

山东大学 计算机科学与技术 学院

大数据分析实践 课程实验报告

学号：202300130028	姓名：苗雨健	班级：数据 23
实验题目：数据质量实践		
实验学时：2	实验日期：2025/9/26	
<div>实验目的：</div> <div>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</div>		
<div>硬件环境：</div> <div>计算机一台</div>		
<div>软件环境：</div> <div>Linux 或 Windows</div>		
<div>实验步骤与内容：</div> <div>首先加载必要库，在正常导入时遇到问题，我们采用 latin-1 格式进行导入</div> <div><pre>1 # 导入必要的库 2 import pandas as pd 3 import numpy as np 4 import matplotlib.pyplot as plt 5 import seaborn as sns 1 2 df = pd.read_csv("C:\\Users\\33566\\Downloads\\Pokemon.csv", encoding='latin-1') 3</pre></div> <div>删去末尾空行</div>		

```

1 # 2. 查看并删除末尾无意义数据
2 print("最后5行数据:")
3 print(df.tail())
4
5 # 删除最后四行无意义数据
6 df_clean = df.iloc[:-4].copy()

```

最后5行数据:

	#	Name	Type 1	Type 2	Total	HP	\
805	721	Volcanion	Fire	Water	600	80	
806	undefined	undefined	undefined	undefined	undefined	undefined	
807	undefined	undefined	undefined	undefined	undefined	undefined	
808	NaN	NaN	NaN	NaN	NaN	NaN	
809	NaN	NaN	NaN	NaN	NaN	NaN	

	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	\
805	110	120	130	90	70	6	
806	undefined	undefined	undefined	undefined	undefined	undefined	

然后对 type2 进行处理

```

1 # 3. 处理Type 2列的异常值
2 # 查看Type 2列的唯一值
3 print("Type 2列唯一值:")
4 print(df_clean['Type 2'].unique())
5
6 # 查看Type 2列的值分布
7 type2_counts = df_clean['Type 2'].value_counts()
8 print("\nType 2值分布:")
9 print(type2_counts)
10
11 # 将错误项替换为NaN
12 df_clean['Type 2'] = df_clean['Type 2'].replace(['A', '273', '0', 'BBB'], np.nan)
13

```

Type 2列唯一值:

```

['Poison' nan 'Flying' 'Dragon' '0' 'Ground' '273' 'Fairy' 'Grass'
 'Fighting' 'Psychic' 'Steel' 'Ice' 'A' 'Rock' 'Dark' 'Water' 'Electric'
 'Fire' 'Ghost' 'Bug' 'BBB' 'Normal']

```

Type 2值分布:

```

Type 2
Flying      98
Poison      37
Ground      35
Psychic     33

```

去除重复

```

1 # 4. 检查并删除重复行
2 duplicates = df_clean.duplicated()
3 print(f"重复行数量: {duplicates.sum()}")
4
5 if duplicates.sum() > 0:
6     print("重复的行:")
7     print(df_clean[duplicates])
8     # 删除重复行, 保留第一个出现的数据
9     df_clean = df_clean.drop_duplicates()
10
11 print(f"删除重复值后数据形状: {df_clean.shape}")

```

```

185  168    Ariados    Bug  Poison    390  7
186  168    Ariados    Bug  Poison    390  7
187  168    Ariados    Bug  Poison    390  7

```

```

      Speed Generation Legendary
15      30           1    FALSE
23      71           1    FALSE
185     40           2    FALSE
186     40           2    FALSE
187     40           2    FALSE

