



涛哥聊Python

...

涛哥

优质资料整理



一、入门Pandas

1.Pandas是什么

Pandas (Python Data Analysis Library的缩写) 是一个在Python编程语言中广泛使用的开源数据处理库。它提供了易于使用的数据结构和数据分析工具，特别适用于数据清洗、数据转换、数据分析和数据可视化任务。Pandas的两个主要数据结构是Series和DataFrame：

1. Series：Series是一维的带标签数据结构，类似于数组或列表。它由一组数据和与之相关的索引标签组成，使你能够快速访问和处理数据。Series可以包含各种数据类型，包括整数、浮点数、字符串等。
2. DataFrame：DataFrame是Pandas中的核心数据结构，它是一个二维表格，类似于电子表格或SQL数据库表。DataFrame由多个列组成，每一列可以包含不同类型的数据，例如数字、字符串、日期等。DataFrame也有行索引和列标签，使数据的选择和处理更加方便。

Pandas提供了丰富的功能和方法，用于数据的导入、导出、清洗、筛选、转换、统计分析和可视化。它通常与其他Python库，如NumPy、Matplotlib和Seaborn，一起使用，以构建数据科学和数据分析应用程序。

2.Pandas的主要功能

- 数据导入和导出：支持各种文件格式，如CSV、Excel、SQL数据库、JSON等。
- 数据清洗：处理缺失值、重复数据和异常值。
- 数据筛选和索引：以不同的方式选择和过滤数据。
- 数据转换：进行数据类型转换、重塑数据、合并和拆分数据。
- 统计分析：计算各种统计指标、汇总数据和执行分组操作。
- 数据可视化：创建各种图表，如折线图、柱状图、散点图等。

3.Pandas的重要性

Pandas在数据科学和数据分析领域具有重要性

1. 数据处理和清洗：Pandas提供了强大的工具，使数据清洗和预处理变得更加容易。它能够处理缺失数据、重复数据、异常值等，以准备数据用于进一步的分析。
2. 数据探索：Pandas允许分析师和数据科学家快速查看数据的摘要、统计信息以及基本属性。这有助于理解数据的特征和结构。
3. 数据转换和重塑：Pandas允许用户在数据帧中执行各种操作，如数据的透视、堆叠、合并和拆分。这对于数据的形状转换和数据重组非常有用。
4. 数据分析：Pandas提供了丰富的统计分析工具，包括分组、聚合、排序和筛选数据。这些功能对于生成报告、洞察数据分布和找到数据之间的关系非常有用。
5. 时间序列分析：Pandas专门设计用于处理时间序列数据，这在金融、经济学和其他领域中非常重要。它允许用户进行时间索引、滚动窗口统计和时间间隔计算。

6. 数据可视化：Pandas与Matplotlib、Seaborn等可视化库集成得很好，使用户能够轻松地创建各种图表和可视化效果，以更好地理解数据。
7. 数据导入和导出：Pandas支持从多种数据源导入数据，包括CSV、Excel、SQL数据库等。同时，它还可以将数据导出为各种文件格式，使数据共享和报告编制更加方便。
8. 与其他数据科学工具的整合：Pandas与其他数据科学库如NumPy、Scikit-learn和TensorFlow等协同工作，使用户能够构建端到端的数据科学和机器学习工作流。
9. 开源和活跃的社区：Pandas是一个开源项目，拥有庞大的社区支持。这意味着用户可以轻松获得文档、教程、问题解答和新功能的开发。

4. 安装Pandas和相关依赖

要安装Pandas和相关依赖，可以使用Python的包管理工具，通常是pip或conda。

以下是如何安装Pandas及其相关依赖的步骤：

(1) 使用pip安装Pandas：

使用pip，可以在命令行中运行以下命令来安装Pandas：

```
1 pip install pandas
```

这将安装最新版本的Pandas。

(2) 使用conda安装Pandas（适用于Anaconda用户）：

使用Anaconda，可以在命令行中运行以下命令来安装Pandas：

```
1 conda install pandas
```

这将使用Anaconda仓库中的Pandas版本进行安装。

(3) 安装相关依赖：

Pandas通常会自动处理其依赖关系，但如果需要额外的功能，可能需要安装其他库。例如，如果你想要在Pandas中使用数据可视化，你可以安装Matplotlib和Seaborn：

使用pip安装Matplotlib和Seaborn：

```
1 pip install matplotlib seaborn
```

在大多数情况下，安装Pandas时会自动处理大部分依赖关系，但有时可能需要根据项目的具体需求安装其他相关库。

5. 导入Pandas库

要导入Pandas库，在Python脚本或Jupyter Notebook中执行以下代码行：

```
1 import pandas as pd
```

这行代码将导入Pandas库，并将其命名为"pd"，这是一种通用的约定，方便更轻松地调用Pandas库的函数和类。

一旦导入了Pandas库，就可以开始使用它的功能进行数据处理、分析和操作。

6. 创建Pandas的Series

要创建Pandas的Series，使用 `pd.Series()` 构造函数，将一维数据结构（如列表、数组、字典等）传递给它。

以下是一些示例：

(1) 从列表创建Series：

```
1 import pandas as pd
2
3 data = [1, 2, 3, 4, 5]
4 series = pd.Series(data)
5 print(series)
```

这将创建一个包含整数的Series，并自动分配索引（默认从0开始的整数索引）。

(2) 从数组创建Series：

```
1 import pandas as pd
2 import numpy as np
3
4 data = np.array([10, 20, 30, 40, 50])
5 series = pd.Series(data)
6 print(series)
```

可以使用NumPy数组来创建Series。这样，Series将继承数组的数据类型。

(3) 自定义索引标签：

还可以为Series提供自定义索引标签，使数据更具有描述性：

```
1 import pandas as pd
2
3 data = [1, 2, 3, 4, 5]
4 index_labels = ['A', 'B', 'C', 'D', 'E']
5 series = pd.Series(data, index=index_labels)
6 print(series)
```

在这个示例中，Series的每个元素都有一个自定义的索引标签。

(4) 从字典创建Series：

你还可以从字典创建Series，字典的键将变为Series的索引，而字典的值将成为Series的数据：

```
1 import pandas as pd
2
3 data = {'Alice': 25, 'Bob': 30, 'Charlie': 35}
4 series = pd.Series(data)
5 print(series)
```

这将创建一个带有人名作为索引标签的Series，年龄作为数据。

7.创建Pandas的DataFrame

要创建Pandas的DataFrame，使用 `pd.DataFrame()` 构造函数，并向它传递数据。

DataFrame通常是二维的，可以包含多个列，并且可以提供自定义的行索引和列标签。

以下是几种创建DataFrame的示例：

(1) 从字典创建DataFrame：

```
1 import pandas as pd
2
3 data = {'Name': ['Alice', 'Bob', 'Charlie'],
4         'Age': [25, 30, 35],
5         'City': ['New York', 'Los Angeles', 'Chicago']}
6 df = pd.DataFrame(data)
7 print(df)
```

这将创建一个包含姓名、年龄和城市信息的DataFrame。列标签从字典的键中自动生成。

(2) 从列表的列表创建DataFrame：

使用包含嵌套列表的列表创建DataFrame，其中每个嵌套列表表示一行数据：

```
1 import pandas as pd
2
3 data = [['Alice', 25, 'New York'],
4         ['Bob', 30, 'Los Angeles'],
5         ['Charlie', 35, 'Chicago']]
6 df = pd.DataFrame(data, columns=['Name', 'Age', 'City'])
7 print(df)
```

在这个示例中，提供了自定义的列标签。

(3) 从CSV文件加载DataFrame：

通常，还会从外部数据源加载数据创建DataFrame，如CSV文件：

```
1 import pandas as pd
2
3 # 从CSV文件加载数据
4 df = pd.read_csv('data.csv')
5 print(df)
```

