

山东大学计算机科学与技术学院

大数据分析与实践课程实验报告

学号：202300130063	姓名：邱琨	班级：23 数据																																																																																																																																																																								
实验题目：实验二 数据质量实践																																																																																																																																																																										
实验学时：2	实验日期：2025/9/26																																																																																																																																																																									
实验目标：本次实验主要围绕宝可梦数据集进行分析，考察在拿到数据后如何对现有的数据进行预处理清洗操作，建立起对于脏数据、缺失数据等异常情况的一套完整流程的认识。																																																																																																																																																																										
实验过程： 1、库的导入与数据的读入 <pre>[7]: import pandas as pd from pandas import DataFrame import numpy as np [9]: primitive_data = pd.read_csv("./Pokemon.csv") primitive_data</pre> <table><tr><th></th><th>#</th><th>Name</th><th>Type 1</th><th>Type 2</th><th>Total</th><th>HP</th><th>Attack</th><th>Defense</th><th>Sp. Atk</th><th>Sp. Def</th><th>Speed</th><th>Generation</th><th>Legendary</th></tr><tr><td>0</td><td>1</td><td>Bulbasaur</td><td>Grass</td><td>Poison</td><td>318</td><td>45</td><td>49</td><td>49</td><td>65</td><td>65</td><td>45</td><td>1</td><td>FALSE</td></tr><tr><td>1</td><td>2</td><td>Ivysaur</td><td>Grass</td><td>Poison</td><td>405</td><td>60</td><td>62</td><td>63</td><td>80</td><td>80</td><td>60</td><td>1</td><td>FALSE</td></tr><tr><td>2</td><td>3</td><td>Venusaur</td><td>Grass</td><td>Poison</td><td>525</td><td>80</td><td>82</td><td>83</td><td>100</td><td>100</td><td>80</td><td>1</td><td>FALSE</td></tr><tr><td>3</td><td>3</td><td>VenusaurMega Venusaur</td><td>Grass</td><td>Poison</td><td>625</td><td>80</td><td>100</td><td>123</td><td>122</td><td>120</td><td>80</td><td>1</td><td>FALSE</td></tr><tr><td>4</td><td>4</td><td>Charmander</td><td>Fire</td><td>NaN</td><td>309</td><td>39</td><td>52</td><td>43</td><td>60</td><td>50</td><td>65</td><td>1</td><td>FALSE</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>805</td><td>721</td><td>Volcanion</td><td>Fire</td><td>Water</td><td>600</td><td>80</td><td>110</td><td>120</td><td>130</td><td>90</td><td>70</td><td>6</td><td>TRUE</td></tr><tr><td>806</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td></tr><tr><td>807</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td><td>undefined</td></tr><tr><td>808</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr><tr><td>809</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td><td>NaN</td></tr></table> <p>810 rows x 13 columns</p>				#	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	TRUE	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
	#	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	TRUE																																																																																																																																																													
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																																																																																																																																																													
2、删除多余的空行并去除 undefined 行																																																																																																																																																																										

```
[12]: primitive_data_1=primitive_data.dropna(how='any')
primitive_data_1 = primitive_data_1[~(primitive_data_1 == 'undefined').all(axis=1)]
primitive_data_1
```

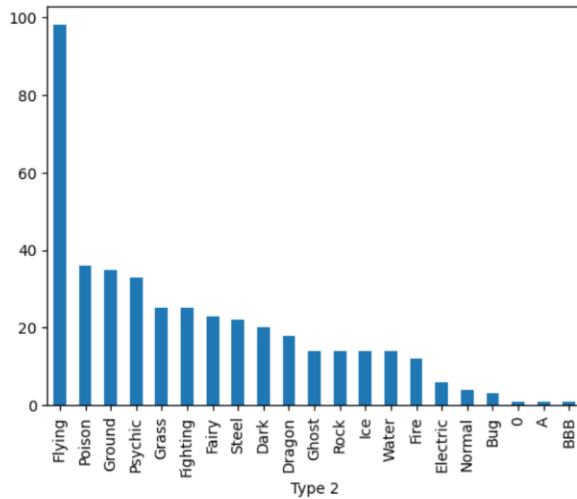
	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49	49	65	65	45	1	FALSE
1	2	Ivysaur	Grass	Poison	405	60	62	63	80	80	60	1	FALSE
2	3	Venusaur	Grass	Poison	525	80	82	83	100	100	80	1	FALSE
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100	123	122	120	80	1	FALSE
6	6	Charizard	Fire	Flying	534	78	84	78	109	85	100	1	FALSE
...
801	719	Diancie	Rock	Fairy	600	50	100	150	100	150	50	6	TRUE
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160	110	160	110	110	6	TRUE
803	720	HoopaHoopa Confined	Psychic	Ghost	600	80	110	60	150	130	70	6	TRUE
804	720	HoopaHoopa Unbound	Psychic	Dark	680	80	160	60	170	130	80	6	TRUE
805	721	Volcanion	Fire	Water	600	80	110	120	130	90	70	6	TRUE

