

## 简介

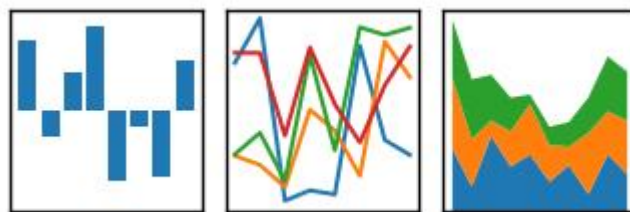
问题：Numpy已经能够帮助我们处理数据，能够结合Matplotlib解决数据分析的展示问题，那么为什么还需要学习Pandas呢？

- numpy 能够帮我们处理处理数值型数据，但是这还不够
- 很多时候，我们的数据除了数值之外，还有**字符串**，还有**时间序列**等

## 1. Pandas是什么？

- Pandas 的名称来自于面板数据（panel data）和 Python 数据分析（data analysis）。
- Pandas 是一个强大的分析结构化数据的工具集，基于 NumPy 构建，提供了 **高级**数据结构和**数据操作工具**，它是使 Python 成为强大而高效的数据分析环境的重要因素之一。
- pandas 是一个开源、BSD 许可的库，为 Python 编程语言提供高性能、易于使用的数据结构 and 数据分析工具。

pandas  
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



## 2. Pandas有哪些特点？

- 灵活的数据结构：支持 Series（一维）、DataFrame（二维），适配 90%+ 结构化数据场景
- 高效数据操作：一行代码实现筛选、排序、分组等复杂操作，比原生 Python 效率提升 10-100 倍
- 丰富的 IO 能力：支持读取 / 保存 CSV、Excel、SQL、JSON 等 10 + 种格式数据
- 无缝集成生态：与 NumPy（数值计算）、Matplotlib（可视化）、Scikit-learn（机器学习）完美兼容

## 安装

可以通过 pip 命令安装 pandas

```
1 # 安装
2 $ pip install pandas
3
4
5 # 导入（惯例用pd作为别名）
6 $ import pandas as pd
```

## 核心数据结构

Pandas 有两大基础数据结构：**Series（一维）** 和 **DataFrame（二维）**，所有操作均围绕二者展开。

### Series

Series：带标签的一维数组。

### 定义与创建

Series 由“索引（index）”和“数据（values）”组成，索引默认从 0 开始递增。

创建主要有以下两种方式：

- 通过列表创建（可以手动指定索引，也可以不指定索引）
- 通过字典创建（key是索引，value是值）

```
1 import pandas as pd
2
3 # Series的定义
4 # 1. 列表创建，默认索引
5 s1 = pd.Series([10, 20, 30, 40, 50])
6 # print(s1)
7
8 # 2. 列表创建，自定义索引
9 s2 = pd.Series([10, 20, 30, 40, 50], index=['a', 'b', 'c', 'd',
10 'e'])
11 # print(s2)
12
13 # 3. 字典创建(key是索引，value是值)
14 s3 = pd.Series({"Alice":85, "Bob":90, "Charlie":78, "David":92})
15 # print(s3)
```

```
15 # 输出:
16 # Alice      85
17 # Bob        90
18 # Charlie    78
19 # David      92
20 # dtype: int64 # 数据类型
```

## 核心属性

核心属性有：

- index: 索引
- values: 数据
- dtype: 数据类型
- shape: 形状

```
1 s3 = pd.Series({"Alice":85, "Bob":90, "Charlie":78, "David":92})
2
3
4 print(s3.index)    # 索引 → Index(['Alice', 'Bob', 'Charlie',
   'David'], dtype='object')
5 print(s3.values)   # 数据 → [85 90 78 92]
6 print(s3.dtype)    # 数据类型 → int64
7 print(s3.shape)    # 形状 (长度) → (4,)
```