上面的代码将加载名为'data.csv'的CSV文件，并将其内容转换为DataFrame。

(4) 从其他数据源加载DataFrame：

Pandas支持从Excel文件、SQL数据库、JSON文件等多种数据源加载数据并创建DataFrame。你可以使用 `pd.read_excel()`、`pd.read_sql()`、`pd.read_json()` 等函数来实现。

二、数据加载和保存

1. 从不同数据源加载数据

(1) 从CSV文件加载数据

要从CSV文件加载数据并创建一个Pandas DataFrame，使用 `pd.read_csv()` 函数。

以下是加载CSV文件的步骤：

1. 确保已经安装了Pandas库，如果没有，请先安装它。
2. 确保CSV文件位于你的工作目录（当前工作目录），或者提供完整的文件路径。
3. 使用 `pd.read_csv()` 函数来加载CSV文件并创建DataFrame。下面是一个示例：

```
1 import pandas as pd
2
3 # 使用pd.read_csv()加载CSV文件，将数据存储在DataFrame中
4 df = pd.read_csv('your_file.csv')
```

替换 `'your_file.csv'` 为你的CSV文件的文件路径。

如果你的CSV文件使用不同的分隔符（例如分号或制表符），可以通过指定 `sep` 参数来指定分隔符。例如，对于使用分号分隔的文件：

```
1 df = pd.read_csv('your_file.csv', sep=';')
```

还可以在 `pd.read_csv()` 中传递其他参数，以根据需要进行自定义数据加载。例如，可以指定编码方式（`encoding` 参数）、处理缺失值（`na_values` 参数）、选择要读取的列（`usecols` 参数）等。

(2) 从Excel文件加载数据

要从Excel文件加载数据并创建一个Pandas DataFrame，可以使用 `pd.read_excel()` 函数。

以下是加载Excel文件的步骤：

1. 确保已经安装了Pandas库，如果没有，请先安装它。
2. 确保Excel文件位于你的工作目录（当前工作目录），或者提供完整的文件路径。
3. 使用 `pd.read_excel()` 函数来加载Excel文件并创建DataFrame。下面是一个示例：

```
1 import pandas as pd
2
3 # 使用pd.read_excel()加载Excel文件，将数据存储在DataFrame中
4 df = pd.read_excel('your_file.xlsx')
```

替换 `'your_file.xlsx'` 为你的Excel文件的文件路径。

如果你的Excel文件包含多个工作表（工作表名称不同），可以使用 `sheet_name` 参数指定要加载的工作表的名称或索引。例如，要加载第一个工作表：

```
1 df = pd.read_excel('your_file.xlsx', sheet_name=0)
```


还可以在 `pd.read_excel()` 中传递其他参数，以根据需要进行自定义数据加载，比如指定行和列的范围、解析日期和时间、处理特殊字符等。

(3) 从SQL数据库加载数据

要从SQL数据库加载数据并创建一个Pandas DataFrame，使用Pandas的 `pd.read_sql()` 函数，该函数允许执行SQL查询并将结果加载到DataFrame中。

以下是加载数据的一般步骤：

1. 确保你已经安装了Pandas库和SQL数据库连接驱动（例如，SQLite、MySQL或PostgreSQL的驱动），如果没有，请先安装它们。
2. 在Python中导入Pandas库。
3. 创建数据库连接。需要指定数据库类型（例如，SQLite、MySQL、PostgreSQL等）、主机名、用户名、密码和数据库名称等连接参数。
4. 使用 `pd.read_sql()` 函数来执行SQL查询并将结果加载到DataFrame中。下面是一个示例：

```
1 import pandas as pd
2 import sqlite3 # 导入数据库连接驱动，这里以SQLite为例
3
4 # 创建数据库连接
5 conn = sqlite3.connect('your_database.db') # 替换为你的数据库信息
6
7 # 编写SQL查询
8 query = 'SELECT * FROM your_table' # 替换为你的表名和查询
9
10 # 使用pd.read_sql()加载数据到DataFrame
11 df = pd.read_sql(query, conn)
12
13 # 关闭数据库连接
14 conn.close()
```

在上面的示例中，使用了SQLite数据库，可以根据数据库类型替换 `sqlite3.connect()` 的部分内容，以适应其他类型的数据库。

确保替换以下内容：

- `'your_database.db'`：数据库文件路径或连接字符串，根据数据库类型。
- `'your_table'`：要查询的表名。

还可以执行复杂的SQL查询，包括联接、过滤、聚合等，然后将结果加载到DataFrame中，以便进一步的数据处理和分析。

2.将数据保存到不同格式的文件

(1) 保存数据为CSV、Excel、JSON等格式

使用Pandas来将数据保存为不同格式的文件，如CSV、Excel、JSON等。以下是如何将数据保存到这些格式的示例：

将数据保存为CSV文件：

使用 `to_csv()` 方法将DataFrame保存为CSV文件。示例如下：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie'],
5         'Age': [25, 30, 35],
6         'City': ['New York', 'Los Angeles', 'Chicago']}
7 df = pd.DataFrame(data)
8
9 # 将数据保存为CSV文件
10 df.to_csv('data.csv', index=False) # index=False表示不保存行索引
```

这将保存DataFrame为名为'data.csv'的CSV文件。

将数据保存为Excel文件：

使用 `to_excel()` 方法将DataFrame保存为Excel文件。示例如下：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie'],
5         'Age': [25, 30, 35],
6         'City': ['New York', 'Los Angeles', 'Chicago']}
7 df = pd.DataFrame(data)
8
9 # 将数据保存为Excel文件
10 df.to_excel('data.xlsx', index=False) # index=False表示不保存行索引
```

这将保存DataFrame为名为'data.xlsx'的Excel文件。

将数据保存为JSON文件：

使用 `to_json()` 方法将DataFrame保存为JSON文件。示例如下：

```
1 import pandas as pd
2
```

```
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie'],
5         'Age': [25, 30, 35],
6         'City': ['New York', 'Los Angeles', 'Chicago']}
7 df = pd.DataFrame(data)
8
9 # 将数据保存为JSON文件
10 df.to_json('data.json', orient='records')
```

这将保存DataFrame为名为'data.json'的JSON文件。`orient='records'` 参数表示将数据以记录（行）的形式保存到JSON文件中。

根据需要，你可以调整这些示例，包括文件名、文件路径、保存的格式和参数。

(2) 控制保存参数

在Pandas中，通过传递不同的参数来控制如何保存DataFrame为CSV、Excel、JSON等文件格式。

以下是一些常用的控制参数示例：

保存为CSV文件时的控制参数：

- `sep`：指定列之间的分隔符，默认为逗号（','）。
- `index`：指定是否包含行索引，默认为True。如果设置为False，不会保存行索引。
- `header`：指定是否包含列标签，默认为True。如果设置为False，不会保存列标签。

```
1 # 示例：将DataFrame保存为CSV文件时的参数控制
2 df.to_csv('data.csv', sep=';', index=False, header=False)
```

保存为Excel文件时的控制参数：

- `sheet_name`：指定要保存到的工作表名称，默认为'Sheet1'。
- `index`：指定是否包含行索引，默认为True。
- `header`：指定是否包含列标签，默认为True。

```
1 # 示例：将DataFrame保存为Excel文件时的参数控制
2 df.to_excel('data.xlsx', sheet_name='Sheet1', index=False, header=True)
```

保存为JSON文件时的控制参数：

- `orient`：指定JSON文件中数据的排列方式。常见选项包括'records'（默认，按行排列）、'split'、'index'等。

```
1 # 示例：将DataFrame保存为JSON文件时的参数控制
2 df.to_json('data.json', orient='split')
```

