

## 04数据处理与统计方法

---

### 一、Pandas数据统计函数

#### 0.读取数据

```
1 df = pd.read_csv('res/beijing_tianqi_2018.csv')
2 # 替换掉温度的后缀℃
3 df.loc[:, "bwendu"] = df["bwendu"].str.replace("℃",
4 "").astype('int32')
5 df.loc[:, "ywendu"] = df["ywendu"].str.replace("℃",
6 "").astype('int32')
```

#### 1.汇总类统计

```
1 # 一下子提取所有数字列统计结果
2 df.describe()
3 ## 查看单个Series的数据
4 df["bwendu"].mean()
5 # 最高温
6 df["bwendu"].max()
7 # 最低温
8 df["bwendu"].min()
```

#### 2.去重与按值计数

```
1 # 1.唯一性去重
2 # 一般不用于数值列，而是枚举、分类列
3 df["fengxiang"].unique()
4 df["tianqi"].unique()
5 df["fengli"].unique()
6 # 2.按值计数
7 df["fengxiang"].value_counts()
8 df["tianqi"].value_counts()
9 df["fengli"].value_counts()
10
11 # 3.相关系数和协方差
12 '''
13 用途：
14 两只股票，是不是同涨同跌？程度多大？正相关还是负相关？
15 产品销量的波动，跟哪些因素正相关、负相关，程度有多大？
16
17 来自知乎，对于两个变量X、Y：
```

```

18 协方差：衡量同向反向程度，如果协方差为正，说明X，Y同向变化，协方差越大说明同向程度越高；如果协方差为负，说明X，Y反向运动，协方差越小说明反向程度越高。
19 相关系数：衡量相似程度，当他们的相关系数为1时，说明两个变量变化时的正向相似程度最大，当相关系数为-1时，说明两个变量变化的反向相似程度最大
20 '''
21 df.cov() # 协方差矩阵
22 df.corr() # 相关系数矩阵
23 df["aqi"].corr(df["bwendu"]) # 空气质量和最高温度的相关系数

```

## 二、字符串处理

Pandas的字符串处理：

1. 使用方法：先获取Series的str属性，然后在属性上调用函数；
2. 只能在字符串列上使用，不能数字列上使用；
3. Dataframe上没有str属性和处理方法
4. Series.str并不是Python原生字符串，而是自己的一套方法，不过大部分和原生str很相似；

Series.str字符串方法列表参考文档：

<https://pandas.pydata.org/pandas-docs/stable/reference/series.html#string-handling>

### 1. 获取Series的str属性，使用各种字符串处理函数

```

1 df["bwendu"].str
2 # 字符串替换函数
3 df["bwendu"].str.replace("C", "")
4 # 判断是不是数字
5 df["bwendu"].str.isnumeric()
6 df["aqi"].str.len() # 报错

```

### 3. 使用str方法得到的bool值可以做条件查询

```

1 condition = df["ymd"].str.startswith("2018-03")
2 df[condition].head()

```

### 4. 需要多次str处理的链式操作

怎样提取201803这样的数字月份？

- 1、先将日期2018-03-31替换成20180331的形式
- 2、提取月份字符串201803

```

1 df["ymd"].str.replace("-", "")
2 # 每次调用函数，都返回一个新series
3 df["ymd"].str.replace("-", "").slice(0, 6) # 报错
4 df["ymd"].str.replace("-", "").str.slice(0, 6)
5 # slice就是切片语法，可以直接用
6 df["ymd"].str.replace("-", "").str[0:6]

```

## 三、数据排序

**Series**的排序:

```
Series.sort_values(ascending=True, inplace=False)
```

参数说明:

- **ascending**: 默认为**True**升序排序, 为**False**降序排序
- **inplace**: 是否修改原始**Series**

**DataFrame**的排序:

```
DataFrame.sort_values(by, ascending=True, inplace=False)
```

参数说明:

- **by**: 字符串或者**List<字符串>**, 单列排序或者多列排序
- **ascending**: **bool**或者**List**, 升序还是降序, 如果是**list**对应**by**的多列
- **inplace**: 是否修改原始**DataFrame**