## 基本操作

```
1 # 通过索引获取值
2 s3 = pd.Series({"Alice":85, "Bob":90, "Charlie":78, "David":92})
3 print(s3["Alice"])
4 print(s3["Bob"])
5
6 # 通过索引修改值
7 s3["Bob"] = 100
8 print(s3)
```

总结：

- Series对象本质上由两个列表构成
- 一个数组构成对象的键(index, 索引)
- 一个数组构建对象的值(values)

# DataFrame

DataFrame可以理解为是：带标签的二维表格。

DataFrame 类似 Excel 表格，由“行索引（index）”、“列名（columns）”和“数据（values）”组成，是 Pandas 最常用的数据结构。

The diagram illustrates a DataFrame as a 2D table. It features a grid of cells. To the left of the grid, row indices 1, 2, 3, and 4 are listed, with a red arrow pointing to the index label '行索引'. Above the grid, column names '姓名', '年龄', '性别', '身高', '体重', and '成绩' are listed, with a red arrow pointing to the column label '列名'. The grid itself contains data, with some cells containing ellipses (...). A red arrow points to one of the data cells, labeled '数据'.

	姓名	年龄	性别	身高	体重	成绩
1	张三	...				
2	李四		...			
3	王五			...		
4	赵六				...	

## 定义与创建

DataFrame本质其实就是二维数组（ndarray），创建方式有两种：

- 通过字典创建（键为列名，值为列数据，行索引可以自定义，也可以默认）
- 通过二维的ndarray创建

```
1 import pandas as pd
2 import numpy as np
3
4 data = {
5     "name": ["Alice", "Bob", "Charlie", "David"],
6     "age": [20, 22, 21, 20],
7     "score": [85, 90, 78, 92],
8     "gender": ["女", "男", "男", "男"]
9 }
10 # 1. 字典创建（字典键为列名，值为列数据）
11 df = pd.DataFrame(data)
12 print(df)
13
14
15 # 2. 自定义行索引
16 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4'])
17 print(df)
18
19
20 # 3. 通过ndarray构建Dataframe
```

```

21 arr1 = np.arange(12).reshape(3,4)
22 df2 = pd.DataFrame(arr1)
23 print(df2)
24
25 # df2输出:
26 #      0  1  2  3          # 行索引
27 # 0      0  1  2  3
28 # 1      4  5  6  7
29 # 2      8  9 10 11
30 # 列索引
31
32
33 # 注意: 依然具有广播的特性
34 data2 = {
35     "name": ['zs', 'ls', 'ww'],
36     "gender": "男",
37     "height": [180, 180, 190]
38 }
39 df2 = pd.DataFrame(data2)

```

## 核心属性

核心属性有:

- index: 行索引
- columns: 列索引 (列名)
- values: 数据
- shape: 形状
- info: 整体概览

```

1  import pandas as pd
2  import numpy as np
3
4  # 1. 创建一个DataFrame
5  data = {
6      "name": ["Alice", "Bob", "Charlie", "David"],
7      "age": [20, 22, 21, 20],
8      "score": [85, 90, 78, 92],
9      "gender": ["女", "男", "男", "男"]
10 }
11
12 df = pd.DataFrame(data)
13
14 # 2. 输出核心属性
15 # 行索引

```

```

16 print(df.index)      # 输出: RangeIndex(start=0, stop=4, step=1)
17
18 # 列索引
19 print(df.columns)    # 输出: Index(['name', 'age', 'score',
20                                'gender'], dtype='object')
21
22 # 数据
23 print(df.values)      # 输出: [['Alice' 20 85 '女'] ['Bob' 22 90
24                                '男'] ['Charlie' 21 78 '男'] ['David' 20 92 '男']]
25
26 # 形状
27 print(df.shape)      # 输出: (4, 4)
28
29 # 数据概览
30 print(df.info())     # 输出包括: 类型、行索引、列索引、数据、数据类型、内存占
                        用等