这些参数能够根据需求自定义保存文件时的行为。可以根据不同的情况选择适当的参数来控制保存的文件格式。

三、基本数据操作

1. 数据查看和浏览

(1) `.head()` 和 `.tail()` 方法

`.head()` 和 `.tail()` 是Pandas中常用的方法，用于查看DataFrame或Series的前几行或后几行数据。它们通常用于快速浏览数据的一部分以了解数据的结构和内容。

- `.head(n)`：这个方法返回DataFrame或Series的前n行数据，默认n为5。
- `.tail(n)`：这个方法返回DataFrame或Series的后n行数据，默认n为5。

以下是这两种方法的示例用法：

DataFrame:

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 使用.head()方法查看前3行数据
9 print(df.head(3))
10
11 # 使用.tail()方法查看后2行数据
12 print(df.tail(2))
```

在上面的示例中，`.head(3)` 将返回DataFrame的前3行数据，而 `.tail(2)` 将返回DataFrame的后2行数据。

Series:

```
1 import pandas as pd
2
3 # 创建一个示例Series
```

```
4 data = [10, 20, 30, 40, 50]
5 series = pd.Series(data)
6
7 # 使用.head()方法查看前2行数据
8 print(series.head(2))
9
10 # 使用.tail()方法查看后3行数据
11 print(series.tail(3))
```

在这个示例中，`.head(2)` 将返回Series的前2行数据，而 `.tail(3)` 将返回Series的后3行数据。这两种方法对于数据探索非常有用，特别是在处理大型数据集时，可以快速查看数据的部分内容。

(2) `.info()` 方法

`.info()` 方法是Pandas中一个常用的方法，用于查看DataFrame的基本信息，包括列名、非空值数量、数据类型和内存使用等。这个方法通常在数据探索的早期用于了解数据的特点。

`.info()` 方法的输出包括以下信息：

- 列名：DataFrame中每一列的名称。
- 非空值数量：每列中非空值的数量。
- 数据类型：每列中的数据类型，如整数、浮点数、字符串等。
- 内存使用：DataFrame占用的内存量。

示例用法：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data)
8
9 # 使用.info()方法查看DataFrame的基本信息
10 df.info()
```

输出：

```
1 <class 'pandas.core.frame.DataFrame'>
2 RangeIndex: 5 entries, 0 to 4
3 Data columns (total 3 columns):
```

```
4 Name      5 non-null object
5 Age       5 non-null int64
6 City      5 non-null object
7 dtypes: int64(1), object(2)
8 memory usage: 385.0+ bytes
```

这个示例中，`.info()` 方法告诉你DataFrame有3列，共有5个条目（行）。它还显示了每列的名称、非空值数量、数据类型以及大约占用的内存量。

`.info()` 方法在数据探索中非常有用，因为它可以帮助你快速了解数据的结构和质量，特别是在处理大型数据集时。

(3) `.shape` 和 `.dtypes` 属性

`.shape` 和 `.dtypes` 是Pandas DataFrame的两个常用属性，用于获取有关DataFrame的基本信息。

1. `.shape` 属性：

- `.shape` 属性返回一个元组，其中包含DataFrame的维度信息。元组的第一个元素表示行数，第二个元素表示列数。这可以帮助你了解DataFrame的大小。

示例用法：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data)
8
9 # 获取DataFrame的维度信息
10 shape = df.shape
11 print(shape) # 输出 (5, 3)，表示DataFrame有5行和3列
```

2. `.dtypes` 属性：

- `.dtypes` 属性返回一个Series，其中包含DataFrame中每列的数据类型。这可以帮助你了解DataFrame中不同列的数据类型，如整数、浮点数、字符串等。

示例用法：

```
1 import pandas as pd
2
```

```
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data)
8
9 # 获取DataFrame中每列的数据类型
10 data_types = df.dtypes
11 print(data_types)
```

输出：

```
1 Name      object
2 Age       int64
3 City      object
4 dtype: object
```

这个示例中，`.dtypes` 属性显示了每列的数据类型。在这里，"Name"和"City"列的数据类型为对象（字符串），"Age"列的数据类型为整数。

`.shape` 和 `.dtypes` 属性对于在数据处理和分析中了解DataFrame的结构和特点非常有用。它们通常是数据探索的起点。

2.数据选择和过滤

(1) 使用中括号 `[]` 选择列

使用中括号 `[]` 来选择Pandas DataFrame中的列。

选择列的方式有多种，下面是一些示例：

选择单列：

要选择单列，只需使用列名作为DataFrame中的键：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie'],
5         'Age': [25, 30, 35],
6         'City': ['New York', 'Los Angeles', 'Chicago']}
7 df = pd.DataFrame(data)
8
9 # 选择单列
```

```
10 name_column = df['Name']
```

在这个示例中，`df['Name']` 选择了名为"Name"的列，将其存储在 `name_column` 中。

选择多列：

要选择多列，可以将列名作为一个列名列表传递给DataFrame：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie'],
5         'Age': [25, 30, 35],
6         'City': ['New York', 'Los Angeles', 'Chicago']}
7 df = pd.DataFrame(data)
8
9 # 选择多列
10 selected_columns = df[['Name', 'Age']]
```

在这个示例中，`df[['Name', 'Age']]` 选择了名为"Name"和"Age"的两列，将它们存储在 `selected_columns` 中。

条件选择列：

还可以使用条件来选择列，只选择满足特定条件的列。这需要使用布尔索引：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie'],
5         'Age': [25, 30, 35],
6         'City': ['New York', 'Los Angeles', 'Chicago']}
7 df = pd.DataFrame(data)
8
9 # 条件选择列
10 selected_columns = df[df['Age'] > 28]
```

在这个示例中，`df[df['Age'] > 28]` 选择了满足"Age"列中大于28的条件的列。

通过使用中括号 `[]` 来选择列，你可以轻松地访问DataFrame中的数据，以便进行进一步的数据处理和分析。

(2) 使用 `.loc[]` 和 `.iloc[]` 选择特定行和列

`.loc[]` 和 `.iloc[]` 是用于选择DataFrame中特定行和列的方法，它们提供了不同的索引方式：

1. `.loc[]`：基于标签的索引，用于根据行和列的标签选择数据。你可以使用列标签和行标签来访问数据。
2. `.iloc[]`：基于整数位置的索引，用于根据行和列的整数位置选择数据。这是与索引标签无关的索引方式。

以下是示例用法：

使用 `.loc[]` 选择特定行和列：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data, index=['A', 'B', 'C', 'D', 'E']) # 自定义行标签
8
9 # 使用 .loc[] 选择特定行和列
10 selected_data = df.loc['B', 'Name'] # 选择行 'B' 和列 'Name' 的数据
```

在上面的示例中，`.loc['B', 'Name']` 选择了行标签为'B'和列标签为'Name'的数据。

使用 `.iloc[]` 选择特定行和列：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data)
8
9 # 使用 .iloc[] 选择特定行和列
10 selected_data = df.iloc[1, 0] # 选择第2行和第1列的数据（从0开始计数）
```

在这个示例中，`.iloc[1, 0]` 选择了第2行和第1列的数据。

`.loc[]` 和 `.iloc[]` 还支持切片操作，以选择一组行或列。例如，可以使用 `.loc[]` 来选择多行和多列，或使用 `.iloc[]` 选择多行和多列的整数位置。这两种方法在数据处理中非常有用，特别是当需要根据标签或整数位置来选择数据时。

(3) 使用条件过滤数据

在Pandas中，可以使用条件过滤数据，即根据特定条件选择DataFrame中的行。这可以通过布尔索引来实现。

以下是一些示例，演示如何使用条件来过滤数据：

示例1：基于列的条件过滤

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 使用条件过滤数据
9 filtered_data = df[df['Age'] > 28] # 选择年龄大于28的行
```

在这个示例中，`df[df['Age'] > 28]` 将选择DataFrame中年龄大于28的行。

示例2：多条件过滤

可以使用多个条件来过滤数据。在这种情况下，需要使用括号将每个条件括起来，并使用逻辑运算符如 `&`（与）和 `|`（或）来组合条件。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data)
8
9 # 使用多条件过滤数据
10 filtered_data = df[(df['Age'] > 25) & (df['City'] == 'New York')]
```

