

Python for Data Analysis

Author : Yunqiu Xu

DataAnalysis Python numpy pandas matplotlib

- References:
 - Book: Python For Data Analysis
 - [python数据分析入门](#)
 - [利用python进行数据分析 FDU](#)
 - [numpy 官方文档](#)
 - [pandas 官方文档](#)
 - [scipy 官方文档](#)
- numpy,pandas,matplotlib,scipy

Chapter 1 Introduction

- Import the data;CSV
- Data transformation;
- Statistical description;
- Hypothesis test (t test);
- Visualization;
- Function definition;

1.1 导入数据

```
1. open(file) #内置函数open()
2.
3. import pandas as pd
4. df=pd.read_table()
5.
6. unames=['uid','gender','age','occupation','zip']
7. users=pd.read_table('users.dat',sep='::',header=None,names=unames)
8. df=pd.read_csv() #csv
9. df=pd.DataFrame() #将已经打开的txt转化为数据框架形式
```

1.2 数据转换

1.2.1 基本转换方式

```
1. df.head()/df.tail() #这个类似R, 只不过R给出6行, Python给出5行
2. df.columns/df.index #相当于R里的colnames/rownames
3. df.T #等价于numpy.transpose
4. df.ix('colname') #提取特定的某列数据, 注意python索引从0开始
5. df.ix[:,0].head() #提取第1列前5行
6. df.ix[10:20,0:3] #11-20行[10,19)的前三列数据[0,3)
7. df.drop(df.columns[[1,2]],axis=1)#舍弃前两列, axis=0为舍弃行
```

1.2.2 透视表 `pandas+numpy`

```
1. pd.pivot_table(df,
2.     index=["Name","Rep"], #设置行名
3.     values=["price"], #设置列名, 未指定则默认其余所有内容为列值
4.     aggfunc=[numpy.sum], #函数, 未指定则默认展示出各列列值
5.     columns=['Product'], #也是设置列名, 但aggfunc只会作用于value
6.     fill_value=0, #将缺失值填充为0
7.     margins=True, #汇总数据, 默认为False
8. )
```

1.3 统计描述

```
1. df.describe() #给出count/mean/std/min/25%/50%/75%/max
```

1.4 假设检验

```
1. from scipy import stats as ss
```

- 举个栗子：检验Abra列均值，假定总产量均值为15000

```
1. ss.ttest_1samp(a = df.ix[:, 'Abra'], popmean = 15000)
2. >>>print后返回tuple, 包括t,p,通过p即可验证假设是否成立
```

1.5 可视化

1.5.1 .plot() in pandas

```
1. table.plot(title='Sum of table1k.prop by year & sex',
2.            yticks=np.linspace(0, 1.2, 13),
3.            xticks=range(1880, 2020, 10))
```

1.5.2 matplotlib

```
1. import matplotlib.pyplot as plt
2. fig=plt.figure #给出画布
3. ax=plt.add_subplot(111) #将画布分成1行1列, 获取从上到下, 从左到右第1块小画布
4. add_subplot(3,4,6) #将画布分成3行4列要在第6块画布作画
5. ax.plot(dfx,dfy,kind='box')
6. plt.show
```

- flatten: 数据扁平化(只可用于array/matrix,list不可用)

```
1. a = array([[1,3],[2,4],[3,5]])
2. a.flatten#[1,3,2,4,3,5]
```

Chapter 2 数据分析实例：1880-2010 年间 全美婴儿姓名的趋势

2.1 导入并处理1880年的数据

- 该数据集可分为三列：名字;性别;每个名字的婴儿数量

```
1. import pandas as pd
2. path='C:/xyq data/py projects/pydata/names/yob1880.txt'
3. names1880 = pd.read_csv(path, names= ['name', 'sex', 'births'])
4. names1880.groupby('sex').births.sum() #统计当年出生的男女数量 (打乱名字只按性别分配, 同时累加第三列)
```

