

实训1：CSV数据清洗与统计

• ***目标***：掌握Pandas和数据预处理。

• ***任务***：读取销售数据CSV（包含日期、产品、销量、价格），清洗数据，生成统计报告。

• ***详细要求***：

- o 读取CSV文件，至少100行，字段：日期（YYYY-MM-DD）、产品（字符串）、销量（整数）、价格（浮点数）。
- o 处理缺失值：销量用均值填充，价格用中位数填充。
- o 移除异常值：销量<0或价格>1000的记录。
- o 计算统计指标：总销量、平均价格、按产品分组的销量总和。
- o 输出：清洗后数据保存为cleaned_sales.csv，统计报告保存为sales_report.txt（包含总销量、平均价格、分组统计）。
- o 代码需包含异常处理（如文件不存在、格式错误）。
- o 添加日志记录，输出每次操作的时间和结果。

• ***技能***：Pandas操作、文件I/O、异常处理、日志记录。

```
1  import pandas as pd
2  import numpy as np
3  import logging
4  import os
5  import re
6  from datetime import datetime
7
8  # 配置日志记录
9  logging.basicConfig(
10     level=logging.INFO,
11     format='%(asctime)s - %(levelname)s - %(message)s',
12     handlers=[
13         logging.FileHandler("sales_data_processing.log"),
14         logging.StreamHandler()
15     ]
16 )
17 logger = logging.getLogger()
18
19
20 def clean_and_process_sales_data(input_file="sales_data.csv"):
21     """
22     处理销售数据的完整流程：
23     1. 读取CSV文件
24     2. 清洗数据（处理缺失值、异常值、格式错误）
25     3. 生成统计报告
26     4. 保存清洗后的数据和统计报告
27     """
28     try:
29         # 1. 读取数据
30         logger.info(f"开始处理销售数据，输入文件：{input_file}")
```

```

31     logger.info(f"当前工作目录: {os.getcwd()}")
32
33     # 检查文件是否存在
34     if not os.path.exists(input_file):
35         raise FileNotFoundError(f"文件不存在: {input_file}")
36
37     # 读取CSV文件
38     df = pd.read_csv(input_file)
39     logger.info(f"成功读取数据, 共 {len(df)} 行记录")
40
41     # 检查数据行数
42     if len(df) < 100:
43         logger.warning(f"警告: 文件仅包含 {len(df)} 行数据, 少于要求的100
行")
44
45     # 2. 数据清洗
46     logger.info("开始数据清洗...")
47
48     # 检查列名
49     required_columns = ['date', 'product', 'sales', 'price']
50     missing_cols = [col for col in required_columns if col not in
df.columns]
51     if missing_cols:
52         raise ValueError(f"CSV文件缺少必要的列: {'',
'.join(missing_cols)}")
53
54     # 转换日期格式
55     date_errors = []
56     for i, date_str in enumerate(df['date']):
57         try:
58             # 尝试解析日期
59             pd.to_datetime(date_str, format='%Y-%m-%d')
60         except:
61             date_errors.append(i)
62             # 尝试修复常见日期格式问题
63             if re.match(r'\d{4}/\d{2}/\d{2}', date_str):
64                 df.at[i, 'date'] = date_str.replace('/', '-')
65             elif re.match(r'\d{2}-\d{2}-\d{4}', date_str):
66                 parts = date_str.split('-')
67                 df.at[i, 'date'] = f"{parts[2]}-{parts[0]}-{parts[1]}"
68
69     if date_errors:
70         logger.warning(f"发现并修复了 {len(date_errors)} 条日期格式问题")
71
72     # 确保日期列是datetime类型
73     df['date'] = pd.to_datetime(df['date'], errors='coerce')
74
75     # 处理缺失值
76     # 销量用均值填充
77     if df['sales'].isna().any():
78         # 计算均值时排除负值
79         valid_sales = df[df['sales'] >= 0]['sales']
80         mean_quantity = int(valid_sales.mean()) if not
valid_sales.empty else 0
81         df['sales'] = df['sales'].fillna(mean_quantity)
82         logger.info(f"填充了 {df['sales'].isna().sum()} 个缺失的销量值, 使
用均值: {mean_quantity}")
83