在这个示例中，`df[(df['Age'] > 25) & (df['City'] == 'New York')]` 选择了DataFrame中年龄大于25且城市为"New York"的行。

示例3：使用 `.query()` 方法

Pandas还提供了 `.query()` 方法，允许使用表达式来过滤数据。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data)
8
9 # 使用 .query() 方法来过滤数据
10 filtered_data = df.query('Age > 28 and City == "New York"')
```

这将选择DataFrame中年龄大于28且城市为"New York"的行。

3.数据排序

(1) 使用 `.sort_values()` 方法进行排序

在Pandas中，使用 `.sort_values()` 方法对DataFrame中的数据进行排序。

`.sort_values()` 方法可以按照一个或多个列的值对DataFrame进行排序，可以按升序或降序排序。

以下是示例用法：

示例1：按单列进行排序

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 按照 'Age' 列升序排序
9 sorted_df = df.sort_values(by='Age')
```

在这个示例中，`.sort_values(by='Age')` 对DataFrame按照'Age'列的值进行升序排序。

示例2：按多列进行排序

还可以按多个列的值进行排序。在这种情况下，需要传递一个列名列表作为 `by` 参数。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24],
6         'City': ['New York', 'Los Angeles', 'Chicago', 'San Francisco', 'Boston']}
7 df = pd.DataFrame(data)
8
9 # 按照 'Age' 和 'City' 列升序排序
10 sorted_df = df.sort_values(by=['Age', 'City'])
```

在这个示例中，`.sort_values(by=['Age', 'City'])` 对DataFrame按照'Age'列的值进行升序排序，如果有相同的'Age'值，则按照'City'列的值进行排序。

示例3：降序排序

按降序排序，可以使用 `ascending` 参数。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 按照 'Age' 列降序排序
9 sorted_df = df.sort_values(by='Age', ascending=False)
```

在这个示例中，`.sort_values(by='Age', ascending=False)` 对DataFrame按照'Age'列的值进行降序排序。

(2) 使用 `.sort_index()` 按索引排序

在Pandas中，使用 `.sort_index()` 方法按行或列的索引进行排序。这对于重新排列DataFrame的行或列顺序非常有用。

`.sort_index()` 方法可以指定升序或降序排序以及按照行索引或列索引进行排序。

以下是示例用法：

示例1：按行索引排序

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
```

```
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24]}
6 df = pd.DataFrame(data, index=['E', 'B', 'C', 'A', 'D']) # 自定义行索引
7
8 # 按行索引升序排序
9 sorted_df = df.sort_index()
```

在这个示例中，`.sort_index()` 按行索引进行升序排序。

示例2：按列索引排序

按列索引排序，可以通过指定 `axis` 参数为1来实现。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 按列索引升序排序
9 sorted_df = df.sort_index(axis=1)
```

在这个示例中，`.sort_index(axis=1)` 按列索引进行升序排序。

示例3：降序排序

还可以通过 `ascending` 参数来指定升序或降序排序。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eve'],
5         'Age': [25, 30, 35, 28, 24]}
6 df = pd.DataFrame(data, index=['E', 'B', 'C', 'A', 'D'])
7
8 # 按行索引降序排序
9 sorted_df = df.sort_index(ascending=False)
```

在这个示例中，`.sort_index(ascending=False)` 按行索引进行降序排序。

`.sort_index()` 方法对于重新排列DataFrame的行或列非常有用，特别是在需要按照索引进行特定顺序的情况下。

4.描述性统计

(1) 使用 `.describe()` 方法生成基本统计信息

`.describe()` 方法是Pandas中用于生成基本统计信息的有用工具。这个方法用于生成关于数值列的统计摘要，包括均值、标准差、最小值、25%分位数、中位数、75%分位数和最大值。它帮助快速了解数据的分布和一些基本统计特征。

以下是 `.describe()` 方法的示例用法：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Age': [25, 30, 35, 28, 24, 32, 38, 41, 21, 28],
5         'Salary': [50000, 60000, 75000, 60000, 55000, 70000, 85000, 90000, 45000, 65000]}
6 df = pd.DataFrame(data)
7
8 # 使用 .describe() 方法生成基本统计信息
9 statistics = df.describe()
```

运行上述示例后，`statistics` 将包含有关"Age"和"Salary"列的基本统计信息，如均值、标准差、最小值、25%分位数、中位数、75%分位数和最大值。

示例输出：

	Age	Salary
count	10.000000	10.000000
mean	31.100000	64750.000000
std	6.210708	15195.213322
min	21.000000	45000.000000
25%	26.750000	55000.000000
50%	28.000000	62500.000000
75%	35.750000	71250.000000
max	41.000000	90000.000000

这些统计信息提供了有关数据的汇总，快速了解数据的分布和范围。`.describe()` 方法对于数据探索和初步分析非常有用。注意，`.describe()` 默认仅应用于数值列，因此非数值列将被忽略。

(2) 使用各种统计函数如 `.mean()`、`.median()`、`.std()` 等

Pandas提供了多种用于计算统计指标的方法，如 `.mean()`、`.median()`、`.std()` 等。这些方法可用于计算数值数据列的各种统计信息。

以下是示例用法：

示例用法 - 计算均值（平均值）：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Age': [25, 30, 35, 28, 24]}
5 df = pd.DataFrame(data)
6
7 # 计算 'Age' 列的均值
8 mean_age = df['Age'].mean()
```

在这个示例中，`df['Age'].mean()` 计算了'Age'列的均值。

示例用法 - 计算中位数：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Age': [25, 30, 35, 28, 24]}
5 df = pd.DataFrame(data)
6
7 # 计算 'Age' 列的中位数
8 median_age = df['Age'].median()
```

在这个示例中，`df['Age'].median()` 计算了'Age'列的中位数。

示例用法 - 计算标准差：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Age': [25, 30, 35, 28, 24]}
5 df = pd.DataFrame(data)
6
7 # 计算 'Age' 列的标准差
8 std_deviation_age = df['Age'].std()
```

在这个示例中，`df['Age'].std()` 计算了'Age'列的标准差。

还可以使用其他统计方法如 `.max()`（最大值）、`.min()`（最小值）、`.sum()`（总和）、`.var()`（方差）等来计算不同的统计指标。这些方法对于在数据分析和探索中计算和了解数

据的特性非常有用。可以根据具体的需求选择适当的统计函数。

四、数据清洗和预处理

1.处理缺失数据

(1) `.isna()` 和 `.fillna()` 方法处理缺失数据

在Pandas中,使用 `.isna()` 和 `.fillna()` 方法来处理缺失数据 (NaN)。缺失数据是数据集中的空值或未定义值,通常由于某种原因缺少或未记录。

以下是这两种方法的示例用法:

示例用法 - `.isna()` 检测缺失数据:

`.isna()` 方法用于检测DataFrame中的缺失数据,它返回一个布尔值的DataFrame,其中 `True` 表示缺失数据, `False` 表示非缺失数据。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame包含缺失数据
4 data = {'Name': ['Alice', 'Bob', 'Charlie', None, 'Eve'],
5         'Age': [25, 30, None, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 使用 .isna() 检测缺失数据
9 missing_data = df.isna()
```

在这个示例中, `missing_data` DataFrame的值将显示哪些单元格包含缺失数据,例如,第三行的"Name"和第三行的"Age"列包含缺失数据。

示例用法 - `**.fillna()**` 填充缺失数据:

`.fillna()` 方法用于填充缺失数据,可以指定要用来填充缺失值的特定值。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame包含缺失数据
4 data = {'Name': ['Alice', 'Bob', 'Charlie', None, 'Eve'],
5         'Age': [25, 30, None, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 使用 .fillna() 填充缺失数据
9 df_filled = df.fillna(0) # 用0填充缺失数据
```

在这个示例中，`.fillna(0)` 将DataFrame中的缺失数据填充为0。

也可以使用其他值来填充缺失数据，如均值、中位数、前一行的值等，具体取决于你的需求。例如，要使用均值填充缺失数据：

```
1 # 使用均值填充缺失数据
2 mean_age = df['Age'].mean() # 计算 'Age' 列的均值
3 df['Age'] = df['Age'].fillna(mean_age)
```

这将使用'Age'列的均值来填充缺失数据。

缺失数据处理对于数据清洗和分析非常重要，可以根据具体情况选择如何检测和填充缺失数据。

(2) `.dropna()` 方法删除缺失数据

`.dropna()` 方法是Pandas中用于删除缺失数据的方法。可以从DataFrame中删除包含缺失值（NaN）的行或列，以便在数据清洗和分析中去除不完整的数据。

以下是 `.dropna()` 方法的示例用法：

示例用法 - 删除包含缺失值的行：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame包含缺失数据
4 data = {'Name': ['Alice', 'Bob', None, 'David', 'Eve'],
5         'Age': [25, 30, None, 28, None]}
6 df = pd.DataFrame(data)
7
8 # 使用 .dropna() 删除包含缺失值的行
9 df_cleaned = df.dropna()
```

在这个示例中，`df.dropna()` 将删除包含缺失值的行，因此只有第1行和第4行会保留。

示例用法 - 删除包含缺失值的列：

还可以使用 `.dropna()` 方法删除包含缺失值的列，通过指定 `axis` 参数为1来实现。

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame包含缺失数据
4 data = {'Name': ['Alice', 'Bob', None, 'David', 'Eve'],
5         'Age': [25, 30, None, 28, None]}
6 df = pd.DataFrame(data)
7
```

```
8 # 使用 .dropna(axis=1) 删除包含缺失值的列
9 df_cleaned = df.dropna(axis=1)
```

在这个示例中，`df.dropna(axis=1)` 将删除包含缺失值的列，因此只有'Name'列会保留。

示例用法 - 自定义删除条件：

还可以使用 `.dropna()` 方法的其他参数来根据自定义条件删除行或列。例如，通过指定 `thresh` 参数来设置删除前需要有多少非缺失值的阈值。

```
1 # 删除至少需要有2个非缺失值的行
2 df_cleaned = df.dropna(thresh=2)
```

`.dropna()` 方法对于数据清洗中的缺失值处理非常有用，可以根据具体需求选择删除缺失数据的方式。

2.数据去重

在Pandas中，使用 `.duplicated()` 方法来识别重复的行，以及 `.drop_duplicates()` 方法来删除这些重复行。这对于数据清洗和去除重复数据非常有用。

以下是这两种方法的示例用法：

示例用法 - `.duplicated()` 方法识别重复行：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame包含重复数据
4 data = {'Name': ['Alice', 'Bob', 'Alice', 'David', 'Eve'],
5         'Age': [25, 30, 25, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 使用 .duplicated() 方法识别重复行
9 duplicates = df.duplicated()
```

在这个示例中，`duplicates` 将包含一个布尔值的Series，其中 `True` 表示该行是重复的。

示例用法 - `.drop_duplicates()` 方法删除重复行：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame包含重复数据
4 data = {'Name': ['Alice', 'Bob', 'Alice', 'David', 'Eve'],
```

```
5     'Age': [25, 30, 25, 28, 24]}
6 df = pd.DataFrame(data)
7
8 # 使用 .drop_duplicates() 方法删除重复行
9 df_cleaned = df.drop_duplicates()
```

在这个示例中，`df.drop_duplicates()` 将删除重复的行，保留第一次出现的行。因此，重复的行'Name'为'Alice'和'Age'为25的行会被删除。

还可以使用 `.drop_duplicates()` 方法的其他参数来自定义删除重复行的方式，如指定特定列来识别重复，或保留最后一次出现的行等。这两种方法对于数据去重和数据清洗非常有用。

3.数据类型转换

(1) `.astype()` 方法改变列的数据类型

`.astype()` 方法用于改变Pandas DataFrame或Series中的数据类型。这对于数据类型转换和数据处理非常有用，例如，将字符串列转换为数值列或将数值列转换为日期时间列。

以下是 `.astype()` 方法的示例用法：

示例用法 - 将列转换为整数数据类型：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Age': ['25', '30', '35', '28', '24']}
5 df = pd.DataFrame(data)
6
7 # 使用 .astype() 将 'Age' 列转换为整数数据类型
8 df['Age'] = df['Age'].astype(int)
```

在这个示例中，`.astype(int)` 将'Age'列的数据类型从字符串转换为整数。

示例用法 - 将列转换为日期时间数据类型：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Date': ['2023-01-15', '2023-02-20', '2023-03-25']}
5 df = pd.DataFrame(data)
6
7 # 使用 .astype() 将 'Date' 列转换为日期时间数据类型
8 df['Date'] = pd.to_datetime(df['Date'])
```

在这个示例中，`pd.to_datetime()` 将'Date'列的数据类型从字符串转换为日期时间。

示例用法 - 将列转换为分类数据类型：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Category': ['A', 'B', 'A', 'C', 'B']}
5 df = pd.DataFrame(data)
6
7 # 使用 .astype() 将 'Category' 列转换为分类数据类型
8 df['Category'] = df['Category'].astype('category')
```

在这个示例中，`.astype('category')` 将'Category'列的数据类型转换为分类数据类型，这在节省内存和加速某些操作时非常有用。

`.astype()` 方法可以根据需要更改数据类型，确保数据适用于特定的分析或可视化任务。注意，数据类型的转换应谨慎进行，因为可能会导致数据损失或不匹配的情况。

(2) 处理日期和时间数据

在Pandas中，处理日期和时间数据非常方便，Pandas提供了丰富的功能和方法来处理日期和时间数据。

以下是一些常见的用例和示例：

示例1：创建日期时间数据

使用 `pd.to_datetime()` 方法将字符串转换为日期时间格式：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame包含日期时间数据
4 data = {'Date': ['2023-01-15', '2023-02-20', '2023-03-25']}
5 df = pd.DataFrame(data)
6
7 # 将 'Date' 列转换为日期时间数据类型
8 df['Date'] = pd.to_datetime(df['Date'])
```

在这个示例中，将'Date'列的字符串转换为日期时间数据类型。

示例2：提取日期和时间组件

一旦你有日期时间列，可以轻松提取日期和时间的不同组件，如年、月、日、小时、分钟等。使用 `.dt` 属性来实现

```
1 # 提取年份
2 df['Year'] = df['Date'].dt.year
3
4 # 提取月份
5 df['Month'] = df['Date'].dt.month
6
7 # 提取日
8 df['Day'] = df['Date'].dt.day
9
10 # 提取小时
11 df['Hour'] = df['Date'].dt.hour
```

这些新列包含了从日期时间列中提取的各种组件。

示例3：计算日期差值

计算日期时间列之间的差值，例如，计算两个日期之间的天数差：

```
1 # 创建一个新列来计算日期差值
2 df['Date1'] = pd.to_datetime(['2023-01-15', '2023-02-20', '2023-03-25'])
3 df['DateDiff'] = (df['Date1'] - df['Date']).dt.days
```

在这个示例中，创建了一个新列"DateDiff"，其中包含两个日期之间的天数差。

示例4：日期时间范围

Pandas可以生成日期时间范围，例如生成一个包含一段时间内每日日期的日期范围：