2.2 整合多个年份的数据

- 将1880 - 2010 的数据统计文件整合进一个数据集

```
1. years = range(1880, 2011 )
```

```

2.
3. pieces = [] #将每个数据集作为pieces的一个元素
4. columns = names = ['name', 'sex', 'births']
5.
6. for year in years:
7.     path = 'names/yob%d.txt' % year
8.     frame = pd.read_csv(path, names = columns)
9.     frame['year'] = year
10.    pieces.append(frame)
11.    #得到的pieces是一个长度为131的list, 每个元素都是一年的数据
12.
13.    #使用pd.concat()将各年数据整合进一个数据集
14.    names = pd.concat(pieces, ignore_index = True) #忽略read_csv返回的原始行号

```

2.3 数据处理

2.3.1 统计每年(行)出生的男婴和女婴数量

```

1. from pandas import pivot_table
2. total_births = pivot_table(names,
3.                             values='births', #列值
4.                             index='year', #行值
5.                             columns='sex', #列值, 但不受aggfunc影响
6.                             aggfunc=sum #求和函数, 作用于values
7.                             )
8. total_births.tail() #数据预览
9.
10. import matplotlib.pyplot as plt
11. fig=plt.figure()
12. ax=fig.add_subplot(111)
13. ax.plot(x,y,title='Total births by sex and year')
14. plt.show()

```

2.3.2 插入新列prop, 记录指定婴儿数量相对于总出生人数的比例

```

1. def add_prop(group):
2.     births = group.births.astype(float) #注意不可用整除
3.     group['prop']=births/births.sum()
4.     return group
5. names = names.groupby(['year', 'sex']).apply(add_prop)
6. #apply函数可用于调用参数已经存在于元组或字典的某个函数

```

2.3.3 选取每年排名前1000的姓进行统计(sex/year的前1000名)

- `sort_index` in `DataFrame`

```
1. from pandas import DataFrame as df
2. def get_top1k(group):
3.     return group.sort_index(by= 'births', ascending = False)[:1000] #前
    1000名, 降序排列
4.
5. grouped = names.groupby(['year', 'sex'])
6. top1k = grouped.apply(get_top1k)
```

2.4 分析命名趋势

2.4.1 按性别分类并建立透视表

```
1. from pandas import pivot_table
2. boys=top1k[top1k.sex=='M']
3. girls=top1k[top1k.sex=='F']
4.
5. total_births = pivot_table(top1k, values='births', index= 'year', columns = 'name', aggfunc = sum)
```

2.4.2 评估某几个姓名使用率随年份的变化

```
1. subset = total_births[['John', 'Harry', 'Mary', 'Marilyn' ]]
2. subset.plot(subplots=True,
3.             figsize = (12,10),
4.             grid= False,
5.             title='Number of births per year')
```

2.4.3 评估命名多样性的增长 (计算前1000个名字在总数中的prop)

```
1. import numpy as np
2. table = pivot_table(top1k, values='prop', index='year', columns='sex', aggfunc = sum)
3. table.plot(title='Sum of table1k.prop by year & sex',
4.            yticks=np.linspace(0, 1.2, 13),
5.            xticks= range(1880, 2020, 10))
```

Chapter 3 Numpy

```
1. from numpy import *
```

3.1 创建Array

```
1. arr=array([[1,2,3],[4,5,6]],dtype=' ') #创建一个2*3 array
2. a=array([[[1,2,3],[4,5,6]],[[7,8,9],[10,11,12]]) #2*2*3
3. zeros((3,4))/empty((3,4)) #零矩阵
4. ones((2,3,4), dtype=int16) #2*3*4阵列, 元素全为1
5. eye()/identity() #单位矩阵
```

```
1. arange(10,30,5) #[10,30)间隔为5, array([10,15,20,25])
2. arange().reshape() #常用创建矩阵方法, 先arange再reshape
3. linspace(0,2,9) #[0,2)中均匀选取9个数字
```

```
1. from np import random #对内置的random进行补充
2. samples = np.random.normal(size=(4, 4)) #4*4 normal distribution
3. random.random((2,12)) #2*12 随机分布
4. #若样本量很大, np.random的速度更快
```