### 1.Series的排序

```
1 df["aqi"].sort_values()  
2 df["aqi"].sort_values(ascending=False)
```

### 2.DataFrame的排序

#### 2.1 单列排序

```
1 df.sort_values(by="aqi")  
2 df.sort_values(by="aqi", ascending=False)
```

#### 2.2 多列排序

```
1 # 按空气质量等级、最高温度排序, 默认升序  
2 df.sort_values(by=["aqiLevel", "bwendu"])  
3 # 两个字段都是降序  
4 df.sort_values(by=["aqiLevel", "bwendu"], ascending=False)  
5 # 分别指定升序和降序  
6 df.sort_values(by=["aqiLevel", "bwendu"], ascending=[True,  
False])
```

## 四、新增数据列

在进行数据分析时, 经常需要按照一定条件创建新的数据列, 然后进行进一步分析。

1. 直接赋值
2. **df.apply()**方法
3. **df.assign()**方法

#### 4. 按条件选择分组分别赋值

### 1.直接赋值

```
1 # 注意, df["bwendu"]其实是一个Series, 后面的减法返回的是Series
2 df.loc[:, "wenchu"] = df["bwendu"] - df["ywendu"]
```

### 2.apply方法

添加一行温度类型:

1. 如果最高温度大于33度就是高温
2. 低于-10度是低温
3. 否则是常温

```
1 def get_wendu_type(x):
2     if x["bwendu"] > 33:
3         return '高温'
4     if x["ywendu"] < -10:
5         return '低温'
6     return '常温'
7
8 # 注意需要设置axis==1, 这是series的index是columns
9 df.loc[:, "wendu_type"] = df.apply(get_wendu_type, axis=1)
10 # 查看温度类型的计数
11 df["wendu_type"].value_counts()
```

### 3.assign方法

将温度从摄氏度变成华氏度

```
1 # 可以同时添加多个新的列
2 # 摄氏度转华氏度
3 df.assign(
4     ywendu_huashi = lambda x : x["ywendu"] * 9 / 5 + 32,
5     bwendu_huashi = lambda x : x["bwendu"] * 9 / 5 + 32
6 )
```

#### 4. 按条件选择分组分别赋值

按条件先选择数据, 然后对这部分数据赋值新列

实例: 高低温差大于10度, 则认为温差大

```

1 # 先创建空列（这是第一种创建新列的方法）
2 df['wenchang_type'] = ''
3
4 df.loc[df["bwendu"]-df["ywendu"]>10, "wenchang_type"] = "温差大"
5
6 df.loc[df["bwendu"]-df["ywendu"]<=10, "wenchang_type"] = "温差正常"
7
8 df["wenchang_type"].value_counts()

```

## 五、缺失值处理

**Pandas** 使用这些函数处理缺失值：

- **isnull** 和 **notnull**：检测是否是空值，可用于 **df** 和 **series**
- **dropna**：丢弃、删除缺失值
- **axis**：删除行还是列，{0 or 'index', 1 or 'columns'}
  - **how**：如果等于any则任何值为空都删除，如果等于all则所有值都为空才删除
  - **inplace**：如果为True则修改当前df，否则返回新的df
- **fillna**：填充空值
- **value**：用于填充的值，可以是单个值，或者字典（key是列名，value是值）
- **method**：等于ffill使用前一个不为空的值填充 forward fill；等于bfill使用后一个不为空的值填充 backward fill
- **axis**：按行还是列填充，{0 or 'index', 1 or 'columns'}
- **inplace**：如果为True则修改当前df，否则返回新的df

```

1 # 1.读取数据
2 df = pd.read_excel('res/student_excel.xlsx',skiprows=2)
3 # 2.检测空值
4 df.isnull()
5 df['分数'].isnull()
6 df['分数'].notnull()
7 df.loc[df['分数'].notnull()] # 筛选没有空分数的所有行
8 # 3.删除全是空值的列
9 df.dropna(axis=1,how='all',inplace=True)
10 # 4.删除全是空值的行
11 df.dropna(axis=0,how='all',inplace=True)
12 # 6.将分数列为空的填充为0分
13 df.fillna({'分数':0},inplace=True)
14 # 等同于
15 df.loc[:, '分数'] = df['分数'].fillna(0)
16 # 7.填充姓名的缺失值
17 # 使用前面的有效值填充，用ffill: forward fill
18 df.loc[:, '姓名'] = df['姓名'].fillna(method="ffill")
19 # 8.将清洗好的数据保存
20 df.to_excel('res/student_excel_clean.xlsx',index=False)

```

