1

806 rows x 13 columns

type2 存在异常的数值取值，可清空

```
[14]: condition=df['Type 2'].apply(lambda x: isinstance(x, (int, float)))
df.loc[condition,'Type 2']=np.nan
df
```

```
[14]:
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49	49	65	65	45	1	False
1	2	Ivysaur	Grass	Poison	405	60	62	63	80	80	60	1	False
2	3	Venusaur	Grass	Poison	525	80	82	83	100	100	80	1	False
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100	123	122	120	80	1	False
4	4	Charmander	Fire	NaN	309	39	52	43	60	50	65	1	False
...
801	719	Diancie	Rock	Fairy	600	50	100	150	100	150	50	6	1
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160	110	160	110	110	6	1
803	720	HoopaHoopa Confined	Psychic	Ghost	600	80	110	60	150	130	70	6	1
804	720	HoopaHoopa Unbound	Psychic	Dark	680	80	160	60	170	130	80	6	1
805	721	Volcanion	Fire	Water	600	80	110	120	130	90	70	6	1

806 rows × 13 columns

数据集中存在重复值

```
[15]: df_clean=df.drop_duplicates(subset=['Name'])
df_clean
```

```
[15]:
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49	49	65	65	45	1	False
1	2	Ivysaur	Grass	Poison	405	60	62	63	80	80	60	1	False
2	3	Venusaur	Grass	Poison	525	80	82	83	100	100	80	1	False
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100	123	122	120	80	1	False
4	4	Charmander	Fire	NaN	309	39	52	43	60	50	65	1	False
...
801	719	Diancie	Rock	Fairy	600	50	100	150	100	150	50	6	1
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160	110	160	110	110	6	1
803	720	HoopaHoopa Confined	Psychic	Ghost	600	80	110	60	150	130	70	6	1
804	720	HoopaHoopa Unbound	Psychic	Dark	680	80	160	60	170	130	80	6	1
805	721	Volcanion	Fire	Water	600	80	110	120	130	90	70	6	1

801 rows × 13 columns

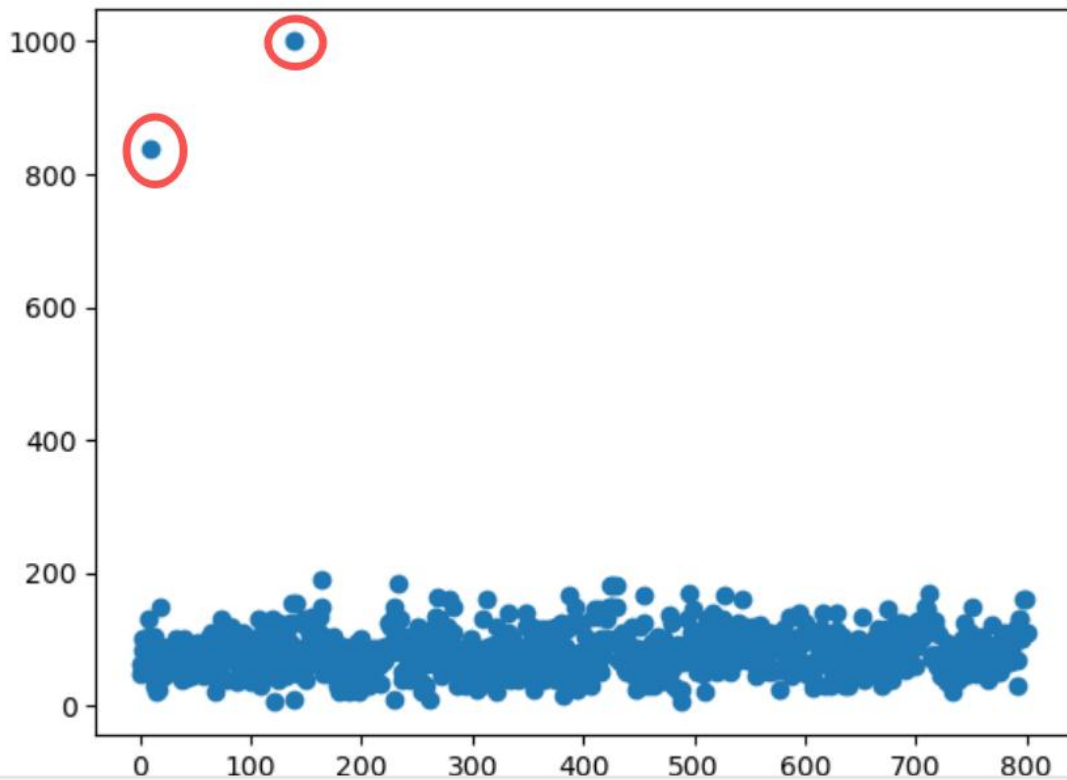
```
[16]: has_dup=df_clean.duplicated().any()
if has_dup:
    print("仍有重复值")
else:
    print("无重复值")
```