3.2 查看矩阵properties

3.2.1 基本查看方式

```
1. arr.ndim #dimension,此处返回2
2. arr.shape #给出矩阵某个维度的长度, 此处返回(2L, 3L)表示2行3列
3. shape[0] #某个矩阵的第一维长度
4. shape(3) #单独数字, 返回为空()
5. shape([3]) #一维矩阵, 返回(1L,)
6. arr.size #元素数量
7.
8. arr.dtype.name #查看数据类型, 默认为int32
9. arr.dtype('float64') #设置矩阵类型
10. arr.itemsize #每个元素的大小(bytes)
11. float64=64/8=8
12. int32=32/8=4
```

3.2.2 nonzero

```
1. nonzero(x) #返回x阵列中非零元素的位置
```

```
2. x[nonzero(x)] #返回x阵列中非零元素值
```

- 举个栗子

```
1. x=eye(3)
2. nonzero(x) #(array([0, 1, 2]), array([0, 1, 2])) 即非零元素位置为(0,0),(1,1),(2,2)
3. x[nonzero(x)] #(array([0, 1, 1])
```

- 另一种形式 `transpose`

```
1. np.transpose(np.nonzero(x)) #array([[0, 0],[1, 1],[2, 2]])
```

3.3 基本操作

```
1. A+B #相加,等价于add(A,B)
2. B**2 #每个元素平方
3. A<35 #各元素返回为boolean
4. arrA*arrB #同位置元素相乘,类似octave A .* B
5. dot(arrA,arrB) #矩阵乘法
6. b+=a #等价于b=b+a
7. a*=3 #a=a*3
8. exp()/sqrt()/abs()/square()/log()/log10()/log2() #每个元素
9. sign() #各元素正负号
10. isnan() #各元素是否为NaN
11. .sum()/ .min()/ .max() #给出指定的元素
12.
13. A.sum(axis=0) #返回每列的sum
14. A.min(axis=1) #返回每行的min
15. A.cumsum(axis=0) #从上到下累加(第二行为新值+上一行值)
```

- 举个栗子:

```
1. b=array([[0,1,2,3],[4,5,6,7],[8,9,10,11]])
2. b.cumsum(axis=1) #行方向:从左到右累加
3. >>>array([[0,1,3,6],[4,9,15,12],[8,17,27,38]])
```

- 再举个栗子:

```
1. xxx=np.arange(0,12).reshape((2,3,2)) #2个3*2矩阵
```

```

2.     xxx
3.     >>>array([[ 0,  1],
4.              [ 2,  3],
5.              [ 4,  5]],
6.
7.              [[ 6,  7],
8.              [ 8,  9],
9.              [10, 11]])
10.
11.     xxx.sum(axis=0)
12.     >>>array([[ 6,  8],
13.              [10, 12],
14.              [14, 16]])
15.
16.     xxx.sum(axis=1)
17.     >>>array([[ 6,  9],
18.              [24, 27]]) #两个矩阵每列的和
19.
20.     xxx.sum(axis=2)
21.     >>>array([[ 1,  5,  9],
22.              [13, 17, 21]]) #两个矩阵每行的和

```

3.4 Indexing, Slicing and Iterating

```

1.     argsort(a) #排序, 返回元素在原数组的位置
2.         argsort([1, 2, 4, 9, 5, 3])->[0,1,5,2,4,3]
3.         argsort(-a) #反向排列(从大到小)
4.
5.     arr[-1]/arr[5:8,1]/arr[:,2] #可以查询或赋值某个元素, 单数字默认为行
6.         a[::2]=1000 #从索引0开始, 每两个元素赋值为1000
7.         a[::-1]#反向排列

