# pandas库

#### readme

Pandas 库在处理excel类型数据非常好用,常见业务对接数据往往通过excel收集,学习之前请确保学习完成 li 正则表达式

#### 1. 导入库

```
1 import pandas as pd bk/7020 bk/7020 bk/7020 bk/7020 bk/7020 bk/7020
```

#### 2. 创建DataFrame

A,B,C为列名,a,b,c为行名

```
ABCax100010by200020cz300030
```

# 3. 查看DataFrame中的值 df.values

```
1 df.values
```

# 4. 获取行列名df.index, df.columns

```
1 df.index

2 df.columns 20 bix 7020 bix 7020 bix 7020 bix 7020 bix 7020 bix 7020
```

```
Index(['a', 'b', 'c'], dtype='object')
Index(['A', 'B', 'C'], dtype='object')
```

• 创建时可以更改列的顺序,只需在 columns 处改变顺序



# 5. 添加、删除一列 del df['D']

```
    1 df['D'] = 10 # 把D列全部赋值为10

    2 del df['D']
```

### 6. 横行合并两个DataFrame

- how = 'left', 'right', 'inner', 'outer'
- 这里和sql的合并表方式类似

```
1 df2 = pd.DataFrame(data = {'x1':['A','B','C'],'x2':[1,2,3]})
2 df3 = pd.DataFrame(data = {'x1':['A','B','D'],'x2':['T','F','T']})
3 pd.merge(df2,df3,how = 'left', on = 'x1')
```

```
      x1
      x2_x
      x2_y

      0
      A
      1
      T

      1
      B
      2
      F

      2
      C
      3
      NaN
```

```
1 pd.merge(df2,df3,how = 'right', on =
  'x1')
```

```
        x1
        x2_x
        x2_y

        0
        A
        1.0
        T

        1
        B
        2.0
        F

        2
        D
        NaN
        T
```

```
1 pd.merge(df2,df3,how = 'outer', on =
  'x1')
```

```
        x1
        x2_x
        x2_y

        0
        A
        1.0
        T

        1
        B
        2.0
        F

        2
        C
        3.0
        NaN

        3
        D
        NaN
        T
```

```
pd.merge(df2,df3,how = 'inner', on =
'x1')
```

```
    x1
    x2_x
    x2_y

    0
    A
    1
    T

    1
    B
    2
    F
```

## 7. 常见统计操作

```
1 import numpy as np
2 df = pd.DataFrame(np.random.randn(8,4),columns = list('ABCD'))
3 df
```

```
0.132863
            0.228572
                       0.874787
                                  -0.261211
0.433415
            0.626160
                       0.060306
                                  -0.144069
0.337980
            0.627468
                       0.795961
-0.652359
           0.298969
                       -1.198896
                                  -0.061287
0.690291
           -0.782332
                       0.846143
                                 -0.026827
0.788389
           -0.096860
                       0.164575
                                  -3.322447
-1.452064
           -0.192849
                       0.190505
                                   1.257995
                       1.736579
-0.997478
           -0.737240
                                   0.041573
```

```
1 df.sum()
2 df.mean()
3 df.var()
4 df.std()
5 df.cumsum() # 从上往下累加
6 df.describe()
7 df.info()
8 df.min()
9 df.max()
10 df['A'].value_counts() # 统计A列的数据分布
11 df.shape
```

### 8. DataFrame切片

```
1 data = {'A':['x','y','z'],
           'B':[1000,2000,3000],
2
 3
           'C':[10,20,30]}
4 df = pd.DataFrame(data,index = ['a','b','c'])
5
6 # 提取名为A的列
7 df['A']
8 df.A
9
10 # 提取名为a的行
11 df.loc['a']
12
13 # 提取a到c行,A和B列
14 df.loc['a':'c',['A','B']]
15
16 # 提取0-2行, 0, 1列
17 df.iloc[[0,1,2],0:2]
```

- df.loc[] 依据行列名, df.iloc[] 依据行列号
- df.iloc[] 使用非常多,经常搭配遍历 for i in range(len(df)): df.iloc[i,1] = '' 使用
- 即从上往下把每一行到第一列进行一个赋值

#### 9. 通过筛选条件切片

```
1 df.C > 10

blkT020
2 # output: T020
3 # a False
4 # b True
blkT020
```

# 10. 改变一列的值

# 11. 读取、存储excel

• 对excel的数据实际落地项目中使用较多,因此读取、存储非常关键

```
1 df = pd.read_csv('name4excel.csv')
2 df = pd.read_excel('name4excel.xlsx')
3
4 df.to_csv('output.csv', encoding = 'utf-8-sig', index = False)
```

#### 12. 删除重复项

```
В
          C
Α
   1000
         10
У
   1000
         10
  1000
         30
Z
   2000 30
        В
            C
     1000
            10
     2000
           30
```

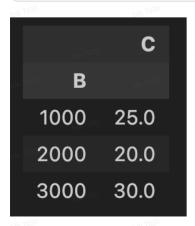
```
C
          В
      1000
              10
а
   Χ
      1000
             10
   У
      1000
             30
C
      2000
             30
            В
                 C
         1000
                10
     Χ
 а
         1000
                30
 С
         2000
                30
 d
```

## 13. 删除存在空值的行

```
В
                  C
   Α
       1000.0
                 10
а
   X
b
       1000.0
                 10
   У
       1000.0
C
   Z
                 30
d
   t
          NaN
                 30
              В
                   C
     X
         1000.0
                   10
 а
         1000.0
 b
                   10
         1000.0
                   30
     Z
```

# 14. 列相同的值合并分组操作

• B列相同的值分为一组,对C列求和



#### 15. 排序

#### 16. 映射到二维平面

• 个人从来没使用过

```
1 df = pd.DataFrame({'client':['John','John','Silvia','Silvia'],'product':
    ['bananas','oranges','bananas','oranges'],'quantity':[5,3,4,2]})
2 print(df)
3 df.pivot(index = 'client',columns = 'product',values = 'quantity')
```

```
client
           product
                    quantity
0
     John
           bananas
1
     John
           oranges
   Silvia bananas
   Silvia oranges
 product
          bananas
                    oranges
   client
    John
                          2
   Silvia
```

# 17. 提取某列带特定字符的列表

```
1 df = pd.DataFrame({'A':['a','b','c'],'B':['dc','ec','fc']})
2 df[df['B'].str.contains('d')]
```

• 思考题:是否可以如法炮制使用其他的字符串函数,比如lower,upper那些

# 18. 将字符'[]'转回列表

```
1 \ \mathsf{df['A']} =  \mathsf{df['A'].apply(eval)}_{\mathsf{bk},\mathsf{TO20}}
```

- 在读取excel时, [ ] 往往被看作是字符,不能被识别为list,需要上述操作更改
- 大家以后一定会遇到这个报错回来看这个解决方案

### 19. 爆炸函数

```
1 a = pd.DataFrame({'A':['a','b','c'],'B':[[1,2,3],[4,5],[6]]})
2 a.explode('B')
```

blk 7020	Α	В	
0	а	1	
0	a	2	
0	а	3	
1	b	4	
1	b	5	
2	С	6	