无重复值

Attack 属性存在过高的异常值

```
[17]: import matplotlib.pyplot as plt
plt.scatter(range(0,df_clean.shape[0]),df_clean.iloc[:,6])
```

```
[17]: <matplotlib.collections.PathCollection at 0x1f9236d8350>
```



```
[21]: df_clean['Attack'] = pd.to_numeric(df_clean['Attack'], errors='coerce')
df_clean_2 = df_clean[df_clean['Attack'] <= 400].dropna(subset=['Attack'])
df_clean_2
```

```
[21]:
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49.0	49	65	65	45	1	False
1	2	Ivysaur	Grass	Poison	405	60	62.0	63	80	80	60	1	False
2	3	Venusaur	Grass	Poison	525	80	82.0	83	100	100	80	1	False
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100.0	123	122	120	80	1	False
4	4	Charmander	Fire	NaN	309	39	52.0	43	60	50	65	1	False
...
801	719	Diancie	Rock	Fairy	600	50	100.0	150	100	150	50	6	1
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160.0	110	160	110	110	6	1
803	720	HoopaaHoopaa Confined	Psychic	Ghost	600	80	110.0	60	150	130	70	6	1
804	720	HoopaaHoopaa Unbound	Psychic	Dark	680	80	160.0	60	170	130	80	6	1
805	721	Volcanion	Fire	Water	600	80	110.0	120	130	90	70	6	1

798 rows x 13 columns

有数据的 generation 与 Legendary 属性被置换


```
[28]: condition=df_clean_2['Generation'].apply(lambda x: isinstance(x,(int)))
row_n=df_clean_2.index[condition==False].tolist()
row_n
```

[28]: [771]

```
[29]: def swap_values(df,row1,col1,col2):
        df.loc[row1,col1],df.loc[row1,col2]=df.loc[row1,col2],df.loc[row1,col1]
swap_values(df_clean_2,771,'Generation','Legendary')
df_clean_2
```

```
[29]:
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49.0	49	65	65	45	1	False
1	2	Ivysaur	Grass	Poison	405	60	62.0	63	80	80	60	1	False
2	3	Venusaur	Grass	Poison	525	80	82.0	83	100	100	80	1	False
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100.0	123	122	120	80	1	False
4	4	Charmander	Fire	NaN	309	39	52.0	43	60	50	65	1	False
...
801	719	Diancie	Rock	Fairy	600	50	100.0	150	100	150	50	6	1
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160.0	110	160	110	110	6	1
803	720	HoopaHoopa Confined	Psychic	Ghost	600	80	110.0	60	150	130	70	6	1
804	720	HoopaHoopa Unbound	Psychic	Dark	680	80	160.0	60	170	130	80	6	1
805	721	Volcanion	Fire	Water	600	80	110.0	120	130	90	70	6	1

798 rows × 13 columns

结论分析与体会：

熟悉了分析数据质量的方法，掌握了多种数据清洗的方法。