```

- (...) 在保证默认条件的前提下包含所有其他维度

```

1.     #假定a为2*2*3array
2.     a[1,...] #数组中第二个矩阵的全部行列
3.     a[...,2] #数组中所有矩阵的第三列
4.
5.     for x in a.flat:
6.         print x #扁平化数组: 将所有元素打印为一列

```


3.5 变换阵列形式

3.5.1 改变行列

```
1.  arr = arange(24).reshape(4,3,2) #将一行向量转化为4*3*2阵列
2.
3.  #另一种方法：先设置shape再transpose()
4.  arr.shape=(4,3,2)
5.  arr.transpose() #这个也可以用于转置
6.
7.  arr.reshape(3,-1) #指定行数，自动计算列数
8.  arr.resize((2,6)) #同reshape，但直接改变原array，类似sort/sorted的关系
```

3.5.2 组合阵列

```
1.  vstack((a,b))/hstack((a,b)) #垂直/水平组合
2.  column_stack((a,b)) #新增一列b
3.  row_stack((a,b)) #新增一行
4.
5.  arr[:,newaxis] #增加新维度
6.  arr=array([[1,2,3],[4,5,6]]) # (2L, 3L)
7.  arr[:,newaxis]=array([[[1,2,3]],[[4,5,6]]]) # (2L, 1L, 3L)
```

- 举个栗子:

```
1.  a=array([1,2,3,4]) #4L
2.  b=array([5,6,7,8])
3.  c=a[:,newaxis] #4*1 array([[1],[2],[3],[4]])
4.  d=b[:,newaxis] #array([[5],[6],[7],[8]])
5.  column_stack((c,d)) #4*2
```

3.5.3 拆分阵列

```
1.  hsplit(a,3) #水平方向均匀拆分a为3个阵列（行不变，拆分行）
2.  hsplit(a,(3,4)) #分割点为第四列（前三列；第四列；剩下列）
3.  vsplit() #垂直方向，列不变，拆分行
```

3.6 Copy and View

- 赋值/浅复制/深复制

```

1.  b=a #赋值：相当于根本没复制，改变了b，a也会变化
2.  c=a.view() #shallow copy只拷贝父对象，不会拷贝对象的内部的子对象（子对象相当于赋值）
3.  d=a.copy() #deep copy什么都复制了，从此再怎么变化都不会对副本造成影响

```

3.7 使用数组查找数组元素

```

1.  a = arange(12)**2
2.  i = array([1,1,3,8,5])
3.  a[i] #a at the positions i
4.  >>>array([1,1,9,64,25])
5.
6.  j = array([[3,4],[9,7]])
7.  a[j]
8.  >>>等价于array([[ a[3], a[4]],[a[9], a[7]])

```

3.8 布尔数组

```

1.  a = arange(12).reshape(3,4)
2.  b = a > 4 #b中元素为True/False
3.  a[b] #得到a中令b为True(>4)的元素并组合成array
4.  a[b]=0 #此时查看a会发现这些元素被赋值为0

```

3.9 ix_() function

- 组合不同向量用于进行矩阵运算
- 个人感觉应该和pandas.ix()类似

3.10 Linear Algebra

```

1.  from numpy import *
2.  from numpy.linalg import *

```

3.10.1 构造数组

```

1.  a = array([[1.0, 2.0], [3.0, 4.0]]) #2*2
2.  print a #展示数组a

```

```
3. a.transpose() #转置, 等价于pandas.T
4. inv(a) #inversal matrix
```

3.10.2 求trace迹

```
1. u=eye(2)
2. trace(u) #trace (主对角线元素之和)
```

3.10.3 求解线性方程组

```
1. y=array([[5.],[7.]])
2. solve(a,y) #解为array([[ -3.],[ 4.]])
```

