

1

10

```
import pandas as pd
## csv\tsv\txt\...
csv = pd.read_csv('path')

## excel\...
excel = pd.read_excel('path')

## sql\...
sql = pd.read_sql('path')

## csv\...\index_col\0\...\...\...\...\...
reviews = pd.read_csv('../input/wine-reviews/winemag-data_first150k.csv', index_col = 0)
```

txt

```
txt = pd.read_csv(  
    'path'  
    # 读取文件  
    sep = '\t'  
    # 分隔符  
    header = None  
    # 不读取首行  
    names = ['pdate', 'pv', 'uv'])  
)
```

1

1

to_csv

```
## animals.csv  
animals = pd.DataFrame({'Cows': [12, 20], 'Goats': [22, 19]}, index=['Year 1', 'Year
```

```
2'])
animals.to_csv('cows_and_goats.csv')
```

动物

- Series :
动物

```
## Series:index
s1 = pd.Series([1,'a',5.2,7],index=['d','b','a','c'])
ingredients = pd.Series(['4 cups','1 cup','2 large','1 can'],index =
['Flour','Milk','Eggs','Spam'],name = 'Dinner')
## 动物
s1.index

## 动物
s1.values

## index[]values
print(s1['a'])
print(type(s1['a']))
```

- DataFrame :
动物

```
## 动物
df.dtypes
## 动物
df.index
## 动物
df.columns

## 动物Series
df['year']
## 动物DataFrame
df[['year','pop']]

## 动物
df.loc[1]
## 动物(动物python动物)
df.loc[1:3]
```

动物DataFrame

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

动物

动物

```
## 评论数据
reviews.dtypes
## 评论点数
reviews.points.astype('float64')
```

评论点数

- `drop(columns=['点数'])`

```
## 去除id列
data = data.drop(columns=['id'])
```

- `drop_duplicates(subset = ['id', '评论', ...], keep = 'first')`
 - `subset`参数
 - `keep`参数
 - `first`保留第一个
 - `last`保留最后一个
 - `False`不保留任何

```
df = df.drop_duplicates(subset=['id', '评论', 'DeptId'], keep='first')
```

索引

索引方法

`==.iloc[索引]`
`==.loc[索引]`

```
## 通过索引
df.loc
## 通过label
df.iloc
```

`.loc` 索引方法

```
## 通过label索引
df.loc['2018-01-03', 'bWendu']
## 通过索引
df.loc[['2018-01-03', '2018-01-04', '2018-01-05'], ['bWendu', 'yWendu']]
## 通过切片
df.loc['2018-01-03':'2018-01-05', 'bWendu']
## 通过筛选
df.loc[df['yWendu'] < -10, :]
```

索引方法

```
## 评论国家
reviews.country
## 评论国家
reviews['country']
## 第一个
reviews['country'][0]
```

评论国家

```
## isin() country Italy France
reviews.loc[reviews.country.isin(['Italy', 'France'])] # 评论国家

## .idxmax() 返回最大值索引
(reviews.points / reviews.price).idxmax()
```

评论类型

评论类型

```
## 评论
df.loc[:, 'Wencha'] = df['bWendu'] - df['yWendu']
## apply()
def get_wendu_type(x):
    if x['bWendu'] > 33:
        return '高'
    if x['bWendu'] < -10:
        return '低'
    return '中'
df.loc[:, 'wendu_type'] = df.apply(get_wendu_type, axis = 1)
## assign()
df.assign(
    yWendu_huashi = lambda x : x['yWendu'] * 9 / 5 + 32.
    bWendu_huashi = lambda x : x['bWendu'] * 9 / 5 + 32
)
```

评论

```
## 以title为索引
reviews.set_index('title')
```

Pandas基础操作

1. 基本操作

```
## 基本操作
df.describe()
## 基本操作
```

```
df['bWendu'].mean()  
## 均值  
df['bWendu'].max()  
## 最大值  
df['bWendu'].median()
```

2. 亂數統計

```
## 獨特值(數)  
df['fengxiang'].unique()  
## 獨特值數量  
df['fengxiang'].value_counts()
```

3. 相關性

```
## 資料框：相關性矩陣  
df.cov()  
## 資料框：相關性  
df.corr()  
## 資料框：相關係數  
df['api'].corr(df['bWendu'])
```