```
1 # 生成日期范围
2 date_range = pd.date_range(start='2023-01-01', end='2023-01-10', freq='D')
```

这将生成一个包含从'2023-01-01'到'2023-01-10'的每日日期的日期范围。

这些示例演示了Pandas用于处理日期和时间数据的基本功能。Pandas还提供了许多其他功能，如日期时间索引、重采样、移动窗口统计等，以满足各种日期和时间数据处理需求。

4.重命名列和索引

在Pandas中，使用 `.rename()` 方法来重命名DataFrame的列和索引。这对于更改列名或索引名以匹配分析需求或使数据更具可读性非常有用。

以下是 `.rename()` 方法的示例用法：

示例用法 - 重命名列：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'OldName1': [1, 2, 3],
5         'OldName2': [4, 5, 6]}
6 df = pd.DataFrame(data)
7
8 # 使用 .rename() 方法重命名列
9 df = df.rename(columns={'OldName1': 'NewName1', 'OldName2': 'NewName2'})
```

在这个示例中，`.rename(columns={'OldName1': 'NewName1', 'OldName2': 'NewName2'})` 将'OldName1'列重命名为'NewName1'，'OldName2'列重命名为'NewName2'。

示例用法 - 重命名索引：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'A': [1, 2, 3],
5         'B': [4, 5, 6]}
6 df = pd.DataFrame(data, index=['X', 'Y', 'Z'])
7
8 # 使用 .rename() 方法重命名索引
9 df = df.rename(index={'X': 'Row1', 'Y': 'Row2', 'Z': 'Row3'})
```

在这个示例中，`.rename(index={'X': 'Row1', 'Y': 'Row2', 'Z': 'Row3'})` 将索引'X'重命名为'Row1'，'Y'重命名为'Row2'，'Z'重命名为'Row3'。

示例用法 - 同时重命名列和索引：

也可以同时重命名列和索引：

```
1 # 同时重命名列和索引
2 df = df.rename(columns={'OldName1': 'NewName1', 'OldName2': 'NewName2'},
3                 index={'X': 'Row1', 'Y': 'Row2', 'Z': 'Row3'})
```

这将同时重命名列和索引。

`.rename()` 方法可以在不修改数据的情况下更改列和索引的名称，以适应分析需求或使数据更具可读性。

五、数据分析和可视化

1. 基本统计分析

(1) `.groupby()` 进行分组统计

`.groupby()` 是Pandas中用于分组数据并进行统计分析的强大工具。可以根据一个或多个列的值将数据分成多个组，并在这些组上执行聚合操作。

以下是 `.groupby()` 方法的示例用法：

示例用法 - 基本的分组和聚合：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Category': ['A', 'B', 'A', 'B', 'A'],
5         'Value': [10, 20, 30, 40, 50]}
6 df = pd.DataFrame(data)
7
8 # 使用 .groupby() 分组并计算每个组的均值
9 grouped = df.groupby('Category')
10 mean_values = grouped['Value'].mean()
```

在这个示例中，首先使用 `.groupby('Category')` 将数据分为两个组，'A'和'B'，然后计算每个组中'Value'列的均值。

示例用法 - 多列分组和多个聚合操作：

还可以使用多个列进行分组，以及执行多个聚合操作：

```
1 # 使用多列分组并执行多个聚合操作
2 grouped = df.groupby(['Category', 'Type'])
3 agg_results = grouped['Value'].agg(['sum', 'mean', 'max'])
```

在这个示例中，首先使用 `.groupby(['Category', 'Type'])` 将数据根据'Category'和'Type'列进行多级分组，然后执行不同的聚合操作，包括总和、均值和最大值。

示例用法 - 自定义聚合函数：

还可以使用自定义的聚合函数来进行更复杂的分组统计：

```

1 # 使用自定义聚合函数
2 def custom_agg_function(x):
3     return x.max() - x.min()
4
5 grouped = df.groupby('Category')
6 custom_agg = grouped['Value'].agg(custom_agg_function)

```

在这个示例中，定义了一个自定义的聚合函数 `custom_agg_function`，用于计算每个分组中'Value'列的范围（最大值减去最小值）。

`.groupby()` 方法非常有用，对数据进行分组分析，了解不同组的统计信息，进行聚合操作，以及执行自定义的聚合函数。这对于数据探索和分析非常重要。

(2) `.agg()` 计算多个统计量

`.agg()` 方法用于计算多个统计量（如均值、总和、中位数、标准差等）并应用这些统计量到分组后的数据。这是Pandas中非常强大和灵活的工具，允许一次性计算多个统计量。

以下是 `.agg()` 方法的示例用法：

示例用法 - 计算多个统计量：

```

1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Category': ['A', 'B', 'A', 'B', 'A'],
5         'Value': [10, 20, 30, 40, 50]}
6 df = pd.DataFrame(data)
7
8 # 使用 .groupby() 分组并计算多个统计量
9 grouped = df.groupby('Category')
10 summary_stats = grouped['Value'].agg(['mean', 'sum', 'median', 'std'])

```

在这个示例中，`.agg(['mean', 'sum', 'median', 'std'])` 计算了每个分组的'Value'列的均值、总和、中位数和标准差，生成了一个包含这些统计量的DataFrame。

示例用法 - 自定义统计函数：

还可以使用自定义的统计函数来计算多个统计量。例如，定义一个函数来计算范围（最大值减去最小值）：

```

1 # 使用自定义统计函数
2 def custom_stats(x):
3     return x.max() - x.min()

```

```
4
5 grouped = df.groupby('Category')
6 custom_summary = grouped['Value'].agg(['mean', 'sum', custom_stats])
```

在这个示例中，定义一个自定义的统计函数 `custom_stats`，用于计算每个分组的'Value'列的均值、总和和范围。

`.agg()` 方法允许一次性计算多个统计量，包括标准统计量和自定义统计函数，以满足数据分析的需求。这对于在数据分组上执行多个操作非常有用，减少了代码的复杂性。

2.数据可视化

(1) 使用Matplotlib和Seaborn库创建各种图表

Matplotlib和Seaborn是Python中常用的数据可视化库，可以创建各种类型的图表来展示和分析数据。

以下是使用Matplotlib和Seaborn创建常见图表的示例：

示例1：创建折线图（Line Plot）

使用Matplotlib创建一个简单的折线图：

```
1 import matplotlib.pyplot as plt
2
3 # 示例数据
4 x = [1, 2, 3, 4, 5]
5 y = [10, 12, 5, 8, 6]
6
7 # 创建折线图
8 plt.plot(x, y)
9 plt.title('Line Plot Example')
10 plt.xlabel('X-Axis')
11 plt.ylabel('Y-Axis')
12 plt.show()
```

这将创建一个简单的折线图，显示x和y之间的关系。

示例2：创建散点图（Scatter Plot）

使用Matplotlib创建一个散点图：

```
1 import matplotlib.pyplot as plt
2
3 # 示例数据
```

```
4 x = [1, 2, 3, 4, 5]
5 y = [10, 12, 5, 8, 6]
6
7 # 创建散点图
8 plt.scatter(x, y)
9 plt.title('Scatter Plot Example')
10 plt.xlabel('X-Axis')
11 plt.ylabel('Y-Axis')
12 plt.show()
```

这将创建一个散点图，显示x和y之间的关系。

示例3：创建直方图（Histogram）

使用Matplotlib创建一个直方图：

```
1 import matplotlib.pyplot as plt
2
3 # 示例数据
4 data = [3, 3, 1, 5, 6, 7, 8, 7, 4, 3, 2, 1, 4, 5, 6]
5
6 # 创建直方图
7 plt.hist(data, bins=5, edgecolor='black')
8 plt.title('Histogram Example')
9 plt.xlabel('Value')
10 plt.ylabel('Frequency')
11 plt.show()
```