3.10.4 直接构造矩阵

```
1. A=matrix('1.0 2.0;3.0 4.0') #2*2
2. type(A)
3. >>> <class 'numpy.matrixlib.defmatrix.matrix'>
4.
5. A.T #transpose()
6. A.I #inv()
7.
8. Y=matrix('5.0;7.0') #2*1
9. print A*Y #矩阵乘法, 最后得到2*1矩阵
10.
11. solve(A,Y) #和之前一样, 求解线性方程组
```

Chapter 4 Pandas

- 两个重要的数据结构：Series and DataFrame

4.1 Series

```
1. from pandas import Series
2. #Series: 一维类似数组的对象, 包含一系列NumPy数据及一系列标签索引
```

4.1.1 创建Series

```
1. obj = Series([4, 7, -5, 3], index=['d', 'b', 'a', 'c'])
2. obj
```

```

3.     >>>d    4
4.     b    7
5.     a   -5
6.     c    3
7.     dtype: int64
8.
9.
10.    #基于dict创建(会自动排序)
11.    obj=Series(dictionary,index=newindex) #索引可以赋值成新索引,如果没有就用原
      dict里的key
12.
13.    obj.values #获取数据,这里可以回顾下Chapter2里的pivot_table
14.    obj.index #未指定则默认[0,1,2,3],也可以赋值新索引
15.    obj[['a','b']] #返回对应值,也可以对该索引进行赋值

```

4.1.2 类numpy操作

- 由于和NumPy共享数据类型,可以延续numpy的操作

```

1.     obj2[obj2 > 0]
2.     obj2 * 2
3.     np.exp(obj2)
4.     'b' in obj2 #True
5.     #算术运算时,不同索引数据会自动对齐

```

4.2 DataFrame

- 包含一个索引以及与这些索引联合在一起的Series

4.2.1 创建DataFrame

```

1.     from pandas import DataFrame
2.     data={'key1':[],'key2':[],'key3':[]}
3.     frame=DataFrame(data)
4.
5.     frame
6.     >>> Empty DataFrame
7.     Columns: [key1, key2, key3]
8.     Index: []
9.     #index是行名, keys是列名!
10.
11.     DataFrame(data,columns=[],index=[]) #重新指定col,index
12.     #DF中的每一个column可以被认为是一个Series

```

4.2.2 查看元素并赋值

```
1. frame['state']/frame.year #查看某列
2. frame.ix['index'] #查看某行
3.
4. #对某列赋值, 如果值数和列长度不等, 值会被随机填充, 空位为NaN
5.
6. frame['eastern'] = frame.state == 'Ohio' #增加一条名为eastern的新列, 列值为boolean
7.
8. del frame['eastern'] #删除列
```

4.3 Index Objects

```
1. index = obj.index #给出每行的名字
2. index
3. >>>index([a,b,c],dtype=object)
4. index[1:]
5. >>>b,c
6. #注意index是不可变的, 赋值会返回错误
7.
8. 2003 in frame3.index #类似查询value, 查询是否存在名为2003的index (注意2003不是value)
```

4.4 Reindexing

- 创建带有新索引的新对象, 相当于增加新行

```
1. obj=Series([4.5,7.2,-5.3,3.6],index=['d','b','a','c'])
2. obj2 = obj.reindex(['a','b','c','d','e']) # 'e'=NaN
3.
4. #替换缺失值为0
5. obj.reindex(['a','b','c','d','e'], fill_value=0)
6. #当然, 可以直接赋值以增加新行
7. obj["e"]="asshole"
8.
9. DataFrame(np.arange(9).reshape((3, 3)),
10.           index=['a','c','d'],
11.           columns=['Ohio','Texas','California'])
12. frame2 = frame.reindex(['a','b','c','d']) #默认为重新索引行
13. frame.reindex(columns=states) #重新索引列
```

```
14. #可以将上面两个结合在一起同时索引行列
```