```

## 基本操作

- 1 获取前几行(后几行)
- 2 通过索引获取列数据
- 3 增加列
- 4 删除列

```

1 dict = {
2     "name": ["Alice", "Bob", "Charlie", "David"],
3     "age": [20, 22, 21, 20],
4     "score": [85, 90, 78, 92],
5     "gender": ["女", "男", "男", "男"]
6 }
7 df = pd.DataFrame(dict)
8
9 # 1. 获取前几行数据
10 print(df.head(2))
11 print('-'*20)
12 # 2. 获取后几行数据
13 print(df.tail(1))
14 print('-'*20)
15
16 # 3. 获取指定列的数据
17 name_data = df["name"]
18 print(name_data)
19 print(type(name_data))      # 单列数据是 Series
20 print('-'*20)
21

```

```

22 # 4. 获取指定多列的数据
23 age_data = df[["age", "gender"]]
24 print(age_data)
25 print(type(age_data))          # 多列数据是 DataFrame
26 print('-'*20)
27
28 # 5. 增加列
29 df["new_column"] = [1, 2, 3, 4]
30 # print(df)
31 df["new_column2"] = df["new_column"] + 20
32 print(df)
33 print('-'*20)
34
35 # 6. 删除列
36 del df["new_column"]
37 print(df)

```

## 索引操作

索引是 Pandas 高效定位数据的核心，是非常重要的。

### 学习目标：

- 1 掌握Series中的索引的操作
- 2 掌握 **label 索引 (loc)** 和 **位置索引 (iloc)** 的使用与区别
- 3 掌握索引切片

## Series中的索引操作

### 1. Index指定索引的名称

```

1 import pandas as pd
2 import numpy as np
3
4 # 创建一个Series
5 s=pd.Series(['A', 'B', 'C', 'D'])
6
7 # 1. Index指定索引的名称
8 s.index = [1,2,3,4]
9 print(s)

```

## 2. 根据索引名取值

```
1 s2 = pd.Series(['A', 'B', 'C', 'D'], index=['001', '002', '003', '004'])
2
3 # 通过索引名取值
4 print(s2['002'])           # B
5 print(s2.loc['002'])       # B      #推荐的方式
```

## 3. 根据位置取值

```
1 s2 = pd.Series(['A', 'B', 'C', 'D'], index=['001', '002', '003', '004'])
2
3 # 通过位置取值
4 # print(s2[1])             # FutureWarning: 已被废弃
5 print(s2.iloc[1])          # 推荐的方式
```

## 4. 取多个值

```
1 # 3. 通过多个索引取多个值
2 # print(s2[['002', '004']])
3 # print('-'*20)
4
5 print(s2.loc[['002', '004']])
6 print('-'*20)
7
8 # print(s2[[1, 3]])
9
10 print(s2.iloc[[1, 3]])
11 print('-'*30)
```

## 5. 根据切片取值

```
1 print(s2.loc['002':'004'])      # 通过索引名取值是【左闭右闭】
2 print('-'*20)
3 print(s2.iloc[0:2])            # 通过索引位置取值是【左闭右开】
```