Pandas資料處理

資料框

```
## DataFrame資料框  
studf.isnull()  
## 檢查空值  
studf['國'].isnull()  
## 判斷某欄位是否為空  
studf.notnull()  
studf['國'].notnull()
```

刪除空值

`dropna` 用法

- `axis` 設定刪除的軸
- `how` 設定刪除的標準，`'any'` 或 `'all'`
- `inplace` 設定 True，將修改原 DataFrame

```
## 刪除空值  
studf.dropna(axis = 1, how = 'all', inplace = True)  
## 刪除空值  
studf.dropna(axis = 0, how = 'any', inplace = True)
```

總結

fillna : 補充

- value : 値
- method : ffill, bfill
- axis : 軸
- inplace : True

```
## 補充
studf.fillna({'Age':0})
## 補充
studf.loc[:, 'Age'] = studf.fillna(0)
## 補充
studf.loc[:, 'Age'] = studf['Age'].fillna(method = 'ffill')
```

□□

```
## 替換
reviews.taster_twitter_handle.replace("@kerinokeefe", "@kerino")
```

□□□□□

□□□

rename

- index/columns : 軸
- フィールド
- inplace : True

```
## points → score
reviews.rename(columns={'points': 'score'})
```

rename_axis

- rows/columns : 軸
- name : 軸名
- rows/columns : 軸

```
## wines → 'wines'
## fields → 'fields'
reviews.rename_axis('wines', axis = 'rows').rename_axis('fields', axis = 'columns')
```

□□

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

- objs : DataFrame オブジェクト
- axis : 0 (rows) / 1 (columns)
- join : outer, inner

- ignore_index 亂^{ignore_index}＼
- concat乱^{concat}＼NaN

```
## 亂
pd.concat([df1,df2])
## 亂
pd.concat([df1,df2],axis = 1)
```

join

- 亂^{join}
- ** other 亂^{DataFrame}
- ** how 亂^{join_type}
 - 'left' 亂^{DataFrame DataFrame} DataFrame 亂^{NaN}
 - 'right' 亂^{DataFrame DataFrame} DataFrame 亂^{NaN}
 - 'outer' 亂^{DataFrame DataFrame} DataFrame 亂^{NaN}
 - 'inner' 亂^{DataFrame}
- ** on 亂^{join_keys}
- ** lsuffix 亂^{DataFrame DataFrame}
- ** rsuffix 亂^{DataFrame DataFrame}
- ** sort 亂^{False}

```
DataFrame.join(other, how='left', on=None, lsuffix='', rsuffix='', sort=False)
## 亂
powerlifting_combined =
owerlifting_meets.set_index('MeetID').join(powerlifting_competitors.set_index('MeetID'))
```

map

```
## 亂points 亂map乱lambda p reviews
review_points_mean = reviews.points.mean()
reviews.points.map(lambda p: p - review_points_mean)
## 亂description 亂tropical fruity
descriptor_counts = pd.Series([reviews.description.map(lambda x : 'tropical' in
x).sum(), reviews.description.map(lambda y: 'fruity' in y).sum()], index =
['tropical','fruity'])
```

apply

- axis＼columns＼＼

```
## 亂
def remean_points(row):
    row.points = row.points - review_points_mean
    return row

reviews.apply(remean_points, axis='columns')
```

```
## 用來根據row的country（Canada或USA）和points（300或以上）來判斷
def method(row):
    if row.country == 'Canada':
        return 3
    elif row.points >= 95:
        return 3
    elif row.points >= 85:
        return 2
    return 1

star_ratings = reviews.apply(method, axis = 'columns')

## Check your answer
q7.check()
```

ANSWER

2

groupby

```
## 计算评论数
reviews.groupby('points').points.sum()
## 按国家和省份分组，计算平均评分
reviews.groupby(['country', 'province']).apply(lambda df:
df.loc[df.points.idxmax()])
## 根据价格分组，计算平均评分
best_rating_per_price = reviews.groupby('price')['points'].max().sort_index()
```

agg□□

-

```
## 中国电影票房  
reviews.groupby(['country']).price.agg([len, min, max])
```

1

```
## 中国的国家
countries_reviewed.reset_index()
```

2

- ascending
 - inplace
 - by

```
## 重新排序
df['tianqi'].sort_values()
## 按len重新排序
countries_reviewed.sort_values(by='len')
## 按country和len重新排序
countries_reviewed.sort_values(by=['country', 'len'])
## 重新索引
countries_reviewed.sort_index()
```