4.5 删除条目

```
1. new_obj = obj.drop(['c', 'd']) #删除'c', 'd'这两行
2.
3. #对于DataFrame, 可以从任意坐标轴删除索引
4. data.drop('Ohio') #删除某行
5. data.drop('two') #删除某列, 等价于del frame[' ']
```

4.6 Indexing , selection and filtering

```
1. #Series索引类似NumPy
2. obj['b']/obj[1]/obj[2:4]/obj[obj<2]
3.
4. #注意, 切片跟python不同! 两边均包含而非()
5. obj['b':'c'] #索引b与c的值都会返回
6.
7. obj.ix[xxx] #DataFrame中某行
8. obj.ix[:,xxx] #某列
9. obj.ix[xxx,yyy] #行列
```

4.7 Sorting and Ranking

4.7.1 按index排序

```
1. seri.sort_index(
2.     axis=1, #按照某个维度排序
3.     ascending=False #降序排列
4. )
5. df.sort_index(by=['a', 'b']) #对于数据框架, 指定按哪一列索引排序
```

4.7.2 按value排序

```
1. obj.order()
```

4.7.3 Rank

- 与排序不同, Rank会把对象的 values 替换成名次(不一定是整数)

- 对于平级项，可选method: average/min/max/first
- 默认 ascending=True 从小到大
- 举个栗子：

```

1.  obj = Series([7, -5, 7, 4, 2, 0, 4])
2.  obj.rank()
3.  >>>6.5 1.0 6.5 4.5 3.0 2.0 4.5
4.
5.  obj.rank(method='first')
6.  >>>6 1 7 4 3 2 5
7.  #7的名次,-5的名次,...
8.
9.  obj.rank(ascending=False, method='max')
10. >>>2 7 2 4 5 6 4
11.
12. #对于DataFrame, 多了个axis项, 其余类似
13. frame.rank(axis=1)

```

4.8 Summarizing and Computing Descriptive Statistics

```

1.  df.sum(axis=1) #numpy.sum
2.  df.mean(axis=1, skipna=False) #每行平均值
3.  df.idxmax() #max的索引值
4.  df.argmax() #max的索引位置
5.  cumsum/cummin/cummax #累计和/最小最大值
6.  var/std/quantile

```

4.9 Handle Missing Values

```

1.  data.isnull() #查询缺失值
2.
3.  data.dropna() #删除缺失值, 等价于data[data.notnull()]
4.      data.dropna(how='all') #删除全是NA的行
5.      data.dropna(axis=1, how='all') #同理, 删除全是NA的列
6.
7.  df.fillna(0) #填充缺失值, 返回新object
8.      df.fillna({1: 0.5, 3: -1}) #不同列不同填充值
9.      data.fillna(data.mean()) #使用平均值来填充NA
10. df.fillna(0,inplace=True) #不返回新object, 直接在原位修改

```

Chapter 5 Data Loading , Storage and File Formats

5.1 Reading

```
1. pd.read_table(path,names,sep='')
2.
3. pd.read_csv(
4.     header=None, #设置默认header,没有则将第一行数据作为header
5.     na_values=['NULL'], #缺失值显示为'NULL'
6.     nrows=5, #仅仅读取前5行
7. )
8.
9. pd.from_csv()
10. np.loadtxt(delimiter) #np.array(pd.read_csv())
```

5.2 Writing

```
1. data.to_csv(
2.     "filename",
3.     sep="|",
4.     index=False,
5.     header=False,
6.     na_rep="NULL",
7.     cols=["a","b","c"]
8. ) #将DataFrame/Series 另存为csv
```

5.3 Other Kinds of Files / Modules

```
1. import csv
2. import json
```

- For HTML/XML : `urllib2/BeautifulSoup/requests`

5.4 Binary Data Formats

```
1. frame=pd.read_csv()
2. frame.save("xxx_pickle") #另存为二进制格式
```



```
3. pd.load("xxx_pickle") #使用python读取该文件
```