总结：Series中的取值有两种方式：

### 1 根据索引的名字取值

- 可以直接 `对象名[索引名]` 来取值
- 建议使用 `对象名.loc[索引名]` 来取值

### 2 根据索引的位置取值

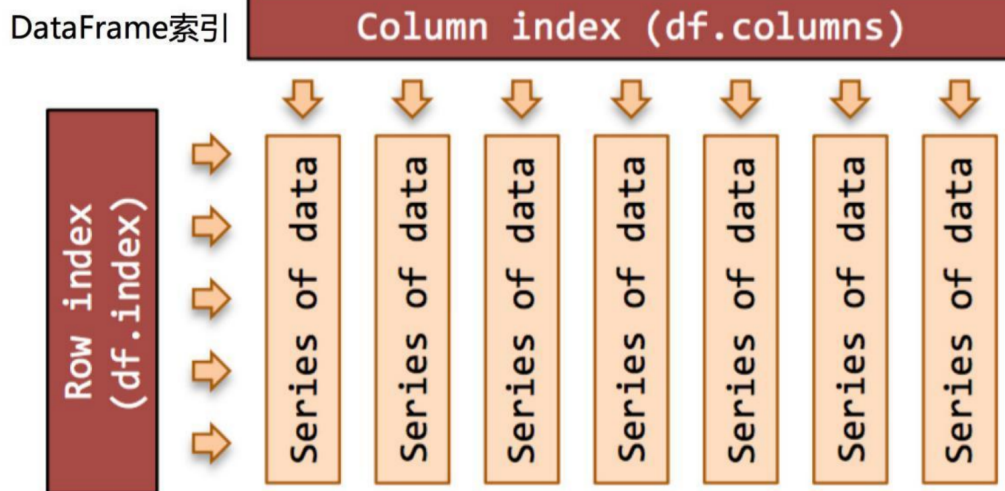
- 可以直接使用 `对象名[位置值]` 来取值
- 建议使用 `对象名.iloc[位置值]` 来取值

### 3 支持使用切片的方式来取值



## DataFrame中的索引操作

首先，认识一下DataFrame的索引，分为行索引和列索引，每一列都构成了一个Series对象。



```
1 df = pd.DataFrame(np.arange(12).reshape(3,4), index=
  ['s1', 's2', 's3'], columns=['c1', 'c2', 'c3', 'c4'])
2 print("原始数据:")
3 print(df)
4
5 # 1. 通过列索引获取值
6 print("通过列索引获取值:")
7 print(df['c2'])
8
9 # 2. 通过多个列索引获取多个列的值
10 print("通过多个列索引获取多个列的值")
11 print(df[['c2', 'c4']])
12
13 # 3. 通过行索引获取行数据
14 print("通过行索引获取行数据:")
15 print(df.loc['s2'])
16
17 # 4. 通过多个行索引获取多个行数据
18 print("通过多个行索引获取多个行数据:")
19 print(df.loc[['s1', 's3']])
20
21 # 5. 通过位置索引获取一行的值
22 print("通过位置索引获取一行的值:")
23 print(df.iloc[1])
24
25 # 6. 通过位置索引获取多行数据
26 print("通过位置索引获取多行数据:")
27 print(df.iloc[[0,2]])
28
29 # 7. 获取具体位置的值
30 print("获取具体位置的值:")
```

```

31 print(df.iloc[1,2])          # 前面是行，后面是列
32 # print(df.iloc[1][2])      # FutureWarning
33
34 # 8. 通过位置索引切片获取数据
35 print("通过位置索引切片获取数据:")
36 print(df.iloc[0:2,1:3])      # 前面是行，后面是列，左闭右开
37
38 # 9. 通过名字索引切片获取数据
39 print("通过名字索引切片获取数据:")
40 print(df.loc['s1':'s3', 'c2':'c4']) # 前面是行，后面是列，左闭右闭
41
42 # 10. 一次取多个不连续的行或者列
43 print(df.iloc[[1,2], [0,3]]) # 前面是行，后面是列
44 #改为loc实现上一行
45 print(df.loc[['s2', 's3'], ['c1', 'c4']])

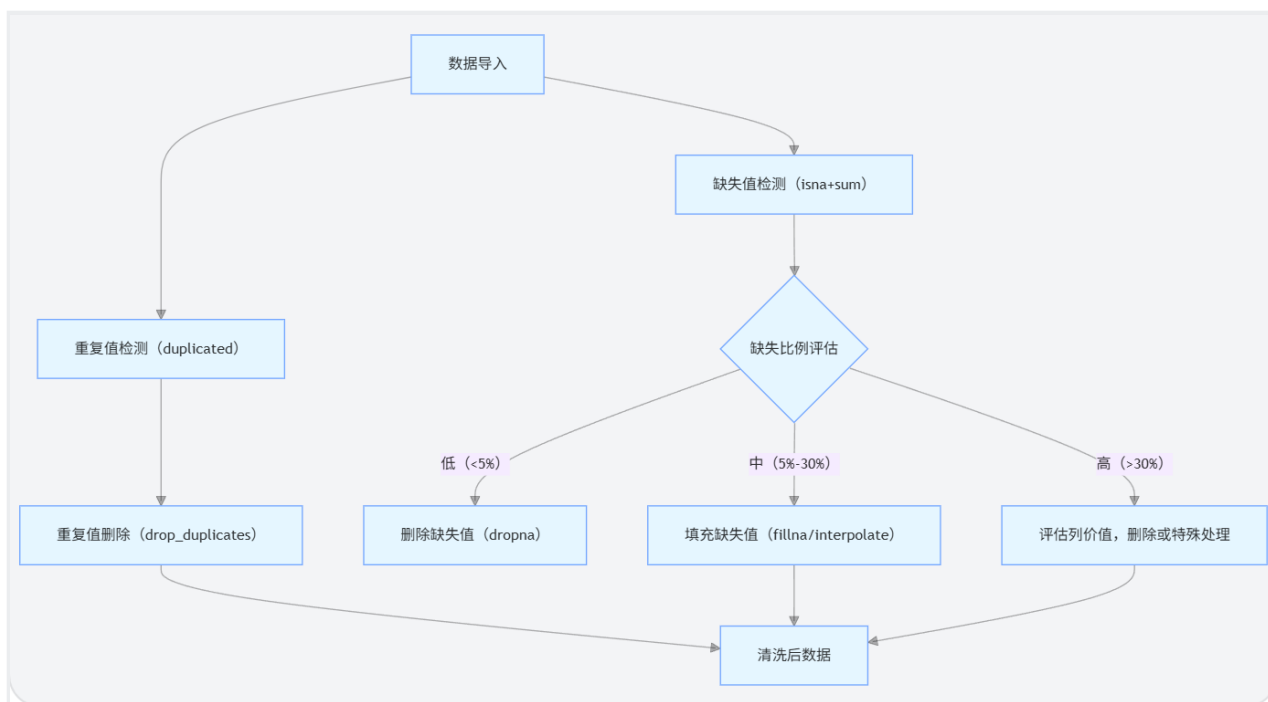
```

