

1:print()

```
1 1.print(x,y,sep=,end=)
2 sep='\n' new line
3 sep='/' same line and sep with / symbol
```

In [4]:

```
1 sx='hello';sy='python';sz=10
2 print(sx,sy)
3 print(sx+sy) # + oly work for str
4 print(sx,sy,sz)
5 print(sx,sy,sz,sep='\n') # break to 3 lines
6
```

```
hello python
hellopython
hello python 10
hello
python
10
```

In [5]:

```
1 print(sx,sy,sz,sep='//',end=';')
```

```
hello//python//10;
```

In [7]:

```
1 #example:repeat str
2 print((sx,sy)*3)
```

```
('hello', 'python', 'hello', 'python', 'hello', 'python')
```

```
In [11]: 1 """example 3: %s, list
2 1.%s =str
3 2.%d =digit
4 """
5 name='john'; age=23
6 #print("%s is %d year old."%(age,name)) #error 順序要跟
7 print("%s is %d year old."%(name,age))
```

john is 23 year old.

```
In [18]: 1 data=("apple", 'bana', 55.43)
2 print(type(data))
3 print("%s are %s fruits,priced %.2f"%data) #%.2f:float to 2 decimal point
```

<class 'tuple'>
apple are bana fruits,priced 55.43

```
1 list=[] #mutable
2 tuple=() #immutable
3 dict={}
```

```
In [20]: 1 """print list"""
2 myList=[10,20,30] #print myList[1]
3 print(' apple pie costs %d.' %myList[1])
```

apple pie costs 20.

2.input()

```
1 <data_type>(input())
```

```
In [21]: 1 bonus=int(input("input the bonus: "))
2 salary=100+bonus
3 print("salary is ", salary,end=';')
```

```
input the bonus: 50
salary is 150;
```

3. break many lines

```
1 slash \
2 blackslash /
```

```
In [22]: 1 """ex1 : break in element"""
2 a=10+20\
3     +30\
4     +40
5 print(a)
```

```
100
```

```
In [26]: 1 """ex2:break in result"""
2 a=1000
3 s="it is too many lines, they are xxxxxxx"
4 print(s,
5     a,
6     sep='++++++',
7     end=';')
```

```
it is too many lines, they are xxxxxxx+++++1000;
```

4 math symbols

```
1 1: divide /
2 2: remainder %
3 3: exponent **; a=a**2, square
4 4: not equal to !=
```

```
5 | 5.used in loc, &(and),(|)or , can't type and or directly
```

In [28]:

```
1 a=10;b=3
2 print(a%b) #10/3=3....1
3 print(a**b) #10*10*10
```

```
1
1000
```

5.string

In [32]:

```
1 """ex1:simple access"""
2 str='python'
3 num=105
4 print(str[1])
5 print(str*3)
6 print(str,"is a ",num,)
```

```
y
pythonpythonpython
python is a  105
```

In []:

```
1 """DML
2 1.string[start:end:step] #substring
3 2.s.replace('a', '')/s.strip('<remove thing>')
4 3.String.find()=java contains(),find() doesn't find a match, it returns -1
5   配合 if else
6 4. remove string space
7 a)x.replace(' ','')
8 b)strip(), 也有lstrip(),rstrip()
9 """
```

In [39]:

```
1 weight='150lbs'
2 w1=weight[0:2]
3 w2=weight.replace('lbs','')
4
5 print(w1) # return 2-0=2 chars,
6 print(w2)
7
8 if weight.find('lbs')!==-1:
9     print("find the lbs char")
10 else:
11     print("not find")
```

```
15
150
find the lbs char
```

In [5]:

```
1 #remove space
2 str=" hello python "
3 str1=str.replace(' ','')
4 print(str1) #replace all space
5
6 str2=str.strip() #remove l/r only
7 print(str2)
8
9 str3=str.lstrip()
10 print(str3)
```

```
hellopython
hello python
hello python
```