这将创建一个直方图，显示数据的分布情况。

示例4：创建箱线图（Box Plot）

使用Seaborn创建一个箱线图：

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # 示例数据
5 data = sns.load_dataset("tips")
6
7 # 创建箱线图
8 sns.boxplot(x="day", y="total_bill", data=data)
9 plt.title('Box Plot Example')
10 plt.xlabel('Day')
11 plt.ylabel('Total Bill')
```

```
12 plt.show()
```

这将创建一个箱线图，显示不同天的总账单分布情况。

示例5：创建热力图（Heatmap）

使用Seaborn创建一个热力图：

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # 示例数据
5 data = sns.load_dataset("flights")
6
7 # 重塑数据以便创建热力图
8 pivot_data = data.pivot_table(index='month', columns='year', values='passengers')
9
10 # 创建热力图
11 sns.heatmap(pivot_data, cmap="YlGnBu", annot=True, fmt="d")
12 plt.title('Heatmap Example')
13 plt.show()
```

这将创建一个热力图，显示不同月份和年份的乘客数量。

（2）绘制柱状图、折线图、散点图等

以下是使用Matplotlib和Seaborn库绘制柱状图、折线图和散点图的示例：

示例1：创建柱状图（Bar Chart）

使用Matplotlib创建一个简单的柱状图：

```
1 import matplotlib.pyplot as plt
2
3 # 示例数据
4 categories = ['Category A', 'Category B', 'Category C', 'Category D']
5 values = [25, 40, 30, 35]
6
7 # 创建柱状图
8 plt.bar(categories, values)
9 plt.title('Bar Chart Example')
10 plt.xlabel('Categories')
11 plt.ylabel('Values')
12 plt.show()
```

这将创建一个柱状图，显示不同类别的值。

示例2：创建折线图（Line Plot）

使用Matplotlib创建一个简单的折线图：

```
1 import matplotlib.pyplot as plt
2
3 # 示例数据
4 x = [1, 2, 3, 4, 5]
5 y = [10, 12, 5, 8, 6]
6
7 # 创建折线图
8 plt.plot(x, y)
9 plt.title('Line Plot Example')
10 plt.xlabel('X-Axis')
11 plt.ylabel('Y-Axis')
12 plt.show()
```

这将创建一个折线图，显示x和y之间的关系。

示例3：创建散点图（Scatter Plot）

使用Matplotlib创建一个散点图：

```
1 import matplotlib.pyplot as plt
2
3 # 示例数据
4 x = [1, 2, 3, 4, 5]
5 y = [10, 12, 5, 8, 6]
6
7 # 创建散点图
8 plt.scatter(x, y)
9 plt.title('Scatter Plot Example')
10 plt.xlabel('X-Axis')
11 plt.ylabel('Y-Axis')
12 plt.show()
```

这将创建一个散点图，显示x和y之间的关系。

以上示例使用Matplotlib绘制了柱状图、折线图和散点图。你可以根据需要自定义图表的样式、颜色、标签等来使其更具吸引力。

六、高级主题

1.时间序列数据分析

(1) 处理时间序列数据

在Pandas中，处理时间序列数据非常常见，因为它提供了丰富的工具和函数来处理日期和时间数据。

以下是处理时间序列数据的一些常见操作和示例：

示例1：创建时间序列

可以使用 `pd.to_datetime()` 来创建时间序列：

```
1 import pandas as pd
2
3 # 创建一个时间序列
4 date_series = pd.to_datetime(['2023-01-15', '2023-02-20', '2023-03-25'])
```

这将创建一个包含日期时间数据的时间序列。

示例2：设置时间序列为索引

时间序列通常用作DataFrame的索引：

```
1 # 创建一个示例DataFrame并将时间序列设置为索引
2 data = {'Value': [10, 20, 30]}
3 df = pd.DataFrame(data, index=date_series)
```

这将创建一个带有日期时间索引的DataFrame。

示例3：重采样和聚合

可以使用 `.resample()` 来重采样时间序列，并进行聚合操作，例如，将每日数据聚合为每月数据：

```
1 # 重采样为每月数据并计算均值
2 monthly_data = df['Value'].resample('M').mean()
```

在这个示例中，`.resample('M').mean()` 将每日数据聚合为每月均值。

示例4：滚动窗口统计

可以使用滚动窗口来计算移动平均或其他滚动统计量：

```
1 # 计算7天滚动窗口的移动平均
2 rolling_mean = df['Value'].rolling(window=7).mean()
```


这将创建一个包含7天滚动窗口移动平均的Series。

示例5：时间序列索引的选择和切片

可以使用时间序列索引来选择和切片时间序列数据：

```
1 # 选择特定日期范围的数据
2 selected_data = df['2023-01-15':'2023-02-28']
```

这将选择'2023-01-15'到'2023-02-28'之间的数据。

(2) 滚动窗口统计

在数据分析中，滚动窗口统计是一种常见的技术，用于计算时间序列或其他有序数据中的滚动汇总统计信息。Pandas提供了用于执行滚动窗口统计的函数，可以计算移动平均、滚动总和、滚动标准差等。

以下是滚动窗口统计的示例用法：

示例1：计算滚动平均（Moving Average）

使用Pandas计算滚动窗口的移动平均：

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # 示例数据
5 data = {'Value': [10, 15, 20, 25, 30, 35, 40, 45, 50]}
6 df = pd.DataFrame(data)
7
8 # 计算5个数据点的滚动平均
9 rolling_mean = df['Value'].rolling(window=5).mean()
10
11 # 绘制原始数据和滚动平均
12 plt.plot(df['Value'], label='Original Data')
13 plt.plot(rolling_mean, label='Rolling Mean (window=5)')
14 plt.legend()
15 plt.title('Moving Average Example')
16 plt.xlabel('Time')
17 plt.ylabel('Value')
18 plt.show()
```

这将计算5个数据点的滚动平均并将其绘制在原始数据上。

示例2：计算滚动总和（Moving Sum）

使用Pandas计算滚动窗口的移动总和：

```
1 # 计算3个数据点的滚动总和
2 rolling_sum = df['Value'].rolling(window=3).sum()
3
4 # 绘制原始数据和滚动总和
5 plt.plot(df['Value'], label='Original Data')
6 plt.plot(rolling_sum, label='Rolling Sum (window=3)')
7 plt.legend()
8 plt.title('Moving Sum Example')
9 plt.xlabel('Time')
10 plt.ylabel('Value')
11 plt.show()
```

这将计算3个数据点的滚动总和并将其绘制在原始数据上。

可以根据需要调整滚动窗口的大小以适应不同的数据和分析需求。滚动窗口统计有助于平滑数据并识别趋势，它在时间序列分析、信号处理和其他领域中都有广泛的应用。

2. 多层索引

在Pandas中，可以创建和操作多层索引（也称为层次化索引）的DataFrame，以便更灵活地表示和处理多维数据。多层索引可以在一个DataFrame中表示具有多个维度的数据，类似于多维数组。

以下是创建和操作多层索引的示例：

示例1：创建具有多层索引的DataFrame

```
1 import pandas as pd
2
3 # 创建一个带有多层索引的DataFrame
4 data = {
5     'A': [1, 2, 3, 4, 5, 6],
6     'B': [7, 8, 9, 10, 11, 12]
7 }
8 index = pd.MultiIndex.from_tuples([('Group1', 'A'), ('Group1', 'B'), ('Group2',
9 df = pd.DataFrame(data, index=index)
```

在这个示例中，创建一个具有两层索引的DataFrame，第一层是'Group'，第二层是'Category'。

示例2：操作多层索引的DataFrame

可以使用多层索引来选择、切片和执行操作：