## 总结

- 1 `df.loc[]`：根据索引名来取值，如果是切片，左闭右闭
- 2 `df.iloc[]`：根据索引的位置来取值，如果是切片，左闭右开
- 3 在DataFrame中，`df.loc()`和`df.iloc()`默认都是取**行数据**

## 数据清洗

当我们采集到数据之后，数据清洗是提高数据质量的关键步骤，包括处理缺失值、重复值等。



## 重复值处理

重复值会影响分析结果准确性，需先检测再处理，核心函数： `duplicated()`（检测）和  `drop_duplicates()`（删除）

### 1. 检测重复值

```

1  # 创建一个数据表
2  data = {
3      "name": ["Alice", "Alice", "Charlie", "David", "Alice",
4      "Charlie"],
5      "age": [20, 22, 21, 20, 20, 21],
6      "score": [85, 90, 78, 92, 85, 78],
7      "gender": ["女", "男", "男", "男", "女", "男"]
8  }
9  df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4', 's5',
10 's6'], columns=['name', 'age', 'score', 'gender'])
11 print(df)
12
13 # 1. 检测重复值
14 print(df.duplicated())      # 返回布尔Series, True表示重复行
15 print('-'*20)
16
17 # 2. 指定列检测重复值(可以指定多列)
18 # 参数 keep, 默认值为'first', 表示只保留第一个重复行, 'last'表示只保留最后
19 # 一个重复行, False表示标记所有重复行
20 df_dup_by_col = df.duplicated(subset=['name', 'age'],
21 keep='first')
  