6 list

```
1 list vs array
2 1.list contain diff datatype, array must a same datatype
3 both can store non-unique elements
4
5 2. array is effient than list, + use numpy
```

```
1 1.[] start from 0, last is -1
2 2.access: list[start:end:step]
3 3.list.count(x) #num of list elements
4 4.list comprehension
5 #Input processed list to another list in a hight-efficiently way
6 link:
7 https://www.learncodewithmike.com/2020/01/python-comprehension.html
```

In [15]:

```
1 list=['A','B',3,'ee',9,'A',[1,2,3]]
2 print(list[-1])
3 print(list[1::2]) #even number
4 print(list.count('A'))
```

```
[1, 2, 3]
['B', 'ee', 'A']
2
```

In [18]:

```
1 """non-list comprehension"""
2 fruits = ["apple", "banana", "cherry", "kiwi", "mango"]
3 newlist = []
4 for x in fruits:
5     if "a" in x:
6         newlist.append(x)
7
8 print(newlist)
9
10 #ex2:
11 numbers = []
12 for x in range(10):
13     numbers.append(x * 3)
14 print(numbers)
```

```
['apple', 'banana', 'mango']
[0, 3, 6, 9, 12, 15, 18, 21, 24, 27]
```

In [19]:

```
1 newlist2 = [x for x in fruits if "a" in x] #can add if in
2 print(newlist2)
3
4 #ex2:
5 numbers = [x * 3 for x in range(10)]
6 print(numbers)
```

```
['apple', 'banana', 'mango']
[0, 3, 6, 9, 12, 15, 18, 21, 24, 27]
```

In [20]:

```
1 """ex3"""
2 numbers = [x * 3 for x in range(10) if x > 4]
3 print(numbers)
```

```
[15, 18, 21, 24, 27]
```

In [24]:

```
1 """ex4"""
2 numbers = [50, 2, 12, 30, 27, 4]
3 result = [number * 2 for number in numbers]
4 print(result)
5
6 """ex5"""
7 numbers2 = [50, 2, 12, 30, 27, 4]
8 result2 = [number for number in numbers2 if number > 10]
9 print(result2)
10
11 """ex6"""
12 numbers3 = [50, 2, 12, 30, 27, 4]
13 result3 = [number if number > 10 else 0 for number in numbers3]
14 #put if to front is to change the list num first
15 print(result3)
16
```

```
[100, 4, 24, 60, 54, 8]
[50, 12, 30, 27]
[50, 0, 12, 30, 27, 0]
```

7: range+random

```
In [32]: 1 import random as rd  
2 import numpy as np #np also has random method
```

```
In [ ]: 1
```

```
In [35]: 1 r1=rd.randrange(2, 20, 2) #return 1 random than are odd number  
2 print(r1)  
3  
4 r2=rd.randint(0, 9) #randin=random integer  
5 print(r2)
```

```
14  
8
```

```
In [39]: 1 """generate a list of sample  
2 1.for loop + above 2 method  
3 2.list comprehension  
4 3.random.sample() function  
5 """  
6 r1=[]  
7 for i in range(0,5): #gen 5 num  
8     n=rd.randint(1,100) #value range 1-100  
9     #can change to randrange  
10    r1.append(n)  
11 print(r1)  
12  
13 #in fact just change 1 to list comprehension  
14 r2=[rd.randrange(1,100,2) for i in range(0,5,1) ] #1-100, odd number  
15 print(r2)  
16  
17 #Generate 5 random numbers between 10 and 30  
18 r3 = rd.sample(range(10, 30), 5)  
19 print(r3)
```

```
[50, 24, 54, 7, 76]  
[15, 49, 39, 19, 59]  
[29, 14, 16, 21, 26]
```

chapter 8: 2-d for loop

In [5]:

```
1 """8.1 understand nested list/array first
2 """
3 #ex1:in theory 多·實際少用,single list
4 list=[[1,2,3,4],['a','b','c','d']]
5 print(list[0][1],list[1][2])
6 print()
7
8 for i in range(len(list)):
9     for j in range(len(list[i])):
10         print(list[i][j],sep=' ')
11         print()
```