5.5 Excel:

```
1. xls_file=pd.ExcelFile("data.xls")  
2. table=xls_file.parse('sheet1') #转化为DataFrame
```

Chapter 6 Data Wrangling: Clean , Transform , Merge , Reshape

6.1 Combining and Merging

```
1. pd.merge()  
2. pd.concat()  
3. combine_first()
```

6.1.1 Database-style DataFrame Merges

```
1. pd.merge(df1,df2,on='key') #基于key组合两个DataFrame  
2. pd.merge(df3,df4,left_on='key1',right_on='key2') #基于不同的key
```

6.1.2 Merging on Index

```
1. pd.merge(df1,df2,left_on='key',right_index=True)
```

6.1.3 Concatenating Along on Axis

```
1. arr=np.arange(12).reshape((3,4))  
2. np.concatenate([arr,arr],axis=1) #横向组合  
3. pd.concat([s1, s2, s3], axis=1) #横向组合,axis=0纵向组合
```

6.1.4 Combining Data with Overlap

```
1. df1.combine_first(df2) #重复数据df1优先
```

6.2 Reshaping and Pivoting

6.2.1 stack/unstack

```
1. data=pd.DataFrame(  
2.     np.arange(6).reshape((2,3)),                                index=  
pd.Index(['Ohio','Colorado'],name='state'),                columns=pd.Index(['  
1','2','3'],name='number')  
3.     ) #构造一个表格  
4. result=data.stack() #解构表格为Series  
5. result.unstack() #返回解构前的DataFrame  
6. result.unstack(0)/result.unstack('state') #返回但行列互换
```

6.2.2 Pivot

```
1. pivotData=data.pivot('date','item','value')  
2. #date为行名(index), item为列名(key/column), value为pivot中各元素的值(value)
```

6.3 Data Transformation

6.3.1 Removing Duplicates

```
1. data.duplicated() #第一次出现False,再次出现True(保留第一次)  
2. data.drop_duplicates() #返回去除重复值的DF  
3.     data.drop_duplicates(['col1']) #基于col1列去除重复值  
4.     data.drop_duplicates(take_last=True) #保留最后一次
```

6.3.2 Mapping

```
1. data['xxx'].map(functionName)
```

6.3.3 Replacing Values

```
1. data.replace('sb','asshole')
```

6.3.4 Renaming Axis Indexes

```
1. data.index.map(str.upper) #行名全大写  
2. data.rename(index=str.title,columns=str.upper) #行名首字母大写, 列名全大写
```

- 后面略了一些内容

6.4 String Manipulation

6.4.1 String Object Methods

```
1. val.split(',')
2. val.strip()
3. "::-".join(pieces)
4. val.count(',')
```

6.4.2 Regular Expressions, 见python笔记

6.4.3 Vectorized string functions in pandas

```
1. data.str.contains('xxx') #类似data.find('xxx')
   data.str.findall(pattern)
2. data.str.match(pattern)
```

Chapter 7 Visualization

7.1 简单作图

7.1.1 事例

```
1. import matplotlib.pyplot as plt
2. fig=plt.figure() #create a new figure
3. ax1=fig.add_subplot(2,2,1) #2行2列, 在位置1添加坐标系
4. ax4=fig.add_subplot(2,2,4) #2行2列, 在位置4添加坐标系
5.
6. from numpy.random import randn
7. ax1.hist(randn(100),bins=20,color='k',alpha=0.3)
8. ax4.scatter(np.arange(30),np.arange(30)+3*randn(30))
```

7.1.2 subplot function

```
1. fig,axes=plt.subplots(2,3) #直接2*3图层, 6个坐标系
```

- other parameters:

```
1. nrows,
2. ncols,
```

```
3. sharex #same X-axis ticks
4. sharey #same Y-axis ticks
```