```

```
18 print(df_dup_by_col)
```

## 2. 删除重复值

```
1 # 创建一个数据表
2 data = {
3     "name": ["Alice", "Alice", "Charlie", "David", "Alice",
4             "Charlie"],
5     "age": [20, 22, 21, 20, 20, 21],
6     "score": [85, 90, 78, 92, 85, 78],
7     "gender": ["女", "男", "男", "男", "女", "男"]
8 }
9 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4', 's5',
10                               's6'], columns=['name', 'age', 'score', 'gender'])
11
12 # 1. 删除完全重复的行
13 df_drop_dup = df.drop_duplicates()
14 print("删除完全重复的行:")
15 print(df_drop_dup)
16
17 # 2. 删除部分列重复的行
18 df_drop_dup_by_col = df.drop_duplicates(subset=['name'],
19                                          keep='first')
20 print("删除部分列重复的行:")
21 print(df_drop_dup_by_col)
22
23 # 3. 直接在原数据上删除
24 df.drop_duplicates(inplace=True)
25 print("直接在原数据上删除:")
26 print(df)
```

## 缺失值处理

缺失值（NaN, Not a Number）是数据采集常见问题（如漏填、传感器故障），处理核心是“先评估缺失程度，再选修复策略”（删除 / 填充 / 插值）。

### 1. 检测缺失值 `isnull()`

```
1 # 1. 创建一个数据表
2 data = {
3     "name": ["Alice", "Bob", "Charlie", "David", "Eva", "Frank"],
4     "age": [20, 22, 21, 20, 20, np.nan],
5     "score": [85, 90, np.nan, 92, 85, np.nan],
6     "gender": ["女", "男", "男", np.nan, np.nan, np.nan]
```

```

7 }
8 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4', 's5',
    's6'], columns=['name', 'age', 'score', 'gender'])
9 print(df)
10
11 # 2. 检测缺失值
12 print(df.isnull())
13
14 # 3. 计算缺失个数
15 print("每一列缺失个数:")
16 print(df.isnull().sum())
17
18 print("每一行缺失个数:")
19 print(df.isnull().sum(axis=1))
20
21 # 4. 计算缺失率
22 print("每一列缺失率:")
23 print(df.isnull().sum() / len(df))
24
25 print("每一行缺失率:")
26 print(df.isnull().sum(axis=1) / len(df.columns))

```

## 2. 删除缺失值

适用于缺失值占比极低（如 <5%），且删除后不影响样本代表性的场景。

```

1 # 1. 创建一个数据表
2 data = {
3     "name": ["Alice", "Bob", "Charlie", "David", "Eva", "Frank"],
4     "age": [20, 22, 21, 20, 20, np.nan],
5     "score": [85, 90, np.nan, 92, 85, np.nan],
6     "gender": ["女", "男", "男", np.nan, np.nan, np.nan]
7 }
8 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4', 's5',
    's6'], columns=['name', 'age', 'score', 'gender'])
9 print("原数据:")
10 print(df)
11
12 # 2. 删除含有缺失值的行或列 (谨慎操作, 可能会丢失大量数据)
13 df_drop_na = df.dropna()
14 print("删除含有缺失值的行:")
15 print(df_drop_na)
16
17 df_drop_na_by_col = df.dropna(axis=1)
18 print("删除含有缺失值的列:")
19 print(df_drop_na_by_col)
20
21 # 3. 删除 score数据缺失的行
22 df_drop_na_by_col = df.dropna(subset=['score'])
23 print("删除 score数据缺失的行:")