2 c

1
2
3
4a
b
c
d

In [16]:

```
1 """
2 ex2: 常見多數系matrix, one to many
3 process: 1a,1b,1c,1d [0][0],[0][1],[0][2],[0][3],[0][4]
4         2a,2b,2c,2d
5 """
6 color=['red','green','blue','yello']
7 fruit=['apple','bal','cherry','nut']
8 for i in color:
9     for j in fruit:
10        print(i+"+"+j)
11 print()
12
13 for i in range(0,len(color),1):
14    print(color[i],fruit[i],sep='---')
15 print()
```

```
red+apple
red+bal
red+cherry
red+nut
green+apple
green+bal
green+cherry
green+nut
blue+apple
blue+bal
blue+cherry
blue+nut
yello+apple
yello+bal
yello+cherry
yello+nut

red---apple
green---bal
blue---cherry
yello---nut
```

In [23]:

```
1 """ex3:最實用 · combine lists into one
2     zip(list_x,list_y,z,...)
3 """
4 name=['peter', 'john', 'amy']
5 age=[23,12,43]
6 score=[90,88,79]
7 info=list(zip(name,age,score))
8 info
9 #list() is a built-in funciton, don't use list=[1,2,4]
10 #del list    for error: List object is not callable
```

Out[23]: [('peter', 23, 90), ('john', 12, 88), ('amy', 43, 79)]

chapter 9: zip()

In [24]:

```
1 """ex1:dic +zip()
2 change list to dict(key:value)
3 """
4
5 city = ["tw", "us", "hk"]
6 salary = [2.1, 3.8, 6.9]
7 d1=dict(zip(city,salary))
8 d1
```

Out[24]: {'tw': 2.1, 'us': 3.8, 'hk': 6.9}

In [25]:

```

1 """ex2:for loop
2 """
3 names = ["A", "B", "C"]
4 level = [11, 23, 46]
5 ages = [45, 67, 82]
6 for x, y, z in zip(names, level, ages):
7     print(x,y,z)

```

A 11 45
 B 23 67
 C 46 82

In [27]:

```

1 """ex3:gen list"""
2 data = list(zip(range(3), 'ABCD')) #must in pair
3 data

```

Out[27]: [(0, 'A'), (1, 'B'), (2, 'C')]

chapter 10: set

```

1 #note: set values will auto distinct
2 1.a - b 集合a中包含而集合b中不包含的元素
3 2.a | b element in set a or b
4 3.a & b intersection, must in a and b
5 4. a ^ b 不同时包含于a和b的元素, 只出现在set a or b
6 5.'orange' in basket # return True/False
7 #####
8 s.add( x ) # add elements
9 s.update( x )
10 s.remove( x )

```

In [5]:

```
1 basket = {'apple', 'orange', 'apple', 'pear', 'orange', 'banana'}
2 print(basket) #auto remove duplicate value
3 print("banana" in basket)
4 #print(basket[1]) # will error
```

```
{'pear', 'orange', 'banana', 'apple'}
True
```

In [6]:

```
1 """print set() value """
2 basket = {'apple', 'orange', 'apple', 'pear', 'orange', 'banana'}
3 # for each Loop. print all
4 for i in basket:
5     print(i)
```

```
pear
orange
banana
apple
```

In [10]:

```
1 #return specific value
2 for i in basket:
3     if i.find('b')>=0:
4         print(i)
```

```
banana
```

In [11]:

```
1 a = set('abracadabra')
2 print(a)
```

```
{'c', 'd', 'b', 'a', 'r'}
```

In [12]:

```

1 """list comprehension
2 for loop + if
3 """
4 list presentation
5
6 a1 = {x for x in 'abracadabra' if x not in 'abc'}
7 a1

```

Out[12]: {'d', 'r'}

In [14]:

```

1 """must use symbol but not text like ''
2 """
3 a = set('abracadabra')
4 b = set('alacazam')
5 print("a-b is:", a-b) # value in a but no in b
6 print("a^b is:", a^b) # data can only in a or b, can not both
7 print("a|b is :", a|b) # union in a or b
8 print("a&b is : ", a&b) #intersection, must in both a and b

```

```

a-b is: {'d', 'r', 'b'}
a^b is: {'r', 'b', 'l', 'd', 'm', 'z'}
a|b is : {'l', 'c', 'd', 'b', 'm', 'z', 'a', 'r'}
a&b is : {'a', 'c'}