```

```

84     # 价格用中位数填充
85     if df['price'].isna().any():
86         # 计算中位数时排除异常值
87         valid_prices = df[df['price'] <= 1000]['price']
88         median_price = valid_prices.median() if not valid_prices.empty
else 0
89         df['price'] = df['price'].fillna(median_price)
90         logger.info(f"填充了 {df['price'].isna().sum()} 个缺失的价格值，使
用中位数: {median_price:.2f}")
91
92     # 移除异常值
93     initial_count = len(df)
94     # 销量 < 0 或 价格 > 1000
95     df = df[(df['sales'] >= 0) & (df['price'] <= 1000) & (df['price'] >
0)]
96     removed_count = initial_count - len(df)
97     logger.info(f"移除了 {removed_count} 条异常记录（销量<0或价格>1000或价格
<=0）")
98
99     # 处理非数值型价格数据
100    if df['price'].dtype != float and df['price'].dtype != int:
101        # 尝试转换价格为数值型
102        df['price'] = pd.to_numeric(df['price'], errors='coerce')
103        # 再次填充转换失败的价格
104        if df['price'].isna().any():
105            valid_prices = df[df['price'] <= 1000]['price']
106            median_price = valid_prices.median() if not
valid_prices.empty else 0
107            df['price'] = df['price'].fillna(median_price)
108            logger.info(f"转换并填充了 {df['price'].isna().sum()} 个非数值
型价格")
109
110    # 3. 生成统计报告
111    logger.info("计算统计指标...")
112
113    # 总销量
114    total_sales = df['sales'].sum()
115
116    # 平均价格
117    avg_price = df['price'].mean()
118
119    # 按产品分组的销量总和
120    grouped_sales = df.groupby('product')['sales'].sum().reset_index()
121    grouped_sales.columns = ['产品', '总销量']
122
123    # 按日期分组的销量总和
124    daily_sales = df.groupby('date')['sales'].sum().reset_index()
125    daily_sales.columns = ['日期', '总销量']
126
127    # 4. 保存清洗后的数据
128    cleaned_file = "cleaned_sales.csv"
129    df.to_csv(cleaned_file, index=False, encoding='utf-8')
130    logger.info(f"清洗后的数据已保存至: {cleaned_file} ({len(df)} 行记录)")
131
132    # 5. 保存统计报告
133    report_file = "sales_report.txt"
134    with open(report_file, 'w', encoding='utf-8') as f:
135        f.write("=" * 50 + "\n")

```

```

136         f.write("销售数据统计报告\n")
137         f.write(f"生成时间: {datetime.now().strftime('%Y-%m-%d
%H:%M:%S')}}\n")
138         f.write(f"原始数据行数: {initial_count}\n")
139         f.write(f"清洗后数据行数: {len(df)}\n")
140         f.write(f"移除异常记录数: {removed_count}\n")
141         f.write("=" * 50 + "\n\n")
142
143         f.write(f"总销量: {total_sales:},\n")
144         f.write(f"平均价格: {avg_price:.2f}\n\n")
145
146         f.write("按产品分组的销量统计:\n")
147         f.write("-" * 50 + "\n")
148         f.write(grouped_sales.to_string(index=False))
149         f.write("\n\n")
150
151         f.write("按日期分组的销量统计:\n")
152         f.write("-" * 50 + "\n")
153         f.write(daily_sales.to_string(index=False))
154         f.write("\n\n")
155
156         f.write("=" * 50 + "\n")
157         f.write("报告结束\n")
158
159         logger.info(f"统计报告已保存至: {report_file}")
160         logger.info("数据处理完成!")
161
162         return df, total_sales, avg_price, grouped_sales, daily_sales
163
164     except FileNotFoundError as e:
165         logger.error(f"文件错误: {str(e)}")
166         return None, None, None, None, None
167     except ValueError as e:
168         logger.error(f"数据格式错误: {str(e)}")
169         return None, None, None, None, None
170     except Exception as e:
171         logger.exception(f"处理过程中发生未预期的错误: {str(e)}")
172         return None, None, None, None, None
173
174
175 def generate_sample_data():
176     """生成样本数据文件"""
177     logger.info("生成样本数据...")
178
179     # 生成日期范围
180     start_date = datetime(2025, 1, 1)
181     end_date = datetime(2025, 4, 10)
182     dates = [start_date + pd.DateOffset(days=i) for i in range((end_date -
start_date).days + 1)]
183
184     # 生成样本数据
185     np.random.seed(42)
186     products = ['Phone', 'Laptop', 'Tablet', 'Monitor', 'Headphones']
187
188     data = {
189         'date': np.random.choice(dates, 150),
190         'product': np.random.choice(products, 150),
191         'sales': np.random.randint(1, 100, 150),

```

```
192         'price': np.random.uniform(50, 800, 150)
193     }
194
195     # 添加一些缺失值和异常值
196     df = pd.DataFrame(data)
197     df.loc[10:15, 'sales'] = np.nan
198     df.loc[20:25, 'price'] = np.nan
199     df.loc[30, 'sales'] = -5
200     df.loc[35, 'price'] = 1500
201     df.loc[40, 'price'] = "invalid" # 非数值型数据
202
203     # 添加日期格式问题
204     df.loc[45, 'date'] = "2025/02/15" # 错误的分隔符
205     df.loc[50, 'date'] = "02-15-2025" # 错误的格式
206
207     # 保存为CSV
208     df.to_csv('sales_data.csv', index=False)
209     logger.info("样本数据文件 'sales_data.csv' 已创建")
210     return df
211
212
213 if __name__ == "__main__":
214     # 如果数据文件不存在, 生成样本数据
215     if not os.path.exists("sales_data.csv"):
216         logger.info("未找到销售数据文件, 生成样本数据...")
217         generate_sample_data()
218
219     # 处理销售数据
220     cleaned_df, total_sales, avg_price, grouped_sales, daily_sales =
clean_and_process_sales_data()
221
222     # 如果处理成功, 打印部分结果
223     if cleaned_df is not None:
224         print("\n数据处理摘要:")
225         print(f"清洗后记录数: {len(cleaned_df)}")
226         print(f"总销量: {total_sales:,}")
227         print(f"平均价格: {avg_price:.2f}")
228         print("\n按产品分组销量统计:")
229         print(grouped_sales.head())
```