```

1 # 选择第一层索引为 'Group1' 的数据
2 group1_data = df.loc['Group1']
3
4 # 选择第一层索引为 'Group1' 和第二层索引为 'A' 的数据
5 group1_a_data = df.loc[['Group1', 'A']]
6
7 # 切片第一层索引
8 group1_to_group2 = df.loc['Group1':'Group2']
9
10 # 使用多层索引进行分组统计
11 group_means = df.groupby('Group').mean()

```

这些示例演示了如何创建和操作具有多层索引的DataFrame，包括选择特定索引级别的数据、使用多层索引进行切片和分组统计等。

多层索引的DataFrame对于表示复杂的多维数据非常有用，例如具有多个组和多个类别的数据。可以根据需要添加、删除或调整索引级别，以满足你的数据分析需求。

3.高级数据过滤和选择

在Pandas中，使用复杂的条件和查询来筛选、过滤和操作DataFrame中的数据。这可以通过多种方法实现，包括使用布尔索引、使用 `.query()` 方法和使用多个条件操作符。

以下是一些示例：

示例1：使用布尔索引进行条件筛选

```

1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'A': [1, 2, 3, 4, 5],
5         'B': [10, 20, 30, 40, 50]}
6 df = pd.DataFrame(data)
7
8 # 使用布尔索引筛选满足条件的行
9 filtered_df = df[(df['A'] > 2) & (df['B'] < 40)]

```

在这个示例中，我们使用布尔索引筛选了'A'列大于2且'B'列小于40的行。

示例2：使用 `.query()` 方法进行条件查询

```

1 # 使用 `.query()` 方法进行条件查询
2 query_result = df.query('A > 2 and B < 40')

```

`.query()` 方法允许你以更简洁的方式执行条件查询。

示例3：使用多个条件操作符

使用多个条件操作符，如 `&`（与）、`|`（或）和 `~`（非），来构建更复杂的条件：

```
1 # 使用多个条件操作符构建复杂条件
2 complex_condition = (df['A'] > 2) & (df['B'] < 40) | (df['A'] == 5)
3 result = df[complex_condition]
```

这将筛选出'A'列大于2且'B'列小于40，或者'A'列等于5的行。

示例4：使用 `** .loc[] **` 和 `** .iloc[] **` 选择特定行和列

使用 `.loc[]` 和 `.iloc[]` 选择特定行和列，进一步细化数据选择：

```
1 # 使用 .loc[] 选择特定行和列
2 selected_data = df.loc[(df['A'] > 2) & (df['B'] < 40), ['A']]
3
4 # 使用 .iloc[] 选择特定行和列（通过整数位置）
5 selected_data = df.iloc[2:4, 0:1]
```

4. 自定义函数应用

在Pandas中，使用 `.apply()` 方法来应用自定义函数到DataFrame的行或列，以实现更灵活的数据处理和转换。

以下是 `.apply()` 方法的示例用法：

示例1：应用自定义函数到列

假设有一个DataFrame，包含成绩数据，你可以使用 `.apply()` 方法来计算每个学生的总分：

```
1 import pandas as pd
2
3 # 创建一个示例DataFrame
4 data = {'Name': ['Alice', 'Bob', 'Charlie'],
5         'Math': [85, 90, 78],
6         'Science': [92, 88, 86]}
7 df = pd.DataFrame(data)
8
9 # 定义一个函数来计算总分
10 def calculate_total(row):
11     return row['Math'] + row['Science']
12
```

```
13 # 使用 .apply() 方法应用函数到列
14 df['Total'] = df.apply(calculate_total, axis=1)
```

在这个示例中，定义一个名为 `calculate_total` 的函数，然后使用 `.apply()` 方法将这个函数应用到DataFrame的每一行，并将结果存储在一个新的'Total'列中。

示例2：应用自定义函数到列并返回Series

使用 `.apply()` 方法将自定义函数应用到列，并返回一个Series：

```
1 # 定义一个函数来计算成绩等级
2 def calculate_grade(score):
3     if score >= 90:
4         return 'A'
5     elif score >= 80:
6         return 'B'
7     elif score >= 70:
8         return 'C'
9     else:
10        return 'D'
11
12 # 使用 .apply() 方法应用函数到列并返回Series
13 df['Math_Grade'] = df['Math'].apply(calculate_grade)
14 df['Science_Grade'] = df['Science'].apply(calculate_grade)
```

在这个示例中，我们定义了一个名为 `calculate_grade` 的函数，将该函数应用到'Math'和'Science'列，并将结果存储在新的'Math_Grade'和'Science_Grade'列中。

`.apply()` 方法可以在DataFrame中应用自定义函数，以根据需求进行数据处理和转换。可以在函数中实现复杂的逻辑来实现定制化的数据操作。

七、实际项目示例

通过一个实际项目案例演示如何应用Pandas进行数据处理和分析当进行实际项目案例时，Pandas可以用来处理和分析数据。

项目案例：一个简单的实际项目案例为例，演示如何应用Pandas进行数据处理和分析。我们将考虑一个销售数据分析项目。

项目背景：假设你是一家电子产品公司的数据分析师，公司销售了不同种类的产品，你需要分析销售数据以了解销售趋势、最畅销的产品、最佳销售渠道等。

步骤1：数据导入和数据初步观察

首先，导入销售数据，通常是一个CSV文件或Excel文件。使用Pandas的 `.read_csv()` 或 `.read_excel()` 方法。然后，通过使用 `.head()`、`.info()` 等方法来初步观察数据的结构。

```
1 import pandas as pd
2
3 # 导入销售数据
4 sales_data = pd.read_csv('sales_data.csv')
5
6 # 初步观察数据
7 print(sales_data.head())
8 print(sales_data.info())
```

步骤2：数据清洗和预处理

在这一阶段，进行数据清洗，包括处理缺失值、重复值、异常值等。也可以将日期列转换为日期时间对象，以便进行时间序列分析。

```
1 # 处理缺失值
2 sales_data.dropna(subset=['Product'], inplace=True)
3
4 # 处理日期列
5 sales_data['OrderDate'] = pd.to_datetime(sales_data['OrderDate'])
6
7 # 删除重复值
8 sales_data.drop_duplicates(inplace=True)
```

步骤3：数据分析和可视化

使用Pandas和其他库（如Matplotlib或Seaborn），可以执行各种分析和可视化操作。例如，可以计算每个产品类别的销售总额，创建销售趋势图，绘制产品销售份额的饼图等。

```
1 # 计算每个产品类别的销售总额
2 category_sales = sales_data.groupby('ProductCategory')['Sales'].sum()
3
4 # 创建销售趋势图
5 import matplotlib.pyplot as plt
6 sales_data.set_index('OrderDate')['Sales'].resample('M').sum().plot()
7 plt.title('Monthly Sales Trend')
8 plt.xlabel('Date')
9 plt.ylabel('Sales')
10 plt.show()
```

```
11
12 # 创建产品销售份额的饼图
13 category_sales.plot(kind='pie', autopct='%1.1f%%')
14 plt.title('Product Sales Share')
15 plt.show()
```

步骤4：得出结论和建议

在分析数据后，可以得出结论并提出建议。例如，可能会发现某个产品类别的销售额在某个季度下降，建议进行促销活动。或者可能会发现某个销售渠道的表现非常出色，建议增加在该渠道的投资。

以上是一个简单的销售数据分析项目案例，演示了如何使用Pandas进行数据处理和分析。实际项目通常更复杂，可能涉及更多的数据清洗、特征工程、统计分析和机器学习等步骤。Pandas是数据分析的重要工具，帮助分析师从数据中提取有价值的信息并支持决策制定。

更多 Python 相关干货 内容，扫码领取！！

公众号：涛哥聊Python



干货资料领取：

- 1、【优质资料】优质资料合集
- 2、【学习路线】全方位知识点框架
- 3、【问题】Python各领域常见问题
- 4、【面试】面试指南

也欢迎大家围观我的朋友圈，日常会分享技术相关、副业与创业思考等！

添加涛哥 VX：2 57735，围观朋友圈，一起学 Python



公众号涛哥聊Python