```

chapter 11 :read files

```

1 the chapters of basic knowledge have gone. Start to 實戰。
2 1.data files: txt xlsx csv
3 2.usually we need to change the txt, csv to xlsx.
4 ~~~~~
5 3.csv=microsoft excel comma seperated values file
6     xlsx=microsoft excel worksheet
7     txt=text document

```

In [32]:

```

1 import pandas as pd
2 # inplace =True meas change to exiting df, default is False

```

In [15]:

```

1 """use pandas module's methods, but not function,
2     so it is pd.to_read(var), but to_read(var)
3 #####
4 1.
5 given path:C:\Users\fengs\Desktop\py_project\csv_file
6 not use \, but /
7 2.or add (r'')
8 """
9 #df_p=pd.read_csv('C:\Users\fengs\Desktop\py_project\csv_file\pokemon_data.csv') #unicode error
10 #df_p=pd.read_csv('C:/Users/fengs/Desktop/py_project/Upy_project/csv_file/pokemon_data.csv')

```

File "C:\Users\fengs\AppData\Local\Temp\ipykernel_3956/338393444.py", line 8

"""

^

SyntaxError: (unicode error) 'unicodeescape' codec can't decode bytes in position 117-118: truncated \UXXXXXXXXX escape

chapter 11.1 foundational info of DF

In [3]:

```

1 df_csv=pd.read_csv(r'C:\Users\fengs\Desktop\py_project\csv_file\pokemon_data.csv')
2 df_csv.head()

```

Out[3]:

#	Name	Type 1	Type 2	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
0 1	Bulbasaur	Grass	Poison	45	49	49	65	65	45	1	False
1 2	Ivysaur	Grass	Poison	60	62	63	80	80	60	1	False
2 3	Venusaur	Grass	Poison	80	82	83	100	100	80	1	False
3 3	VenusaurMega Venusaur	Grass	Poison	80	100	123	122	120	80	1	False
4 4	Charmander	Fire	NaN	39	52	43	60	50	65	1	False

In [63]:

```
1 #find df's row &columns
2 #(number of rows, number of columns).
3 print(df_csv.shape)
4 print(df_csv.shape[0],df_csv.shape[1])
5 print(df_csv.count(0))
```

```
(800, 12)
800 12
#          800
Name        800
Type 1      800
Type 2      414
HP          800
Attack       800
Defense      800
Sp. Atk      800
Sp. Def      800
Speed         800
Generation    800
Legendary     800
dtype: int64
```

In [19]:

```
1 df_csv.columns #it is an attribute but not a function or method
```

```
Out[19]: Index(['#', 'Name', 'Type 1', 'Type 2', 'HP', 'Attack', 'Defense', 'Sp. Atk',
   'Sp. Def', 'Speed', 'Generation', 'Legendary'],
              dtype='object')
```

In [21]:

```
1 #df.info()=show table DDL
2 df_csv.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 800 entries, 0 to 799
Data columns (total 12 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   #           800 non-null    int64  
 1   Name         800 non-null    object  
 2   Type_1       800 non-null    object  
 3   Type_2       414 non-null    object  
 4   HP           800 non-null    int64  
 5   Attack        800 non-null    int64  
 6   Defense       800 non-null    int64  
 7   Sp. Atk       800 non-null    int64  
 8   Sp. Def       800 non-null    int64  
 9   Speed          800 non-null    int64  
 10  Generation     800 non-null    int64  
 11  Legendary      800 non-null    bool   
dtypes: bool(1), int64(8), object(3)
memory usage: 69.7+ KB
```

In [22]:

```
1 #df_describe()=show stat data
2 df_csv.describe()
```

Out[22]:

	#	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation
count	800.000000	800.000000	800.000000	800.000000	800.000000	800.000000	800.000000	800.000000
mean	362.813750	69.258750	79.001250	73.842500	72.820000	71.902500	68.277500	3.32375
std	208.343798	25.534669	32.457366	31.183501	32.722294	27.828916	29.060474	1.66129
min	1.000000	1.000000	5.000000	5.000000	10.000000	20.000000	5.000000	1.00000
25%	184.750000	50.000000	55.000000	50.000000	49.750000	50.000000	45.000000	2.00000
50%	364.500000	65.000000	75.000000	70.000000	65.000000	70.000000	65.000000	3.00000
75%	539.250000	80.000000	100.000000	90.000000	95.000000	90.000000	90.000000	5.00000
max	721.000000	255.000000	190.000000	230.000000	194.000000	230.000000	180.000000	6.00000

In [20]:

```
1 print(df_csv.head(2))
2 print(df_csv.tail(2))
```

	#	Name	Type 1	Type 2	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	\
0	1	Bulbasaur	Grass	Poison	45	49	49	65	65	45	
1	2	Ivysaur	Grass	Poison	60	62	63	80	80	60	
		Generation	Legendary								
0		1	False								
1		1	False								
	#	Name	Type 1	Type 2	HP	Attack	Defense	Sp. Atk	\		
798	720	Hoopa	Hoopa	Unbound	Psychic	Dark	80	160	60	170	
799	721		Volcanion		Fire	Water	80	110	120	130	
	Sp. Def	Speed	Generation	Legendary							
798	130	80	6	True							
799	90	70	6	True							

Chapter 11.2 query DF

In [29]:

```
1 """select cols,
2 1.df[] only for single col
3 2.df[[x,y]] for multi-col, as it only return one arg
4 """
5 print(df_csv[['#', 'Name', 'Type 1']].head(2))
6 print(df_csv[['#', 'Name', 'HP']].tail(2))
```

```
#      Name Type 1
0 1 Bulbasaur Grass
1 2 Ivysaur  Grass
#
#          Name   HP
798 720 HoopaHoopa Unbound  80
799 721 Volcanion  80
```

In [37]:

```
1 #between 2 values
2 df_csv.loc[df_csv['#'].between(99, 101)]
```

Out[37]:

	#	Name	Type 1	Type 2	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary	
	107	99	Kingler	Water	NaN	55	130	115	50	50	75	1	False
	108	100	Voltorb	Electric	NaN	40	30	50	55	55	100	1	False
	109	101	Electrode	Electric	NaN	60	50	70	80	80	140	1	False

In [4]:

```

1 """loc[row,cols]
2 loc[].range(start,end,step)
3 most practical function of DF
4 2.loc=location
5 iloc=index location
6 """
7 r1=df_csv.loc[range(0,len(df_csv),2),['#','Name','Sp. Atk']].head(3) #even 雙·odd 單
8 r1

```

Out[4]:

	#	Name	Sp. Atk
0	1	Bulbasaur	65
2	3	Venusaur	100
4	4	Charmander	60

```

1 ref:
2 df3 = df.loc[(df['Date'] > 'Feb 06, 2019') & (df['Open'] > 62), ['Date', 'Open']]

```

In [5]:

```

1 """loc + where condition
2 df.loc[(c1)&(c2)&(c3)...,['col1','col2']]
3 1.each condition needs to be bracketed
4 2.use []=='string', but not =, '=' means assign to
5
6 """
7 r2=df_csv.loc[(df_csv['#'].between(3,30)) & (df_csv['Type 1']=='Fire'), ['#','Name','Type 1']]
8 r2.head(3)

```

Out[5]:

	#	Name	Type 1
4	4	Charmander	Fire
5	5	Charmeleon	Fire
6	6	Charizard	Fire

```
In [6]: 1 r3=df_csv.loc[((df_csv['#'].between(3,30)) & (df_csv['Type 1']=='Fire'))|(df_csv['Speed']>150),
2                               ['#','Name','Type 1','Speed'])
3 r3
```

Out[6]:

	#	Name	Type 1	Speed
4	4	Charmander	Fire	65
5	5	Charmeleon	Fire	80
6	6	Charizard	Fire	100
7	6	CharizardMega Charizard X	Fire	100
8	6	CharizardMega Charizard Y	Fire	100
315	291	Ninjask	Bug	160
431	386	DeoxysSpeed Forme	Psychic	180

```
In [7]: 1 """more advaced query
2 1.df[].str.contains('A|B') # like
3 2.~df[].str.contains('A|B')
4 3.join
5 """
6 r4=df_csv.loc[(df_csv['#'].between(20,120)) & (df_csv['Attack']>=100) & (df_csv['Name'].str.contains('s')),
7                               ['Name','#','Attack','Type 1']]
8 r4
```

