

Pandas

xbZhong

2024-08-19

[本页PDF](#)

数据读取

常用类型的读取

```
1 import pandas as pd
2 ## 读取csv、tsv、txt数据
3 csv = pd.read_csv('path')
4
5 ## 读取excel文件
6 excel = pd.read_excel('path')
7
8 ## 读取sql文件
9 sql = pd.read_sql('path')
10
11 ## 读取csv文件时index_col为0则可以忽略多余的索引列
12 reviews = pd.read_csv('../input/wine-reviews/winemag-
    data_first150k.csv', index_col = 0)
```

读取txt文件，自己指定分隔符、列名

```
1 txt = pd.read_csv(
2     'path'
3     # 指定分割字符
4     sep = '\t'
5     # 无标题行
6     header = None
7     # 自定义列名
8     names = ['pdate', 'pv', 'uv']
9 )
```

常用操作

```
1 ## 查看前几行数据  
2 csv.head()  
3 ## 查看数据形状  
4 csv.shape  
5 ## 查看列名列表  
6 csv.columns  
7 ## 查看索引列  
8 csv.index  
9 ## 查看每列的数据类型  
10 csv.dtypes
```

导出文件

to_csv方法

```
1 ## 将animals导出为csv  
2 animals = pd.DataFrame({'Cows': [12, 20], 'Goats': [22, 19]}, index=['Year 1',  
'Year 2'])  
3 animals.to_csv('cows_and_goats.csv')
```

数据结构

- `Series` :一维数据，一行一列

```
1 ## 自定义Series与index  
2 s1 = pd.Series([1,'a',5.2,7],index=['d','b','a','c'])  
3 ingredients = pd.Series(['4 cups','1 cup','2 large','1 can'],index =  
['Flour','Milk','Eggs','Spam'],name = 'Dinner')  
4 ## 获取索引  
5 s1.index  
6  
7 ## 获取值  
8 s1.values  
9  
10 ## 根据index访问values  
11 print(s1['a'])  
12 print(type(s1['a']))  
13
```

- `DataFrame` :二维数据，多行多列

```
1 ## 查看列类型
2 df.dtypes
3 ## 查看行索引
4 df.index
5 ## 查看列索引
6 df.columns
7
8 ## 查询一列, 结果是Series
9 df['year']
10 ## 查询多列, 结果是DataFrame
11 df[['year', 'pop']]
12
13 ## 查询第一行
14 df.loc[1]
15 ## 查询多行(包含第三行, 和python语法有差别)
16 df.loc[1:3]
```

创建DataFrame

```
1 import pandas as pd
2 fruits = pd.DataFrame({'Apple':[30], 'Bananas':[21]})
```

数据类型

几种方法

```
1 ## 查看每一列数据类型
2 reviews.dtypes
3 ## 转换数据类型
4 reviews.points.astype('float64')
```

数据删除

- `drop(columns=['列名'])`

```
1 ## 删除主键id的列
2 data = data.drop(columns=['主键id'])
```

- `drop_duplicates(subset = ['列1', '列2', ...], keep = 'first')`
 - 根据subset里面的列去重
 - keep: 控制如何保留重复值的哪一行, 常见的选项有:
 - first: 保留第一次出现的重复值

- last: 保留第一次出现的重复值
- False: 删除所有重复项

```
df = df.drop_duplicates(subset=['监测点id', '监测时间', '企业DeptId'], keep='first')
```

数据查询

几种查询方法

==.iloc用数字索引左闭右开==

==.loc用数字索引左闭右闭==

```
1 ## 根据行、列的标签值查询
2 df.loc
3 ## 根据行、列的数字位置查询
4 df.iloc
```

.loc 的几种查询方法

```
1 ## 使用单个label的值来查询数据
2 df.loc['2018-01-03','bWendu']
3 ## 使用值列表批量查询
4 df.loc[['2018-01-03','2018-01-04','2018-01-05'],['bWemdu','yWendu']]
5 ## 使用数值区间进行查询
6 df.loc['2018-01-03':'2018-01-05','bWendu']
7 ## 使用条件表达式查询
8 df.loc[df['yWendu'] < -10,:]
```

用列名读取数据

```
1 ## 属性值查询列
2 reviews.country
3 ## 字典方式查询列
4 reviews['country']
5 ## 查询特定值
6 reviews['country'][0]
```