419 rows × 13 columns

3、删除异常值

```
[9]: data["Type 2"].value_counts().plot(kind="bar")
```

[9]: <Axes: xlabel='Type 2'>



```
[12]: data = data[~data["Type 2"].isin(["0", "A", "BBB"])]
data
```

	#	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
6	6	Charizard	Fire	Flying	534	78	84	78	109	85	100	1	FALSE
...
801	719	Diancie	Rock	Fairy	600	50	100	150	100	150	50	6	TRUE
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160	110	160	110	110	6	TRUE
803	720	HoopaHoopa Confined	Psychic	Ghost	600	80	110	60	150	130	70	6	TRUE
804	720	HoopaHoopa Unbound	Psychic	Dark	680	80	160	60	170	130	80	6	TRUE
805	721	Volcanion	Fire	Water	600	80	110	120	130	90	70	6	TRUE

416 rows × 13 columns

4、删除重复值

```
[13]: data[data.duplicated()]
```

	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
23	17	Pidgeotto	Normal	Flying	349	63	60	55	50	50	71	1	FALSE
185	168	Ariados	Bug	Poison	390	70	90	70	60	60	40	2	FALSE
186	168	Ariados	Bug	Poison	390	70	90	70	60	60	40	2	FALSE
187	168	Ariados	Bug	Poison	390	70	90	70	60	60	40	2	FALSE

```
[14]: data = data.drop_duplicates()  
data
```

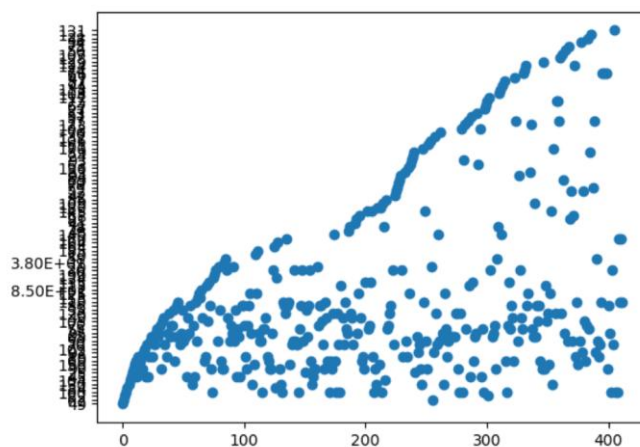
	#	Name	Type 1	Type 2	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0	1	Bulbasaur	Grass	Poison	318	45	49	49	65	65	45	1	FALSE
1	2	Ivysaur	Grass	Poison	405	60	62	63	80	80	60	1	FALSE
2	3	Venusaur	Grass	Poison	525	80	82	83	100	100	80	1	FALSE
3	3	VenusaurMega Venusaur	Grass	Poison	625	80	100	123	122	120	80	1	FALSE
6	6	Charizard	Fire	Flying	534	78	84	78	109	85	100	1	FALSE
...
801	719	Diancie	Rock	Fairy	600	50	100	150	100	150	50	6	TRUE
802	719	DiancieMega Diancie	Rock	Fairy	700	50	160	110	160	110	110	6	TRUE
803	720	HoopaHoopa Confined	Psychic	Ghost	600	80	110	60	150	130	70	6	TRUE
804	720	HoopaHoopa Unbound	Psychic	Dark	680	80	160	60	170	130	80	6	TRUE
805	721	Volcanion	Fire	Water	600	80	110	120	130	90	70	6	TRUE

412 rows x 13 columns

5、Attack 属性存在过高的异常值

```
[21]: plt.scatter(range(0, data.shape[0]), data.iloc[:, 6])
```

```
[21]: <matplotlib.collections.PathCollection at 0x279e02ff770>
```



结论分析：

本实验围绕数据质量实践，旨在掌握数据预处理与脏数据清洗核心流程，建立对数据质量问题的系统认知。实验先加载原始数据集并探索结构，明确问题后开展针对性清洗——剔除无效记录、修正异常值、去重、调整字段错位，还补充处理隐性异常以验证清洗完整性。通过实验，加深了对数据预处理基础保障作用的理解，掌握核心技巧，为后续分析建模奠定基础。