Out[7]:

	Name	#	Attack	Type 1
33	Sandslash	28	100	Ground
84	Rapidash	78	100	Fire
124	KangaskhanMega Kangaskhan	115	125	Normal

```
In [ ]: 1 #add new columns base on conditon
2 df['Is_eligible'] = np.where(df['Age'] >= 18, True, False)
```

chapter 11.3 sorting DF

```
1 """sorting
2 DataFrame.sort_values(by, axis=0, ascending=True, inplace=False, kind='quicksort',
3                      na_position='last', ignore_index=False, key=None)[source]
4 """
5 ref:df.sort_values(by=['Brand'], inplace=True, ascending=False)
```

In []:

```
1 #example: sort many cols
2 df.sort(['a', 'b'], ascending=[True, False]) #or use 0,1
```

In [8]:

```
1 r5=df_csv.loc[df_csv['#'].between(50,100),['#','Name','Type 1','HP','Attack']]
2 r5.sort_values(by=['#'],inplace=True,ascending=False)
3 r5.head(5)
```

Out[8]:

#	Name	Type 1	HP	Attack
108	Volterb	Electric	40	30
107	Kingler	Water	55	130
106	Krabby	Water	30	105
105	Hypno	Psychic	85	73
104	Drowzee	Psychic	60	48

In [9]:

```
1 sales_data = pd.DataFrame({"name": ["William", "Emma", "Sofia", "Markus", "Edward", "Thomas", "Ethan", "Olivia", "Arun", "Anika"],  
2 , "region": ["East", "North", "East", "South", "West", "West", "South", "West", "East", "South"],  
3 , "sales": [50000, 52000, 90000, 34000, 42000, 72000, 49000, 55000, 67000, 65000, 67000],  
4 , "expenses": [42000, 43000, 50000, 44000, 38000, 39000, 42000, 60000, 39000, 44000, 45000]}))  
5 sales_data
```

Out[9]:

	name	region	sales	expenses
0	William	East	50000	42000
1	Emma	North	52000	43000
2	Sofia	East	90000	50000
3	Markus	South	34000	44000
4	Edward	West	42000	38000
5	Thomas	West	72000	39000
6	Ethan	South	49000	42000
7	Olivia	West	55000	60000
8	Arun	West	67000	39000
9	Anika	East	65000	44000
10	Paulo	South	67000	45000

In [10]:

```
1 sales_data.sort_values(by=['sales', 'expenses'], inplace=True, ascending=[False, True])  
2 sales_data
```

Out[10]:

	name	region	sales	expenses
2	Sofia	East	90000	50000
5	Thomas	West	72000	39000
8	Arun	West	67000	39000
10	Paulo	South	67000	45000
9	Anika	East	65000	44000
7	Olivia	West	55000	60000
1	Emma	North	52000	43000
0	William	East	50000	42000
6	Ethan	South	49000	42000
4	Edward	West	42000	38000
3	Markus	South	34000	44000

11.5 custom ordering

In [12]:

```
1 #try to sorting by type 1
2 r6=df_csv.loc[df_csv['#'].between(25,100),['Name','#','HP','Type 1']]
3 r6.head()
```

Out[12]:

	Name	#	HP	Type 1
30	Pikachu	25	35	Electric
31	Raichu	26	60	Electric
32	Sandshrew	27	50	Ground
33	Sandslash	28	75	Ground
34	Nidoran (Female)	29	55	Poison

In [14]:

```
1 #find distinct value of a column,
2 #sort by this custom cols value
3 df_csv['Type 1'].unique()
```

Out[14]: array(['Grass', 'Fire', 'Water', 'Bug', 'Normal', 'Poison', 'Electric',
 'Ground', 'Fairy', 'Fighting', 'Psychic', 'Rock', 'Ghost', 'Ice',
 'Dragon', 'Dark', 'Steel', 'Flying'], dtype=object)

```
1 pandas not equal to list of values
2 df[~df['column_name'].isin(['value_1', 'value_2'])]
3
4 #failed
5 r6['Type 1']!='Water','Normal','Grass')
```