补充方法

```
1 ## isin方法查询country是Italy和France的行
2 reviews.loc[reviews.country.isin(['Italy', 'France'])] # 返回的是一行
3
4 ## .idxmax()返回具有最大值元素的索引
5 (reviews.points / reviews.price).idxmax()
```

新增数据列

几种方法

```
1 ## 直接赋值
2 df.loc[:, 'Wencha'] = df['bWendu'] - df['yWendu']
3 ## 使用apply方法
4 def get_wendu_type(x):
5     if x['bWendu'] > 33:
6         return '高温'
7     if x['bWendu'] < -10:
8         return '低温'
9     return '常温'
10 df.loc[:, 'wendu_type'] = df.apply(get_wendu_type, axis = 1)
11 ## assign方法，不会修改原表，返回的是新的表
12 df.assign(
13     yWendu_huashi = lambda x : x['yWendu'] * 9 / 5 + 32.
14     bWendu_huashi = lambda x : x['bWendu'] * 9 / 5 + 32
15 )
```

新增行

```
1 ## 新增title行
2 reviews.set_index('title')
```

Pandas数据统计函数

1. 汇总类统计

```
1 ## 提取所有数字列统计结果，非数字列也可以统计
2 df.describe()
3 ## 平均值
4 df['bWendu'].mean()
5 ## 最大值
6 df['bWendu'].max()
7 ## 中位数
8 df['bWendu'].median()
```

2. 唯一去重和按值计数

```
1 ## 得到不同类别(去重)
2 df['fengxiang'].unique()
3 ## 根据不同类别计数
4 df['fengxiang'].value_counts()
```

3. 相关系数和协方差

```
1 ## 协方差矩阵：衡量同向反向程度
2 df.cov()
3 ## 相关系数矩阵：衡量相似度程度
4 df.corr()
5 ## 单独查看两变量的相关系数
6 df['api'].corr(df['bWendu'])
```

Pandas缺失值数据处理

检测是否是空值

```
1 ## 查询整个DataFrame是否为空值
2 studf.isnull()
3 ## 查询某个列是否为空值
4 studf['分数'].isnull()
5 ## 与isnull相反
6 studf.notnull()
7 studf['分数'].notnull()
```

丢弃，删除缺失值

`dropna` : 丢弃、删除缺失值

- `axis`: 删除行还是列
- `how`: 为`any`则任何值为空都删除，为`all`则所有值都为空才删除

- inplace: 为True则修改当前df, 否则返回新的df

```

1 ## 删除全是空值的列
2 studf.dropna(axis = 1, how = 'all', inplace = True)
3 ## 删除掉只要有空值的行
4 studf.dropna(axis = 0, how = 'any', inplace = True)

```

填充空值

fillna :填充空值

- value: 用于填充的值, 可以是字典或单个值
- method: 填充方式, ffill为使用前一个不为空的值填充, bfill为使用后一个不为空的值填充
- axis: 按行还是列填充
- inplace: 为True修改当前df, 否则返回新的df

```

1 ## 将分数列为空的填充为0分
2 studf.fillna({'分数':0})
3 ## 等同于
4 studf.loc[:, '姓名'] = studf.fillna(0)
5 ## 将姓名的缺失值填充
6 studf.loc[:, '姓名'] = studf['姓名'].fillna(method = 'ffill')

```

替换

```

1 ## 用replace方法把taster_twitter_handle列的@kerinokeefe替换为@kerino
2 reviews.taster_twitter_handle.replace("@kerinokeefe", "@kerino")

```

重命名和合并

重命名

rename

- index或columns: 重命名行或列
- 用字典来实现
- inplace: 为True表示在原DataFrame上进行修改

```

1 ## 将列points改成score
2 reviews.rename(columns={'points': 'score'})

```

rename_axis

- 给行或列增加name属性

- name: name属性名字
- rows或columns: 属性是在行索引还是列索引

```
1 ## 给行索引添加name属性'wines'，给列索引添加name属性'fields'
2 reviews.rename_axis('wines',axis = 'rows').rename_axis('fields',axis = 'columns')
```

合并

`concat` 语法: `pandas.concat(objs, axis = 0, join = 'outer', ignore_index = False)`