7.1.3 subplot_adjust

7.2 Colors, Markers and Line Styles

```
1. plt.plot(randn(30).cumsum(), color='k', linestyle='dashed', marker='o')
2. plt.plot(data, 'k-', drawstyle='steps-post', label='steps-post')
3. plt.legend(loc='best') #add a legend
```

7.3 Add a shape to a plot

```
1. rect=plt.Rectangle((0.2, 0.75), #起点坐标
2.     0.4, 0.15, #长度x,高度y
3.     color='k', alpha=0.3)
4. circ = plt.Circle((0.7, 0.2), #圆心坐标
5.     0.15, #短轴/长轴
6.     color='b', alpha=0.3)
7. pgon = plt.Polygon([[0.15, 0.15], [0.35, 0.4], [0.2, 0.6]], color='g',
8.     alpha=0.5) #三顶点坐标
9.
10. ax.add_patch(rect)
11. ax.add_patch(circ)
12. ax.add_patch(pgon)
```

7.4 Save to a file

```
1. plt.savefig('filename.png',
2.     dpi=400,
3.     facecolor/edgecolor,
4.     format='png'/'pdf'/'svg'...
5.     bbox_inches='tight' #the portion of the figure to save
6.     )
```

- P.S. There are also some plot functions in pandas

Chapter 8 Data Aggregation and Group Operations

8.1 GroupBy Mechanics

- split into subgroups-->apply function-->combine into group

8.2 GroupBy Example

```
1. df=pd.DataFrame({'key1':['a','a','b','b','a'],'key2':['1','2','1','2','1'], 'data1':np.random.randn(5), 'data2':np.random.randn(5)})
2.
3. means = df['data1'].groupby([df['key1'], df['key2']]).mean()
4.
5. means.unstack() #得到split但未进行apply function的数据
```

8.3 Iterating Over Groups

```
1. for name,group in df.groupby('key1'):
2.     print name
3.     print group
```

8.4 Aggregate function

```
1. quant=df.groupby('key1').quantile(0.9)
2. ag=df.groupby('key1').agg(yourOwnFunction)
```

8.5 Pivot Tables and Cross-Tabulation

```
1. tips.pivot_table('tip_pct',
2.     rows=['sex', 'smoker'],
3.     cols='day',
4.     aggfunc=len,
5.     margins=True)
```

- Cross-Tabulation:交叉分析（就是计数的透视表）

```
1. pd.crosstab(data.Gender,data.Handedness,margins=True)
2. #行为gender,列为handedness,数据为count
```

Chapter 9 Time Series

9.1 概念:

- 时间序列: 用时间排序的一组随机变量
- 特征:unstable, 波动幅度随时间变化

9.2 Date and Time Data Types and Tools

```
1. from datetime import datetime
2. now = datetime.now()
3. now.year/month/day
4.
5. delta=datetime(2011,1,7)-datetime(2008,6,24,8,15)
6. #返回天数差.days与秒数差.seconds
7.
8. starttime=datetime(2016,5,9)
9. start+timedelta(12) #+12 days
```

9.3 Converting between string and datetime

- datetime --> strings

```
1. stamp=datetime(2016,5,9)
2. str(stamp) #convert to default strings "2016-05-09 00:00:00"
3. stamp.strftime("%Y-%m-%d") #"2011-01-03"
```

- strings --> dates

```
1. datetime.strptime(strings,"%Y-%m-%d")
```

- parse:更改不同类型的日期格式

```
1. parse('6/12/2011',dayfirst=True) # (2011,12,6,0,0)
```

9.4 Time Series Basics

```
1. dates=[datetime(xxxx),datetime(xxxx),...]
2. timeSeries=pd.Series(np.random.randn(x),index=dates)
```

9.4.1 Indexing,Selection and Subsetting

9.4.2 Time Series with Duplicate Indices

9.5 Date Ranges,Frequencies and Shifting

Chapter 10 Financial and Economic Data Applications

Chapter 11 Advanced Numpy