```

```
24 print(df_drop_na_by_col)
```

### 3. 填充缺失值

```
1 # 1. 创建一个数据表
2 data = {
3     "name": ["Alice", "Bob", "Charlie", "David", "Eva", "Frank"],
4     "age": [20, 22, 21, 20, 20, np.nan],
5     "score": [85, 90, np.nan, 92, 85, np.nan],
6     "gender": ["女", "男", "男", np.nan, np.nan, np.nan]
7 }
8 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4', 's5', 's6'], columns=['name', 'age', 'score', 'gender'])
9 print("原数据:")
10 print(df)
11
12 # 2. 填充全局缺失值
13 df_fill_na = df.fillna(0)
14 print("填充缺失值为0:")
15 print(df_fill_na)
16
17 # 3. 填充某一列的缺失值
18 df_fill_na_by_col = df.fillna({'gender': '男'})
19 print("填充【gender】缺失值为男:")
20 print(df_fill_na_by_col)
```

## 数据转换：函数应用

当内置函数无法满足需求时，需用自定义函数批量处理数据，核心工具：`apply()`、`map()`、`applymap()`等

## apply

`apply` 是 Pandas 中最灵活的转换函数，**既支持 Series（一维），也支持 DataFrame（二维）**，核心作用是“按规则批量处理数据”，可处理复杂逻辑（多参数、多列联动）。

### 1. 处理Series数据

```
1 # 语法格式
2 Series.apply(func, args=(), **kwargs)

1 # 使用案例
2 # 创建一个学生表
3 data = {
```

```

4     "name": ["Alice", "Bob", "Charlie", "David"],
5     "age": [20, 22, 21, 20],
6     "score": [85, 90, 78, 92],
7     "gender": ["女", "男", "男", "男"]
8 }
9 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4'], columns=
    ['name', 'age', 'score', 'gender'])
10 print(df)
11
12 # 比如需要把年龄数据乘以2
13 def mul_two(x):
14     return x * 2
15 df['age'] = df['age'].apply(mul_two)
16 # 再次打印一下看看变化
17 print(df)
18
19 # 也可以使用匿名函数
20 df['age'] = df['age'].apply(lambda x: x * 2)
21 print(df)

```

## 2. 处理DataFrame数据

```

1 # 语法格式
2 DataFrame.apply(func, axis=0, args=(), **kwargs)

```

```

1 # apply也可以处理DataFrame数据
2 # 和Series的区别在于，操作DataFrame的时候要指定轴
3
4 df = pd.DataFrame(np.arange(12).reshape(3,4), index=['s1', 's2',
    's3'], columns=['c1', 'c2', 'c3', 'c4'])
5 print(df)
6
7 # 1. 统计每个列的数据个数
8 print(df.apply(lambda x: x.count(), axis=0))           # 默认轴是0，即
    列
9
10 # 2. 统计每个列的最大值
11 print(df.apply(lambda x: x.max(), axis=0))           # 默认轴是0，即
    列
12
13 # 3. 统计每个行的数据个数
14 print(df.apply(lambda x: x.count(), axis=1))         # 轴1，是行

```

## map

map 映射函数，核心作用是“按规则逐元素转换数据”，一般用于Series

```
1 # 语法格式
2 # arg: 可传入字典, 函数(匿名函数)
3 # na_action: 缺失值处理参数, 默认为None
4 Series.map(arg, na_action=None)
```

```
1 # 使用案例
2 # 创建一个学生二维表
3 data = {
4     "name": ["Alice", "Bob", "Charlie", "David"],
5     "age": [20, 22, 21, 20],
6     "score": [85, 90, 78, 92],
7     "gender": ["女", "男", "男", "男"]
8 }
9 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4'], columns=
10 ['name', 'age', 'score', 'gender'])
11 print(df)
12
13 # 案例1: 将性别数据映射为male和female
14 # 通过字典映射
15 gender_map = {
16     "男": "male",
17     "女": "female"
18 }
19 df['gender'] = df['gender'].map(gender_map)
20 print(df)
21
22 # 案例2: 将性别数据映射为1和0
23 # 通过函数映射
24 df['gender'] = df['gender'].map(lambda x: 1 if x == "male" else 0)
25 print(df)
26
27 # 案例3: 给所有的数据乘以2
28 df2 = pd.DataFrame(np.arange(12).reshape(3,4), index=['s1', 's2',
29 's3'], columns=['c1', 'c2', 'c3', 'c4'])
30 print(df2)
31 df2 = df2.map(lambda x: x * 2)
32 print(df2)
```