- `objs` : 一个列表, 内容可以是 `DataFrame` 或者 `Series`
- `axis` : 0为按行合并, 1为按列合并
- `join` : 合并时索引对齐方式, 默认为outer,inner会过滤掉不匹配的列
- `ignore_index` : 是否忽略掉原来的位置索引
- `concat`会自动匹配索引的值, 但是没有定义的会填充NaN

```
1 ## 合并
2 pd.concat([df1,df2])
3 ## 按列合并
4 pd.concat([df1,df2],axis = 1)
```

join方法

- 根据索引自动匹配并合并值, 未匹配成功的会丢弃
- `other` : 要连接的另一个 `DataFrame`。
- `how` : 连接方式, 可以选择以下几种:
 - `'left'` : 使用左边的 `DataFrame` 的索引, 右边的 `DataFrame` 会填充 NaN。
 - `'right'` : 使用右边的 `DataFrame` 的索引, 左边的 `DataFrame` 会填充 NaN。
 - `'outer'` : 保留两个 `DataFrame` 的所有索引, 未匹配的地方填充 NaN。
 - `'inner'` : 只保留索引匹配的行。
- `on` : 指定用于连接的列 (如果未设置为索引)。
- `lsuffix` : 当两个 `DataFrame` 中存在重名列时, 为左边 `DataFrame` 的列添加后缀。
- `rsuffix` : 当两个 `DataFrame` 中存在重名列时, 为右边 `DataFrame` 的列添加后缀。
- `sort` : 是否根据连接后的索引排序, 默认为 False。

```
1 DataFrame.join(other, how='left', on=None, lsuffix='', rsuffix='', sort=False)
2 ## 根据两个DataFrame的MeetID进行合并
3 powerlifting_combined =
4 owerlifting_meets.set_index('MeetID').join(powerlifting_competitors.set_index('MeetI
```

映射

map

```
1 ## 求得points列的绝对值，用map定义lambda函数自减，p指的是reviews自己
2 review_points_mean = reviews.points.mean()
3 reviews.points.map(lambda p: p - review_points_mean)
4 ## 计算description列中tropical和fruity出现的次数
5 descriptor_counts = pd.Series([reviews.description.map(lambda x : 'tropical' in
6 x).sum(), reviews.description.map(lambda y: 'fruity' in y).sum()], index =
7 ['tropical','fruity'])
```

apply

- axis为columns表示

```
1 ## 可以用于行
2 def remean_points(row):
3     row.points = row.points - review_points_mean
4     return row
5
6 reviews.apply(remean_points, axis='columns')
7 ## 定义函数，将row中的country为Canada的星级设为3颗星，然后根据分数来评星级
8 def method(row):
9     if row.country == 'Canada':
10         return 3
11     elif row.points >= 95:
12         return 3
13     elif row.points >= 85:
14         return 2
15     return 1
16
17 star_ratings = reviews.apply(method, axis = 'columns')
18
19 ## Check your answer
20 q7.check()
```

数据分组及排序

分组

groupby分组

```
1 ## 按照points进行分组并对points进行数量计算，具有相同points的行会被分到同一组
2 reviews.groupby('points').points.sum()
3 ## 可以接受多个参数并使用.apply()方法
4 reviews.groupby(['country', 'province']).apply(lambda df:
5     df.loc[df.points.idxmax()])
6
7 ## 按照price来分组，选取最大的points并按照points来排序
8 best_rating_per_price = reviews.groupby('price')['points'].max().sort_index()
```

agg方法

- 参数为一个列表，里面的元素是要调用的函数

```
1 ## 可以接收多个参数，运行一系列不同的函数
2 reviews.groupby(['country']).price.agg([len, min, max])
```

多索引问题

```
1 ## 转换回单索引，但前面的多索引会变成新的列
2 countries_reviewed.reset_index()
```

排序

- ascending: 默认为True升序排序，为False降序排序
- inplace: 是否替换原始Series
- by: 按照哪一列进行排序，可以接收一个列表

```
1 ## 按值进行排序
2 df['tianqi'].sort_values()
3 ## 按照len列进行排序
4 countries_reviewed.sort_values(by='len')
5 ## 接受一个列表进行排序
6 countries_reviewed.sort_values(by=['country', 'len'])
7 ## 按索引进行排序
8 countries_reviewed.sort_index()
```