```

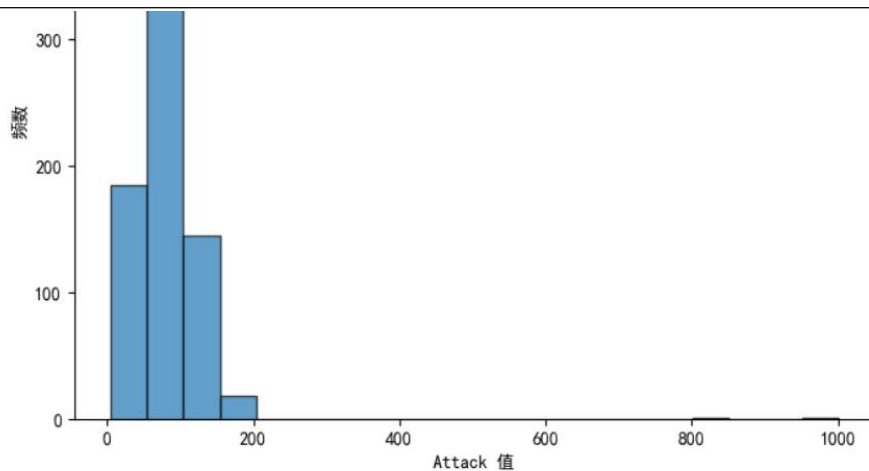
删除重复值后数据形状: (801, 13)

处理 attack 异常值

```

1 # 5. 分析Attack属性分布并处理异常值
2 # 设置中文字体
3 plt.rcParams['font.sans-serif'] = ['SimHei', 'Microsoft YaHei', 'DejaVu Sans']
4 plt.rcParams['axes.unicode_minus'] = False
5
6 # 转换数据类型为数值型
7 df_clean['Attack'] = pd.to_numeric(df_clean['Attack'], errors='coerce')
8
9 # 绘制Attack属性直方图
10 plt.figure(figsize=(8, 6))
11 plt.hist(df_clean['Attack'].dropna(), bins=20, alpha=0.7, edgecolor='black')
12 plt.title('Attack 属性分布')
13 plt.xlabel('Attack 值')
14 plt.ylabel('频数')
15 plt.show()
16
17 # 使用IQR方法识别异常值
18 Q1 = df_clean['Attack'].quantile(0.25)
19 Q3 = df_clean['Attack'].quantile(0.75)
20 IQR = Q3 - Q1
21 lower_bound = Q1 - 1.5 * IQR
22 upper_bound = Q3 + 1.5 * IQR
23
24 print(f"Q1: {Q1}, Q3: {Q3}, IQR: {IQR}")
25 print(f"异常值边界: [{lower_bound}, {upper_bound}]")

```



Q1: 55.0, Q3: 100.0, IQR: 45.0

异常值边界: [-12.5, 167.5]

Attack 异常值数量: 9

异常值数据:

	Name	Attack
9	Squirtle	840.0
140	Tauros	1000.0
165	MewtwoMega Mewtwo X	190.0
237	HonooMega Honoo	185.0

修正 Generation 和 Legendary 属性的置换错误

```

10 1 # 6. 修正Generation和Legendary属性的置换错误
2 # 查找并交换置换的属性值
3 for idx, row in df_clean.iterrows():
4     gen_val = str(row['Generation'])
5     leg_val = str(row['Legendary'])
6     # 判断是否存在属性置换 (Generation为文本, Legendary为数字)
7     if gen_val.isalpha() and leg_val.isdigit():
8         print(f"发现置换行: {row['Name']}")
9         print(f"  原始: Generation={gen_val}, Legendary={leg_val}")
10        # 交换值
11        df_clean.at[idx, 'Generation'] = leg_val
12        df_clean.at[idx, 'Legendary'] = gen_val
13        print(f"  修正后: Generation={leg_val}, Legendary={gen_val}")
14
15 # 转换Generation列的数据类型为数值型
16 df_clean['Generation'] = pd.to_numeric(df_clean['Generation'], errors='coerce')

```

发现置换行: Blastoise
 原始: Generation=FALSE, Legendary=1
 修正后: Generation=1, Legendary=FALSE
 发现置换行: Pikachu
 原始: Generation=FALSE, Legendary=0
 修正后: Generation=0, Legendary=FALSE

结论分析与体会:

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, 建立了对于脏数据、缺失数据等异常情况的一套完整流程的认识