# 数据分析：统计与计算

Pandas 提供丰富的统计函数，覆盖“描述性统计”“聚合计算”核心需求。

重点掌握 `df.describe()`

## 描述性统计

```
1  # 创建一个数据表
2  data = {
3      "name": ["Alice", "Bob", "Charlie", "David", "Eva", "Frank"],
4      "age": [20, 22, 21, 20, 20, 30],
5      "score": [85, 90, 88, 92, 85, 72],
6      "gender": ["女", "男", "男", '男', '女', '男']
7  }
8  df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4', 's5',
9      's6'], columns=['name', 'age', 'score', 'gender'])
10
11 # 1. 单列统计
12 print("单列统计:")
13 print(df['age'].describe())
14 print('-'*20)
15
16 # 2. 多列统计
17 print("多列统计:")
18 print(df[['age', 'score']].describe())
19 print('-'*20)
20
21 # 3. 所有列进行统计(只会对,也只能对数值列进行统计)
22 print(df.describe())
23
24 # 4. 常用统计函数
25 print("常用统计函数:")
26 print(df['age'].mean())          # 平均数
27 print(df['age'].median())        # 中位数
28 print(df['age'].min())           # 最小值
29 print(df['age'].max())           # 最大值
30 print(df['age'].std())            # 标准差
31 print(df['age'].sum())            # 求和
32
```

## 聚合计算（自己了解）

对同一列 / 多列同时应用多个统计函数，高效生成汇总结果。

```
1 # 1. 单列多函数聚合
2 score_stats = df['score'].agg([
3     "mean", "max", "min"
4 ])
5 print(score_stats)
6
7 # 2. 多列分别聚合
8 multi_col_stats = df.agg({
9     "age": ["mean", "count"],      # age列: 均值、非空计数
10    "score": ["median", "std"]     # score列: 中位数、标准差
11 })
12 print(multi_col_stats)
```

## 进阶操作: 分组

分组是“按类别拆分数据→对每组应用操作→合并结果”的过程，是数据分析的核心手段。



```
1 # 创建一个数据表
2 data = {
3     "name": ["Alice", "Bob", "Charlie", "David", "Eva", "Frank"],
4     "age": [20, 22, 21, 20, 20, 30],
5     "score": [85, 90, 88, 92, 82, 72],
6     "gender": ["女", "男", "男", "男", "女", "男"]
7 }
8 df = pd.DataFrame(data, index=['s1', 's2', 's3', 's4', 's5', 's6'], columns=['name', 'age', 'score', 'gender'])
```

```
9 print(df)
10
11 # 1. 分组
12 group_by_gender = df.groupby('gender')
13 print(group_by_gender)
14
15 # 2. 分组之后统计
16 print("分组之后统计:")
17 print(group_by_gender.describe())
18
19 # 3. 分组之后求均值
20 print("分组之后求均值:")
21 mean_score = df.groupby('gender')
22     ['score'].mean().add_prefix('mean_')
23 print(mean_score)
24
25 # 4. 分组之后进行转换
26 # 需求: 按性别分组, 求每个组的最高分 (传入已有聚合函数)
27 df["max_score"] = df.groupby("gender")["score"].transform(max)
28 print(df[["name", "gender", "score", "max_score"]])
29
30 # 需求: 按性别分组, 计算每个学生分数在组内的降序排名 (传入自定义函数)
31 # ascending=False 表示降序
32 df["score_rank"] = df.groupby("gender")["score"].transform(lambda
33     group: group.rank(ascending=False))
34 print(df[["name", "gender", "score", "score_rank"]])
```