In [58]:

```

1 """1.map a new cols
2 Use CategoricalDtype Cast data to an ordered category type
3 As commented, in newer pandas,
4 Series has a replace method to do this more elegantly:
5 s = df['m'].replace({'March':0, 'April':1, 'Dec':3})
6
7 link: df.map() vs df.apply()
8 https://zhuanlan.zhihu.com/p/100064394
9 https://chowdera.com/2020/11/20201106205416849i.html
10 """
11 main=['Water', 'Normal', 'Grass']
12 #link:https://www.statology.org/case-statement-pandas/
13 r6['order']=np.where(r6['Type 1']=='Water',1,
14                      np.where(r6['Type 1']=='Normal',2,
15                                np.where(r6['Type 1']=='Grass',3,
16                                         4 #mean not in 3 above
17                                         ))))
18 r6

```

Out[58]:

	Name	#	HP	Type 1	order
30	Pikachu	25	35	Electric	4
31	Raichu	26	60	Electric	4
32	Sandshrew	27	50	Ground	4
33	Sandslash	28	75	Ground	4
34	Nidoran (Female)	29	55	Poison	4
...
104	Drowzee	96	60	Psychic	4
105	Hypno	97	85	Psychic	4
106	Krabby	98	30	Water	1
107	Kingler	99	55	Water	1
108	Voltorb	100	40	Electric	4

79 rows × 5 columns

In [59]:

```
1 r6.sort_values(by=['order', 'Name'], inplace=True, ascending=[True, True])
2 r6
```

Out[59]:

	Name	#	HP	Type 1	order
98	Cloyster	91	50	Water	1
94	Dewgong	87	90	Water	1
60	Golduck	55	80	Water	1
107	Kingler	99	55	Water	1
106	Krabby	98	30	Water	1
...
54	Venomoth	49	70	Bug	4
53	Venonat	48	60	Bug	4
108	Voltorb	100	40	Electric	4
42	Vulpix	37	38	Fire	4
46	Zubat	41	40	Poison	4

79 rows × 5 columns

In [69]:

```
1 #####show all columns except one
2 #r6.Loc[:,r6.columns != 'order']
3
4 #not show some columbs
5 #df.Loc[:, ~df.columns.isin(['col1', 'col2'])]
6 r6.loc[:,~r6.columns.isin(['order', 'HP'])].head()
```

Out[69]:

	Name	#	Type 1
98	Cloyster	91	Water
94	Dewgong	87	Water
60	Golduck	55	Water
107	Kingler	99	Water
106	Krabby	98	Water

11.5.2 :create a ordering function

```
1 if row A == B: 0
2
3 if rowA > B: 1
4
5 if row A < B: -1
```

In [68]:

```
1 #create lists to DF
2 df=pd.DataFrame({
3     'A':[2,3,1],
4     'B':[2,1,3]
5 })
6 df
```

Out[68]:

	A	B
0	2	2
1	3	1
2	1	3

In [53]:

```
1 def f(row):
2     if row['A'] == row['B']:
3         val = 0
4     elif row['A'] > row['B']:
5         val = 1
6     else:
7         val = -1
8     return val
9 df['C'] = df.apply(f, axis=1) #axis=1 AS diff ROW compare
10 df
```

Out[53]:

	A	B	C
0	2	2	0
1	3	1	1
2	1	3	-1

In [23]:

```
1 import numpy as np
```

In [43]:

```

1 r7=df_csv.loc[(df_csv['#'].between(1,10)) & (~df_csv['Type 1'].isin(['Grass','Water']))]
2 r7

```

Out[43]:

#		Name	Type 1	Type 2	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed	Generation	Legendary
4	4	Charmander	Fire	NaN	39	52	43	60	50	65	1	False
5	5	Charmeleon	Fire	NaN	58	64	58	80	65	80	1	False
6	6	Charizard	Fire	Flying	78	84	78	109	85	100	1	False
7	6	CharizardMega Charizard X	Fire	Dragon	78	130	111	130	85	100	1	False
8	6	CharizardMega Charizard Y	Fire	Flying	78	104	78	159	115	100	1	False
13	10	Caterpie	Bug	NaN	45	30	35	20	20	45	1	False

chapter 13.numpy

chapter 13.1 : background

```

1 numpy use n-dimension array to replace python list,MAIN concepts include{ ndim,shape,dtype}
2 1234 :shape(4,); axis=0
3
4 1234
5 5678
6 shape(2,4); axis has 0,1 0:y-aixs(cols) 1:x-aix(rows)

```

```

1 Q1:why use numpy???
2 In Python we have lists that serve the purpose of arrays, but they are slow to process.
3 NumPy aims to provide an array object that is up to 50x faster than traditional Python lists.
4 The array object in NumPy is called ndarray, it provides a lot of supporting functions that make working with
ndarray very easy.
5 Arrays are very frequently used in data science, where speed and resources are very important.
6
7 Q2:Why is NumPy Faster Than Lists?

```

- 8 NumPy arrays are stored at one continuous place in memory unlike lists, so processes can access and manipulate them very efficiently. This behavior is called locality of reference in computer science.
- 9 This is the main reason why NumPy is faster than lists. Also it is optimized to work with latest CPU architectures.
- 10

Chapter 13.2 basic syntax

In [3]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import os
5 import sys
```

In []:

```
1 numpy.array(object, dtype=None, copy=True, order='K', subok=False, ndmin=0)
2
3 link:
4 https://thispointer.com/python-numpy-create-a-numpy-array-from-list-tuple-or-list-of-lists-using-numpy-array/
```

In [5]:

```
1 #1.create list/tuple array
2 np1=np.array([1,2,3]) #list[]
3 np2=np.array((4,5,6)) #tuple(), immutable, one datatype
4 print(np1)
5 print(np2)
6 print(type(np1),type(np2)) # bot convert to ndarray
```

```
[1 2 3]
[4 5 6]
<class 'numpy.ndarray'> <class 'numpy.ndarray'>
```

In [8]:

```
1 a1=np.array([1,2,3],dtype=int)
2 a2=np.array([1,2,3],dtype=float)
3 print(a1)
4 print(a2)
```

```
[1 2 3]
[1. 2. 3.]
```

In [7]:

```
1 #2.arrang(start,stop,step,dtype=None)
2 """similar list range,range is built-in python class
3      arange is a function of third-party library(Numpy)
4 """
5 print(np.arange(1, 8, 3))
6 print( np.arange(1, 10.1, 3)) #get 4f, include 10.1
```

```
[1 4 7]
[ 1.  4.  7. 10.]
```

In [11]:

```
1 """
2 numpy.linspace(start, stop, num=50,
3                 endpoint=True, retstep=False,
4                 dtype=None, axis=0)
5 -> start : [optional] start of interval range. By default start = 0
6 -> stop   : end of interval range
7 -> retstep : If True, return (samples, step). By deflut retstep = False
8 -> num    : [int, optional] No. of samples to generate
9 -> dtype   : type of output array
10 """
11 x = np.linspace(0, 2, 5) #5 is return 5 num
12 print(x) #default return float
13
14 y=np.linspace(0,8,2) #return 2 nums
15 print(y)
```

```
[0.  0.5 1.  1.5 2. ]
[0. 8.]
```

In [23]:

```
1 #3.create nd-array
2 """
3 arr1 = np.array([1, 2, 3, 4, 5]) #1d
4 arr2 = np.array([[1, 2, 3], [4, 5, 6]]) #2d
5
6 arr3 = np.array([[1, 2, 3], [4, 5, 6]],
7                 [[1, 2, 3], [4, 5, 6]])
8
9 """
10 print('arr3 is :',arr3)
11 print()
12
13 #error:if [1,2,3,4]
14 #3*4 matrix
15 arr4=np.array([[1,2,4,3],['d','e','f','G'],['ten','ele','twe','13th']])
16 print('arr4 is:',arr4)
17 print("dimentsion: ",arr4.ndim) #as not more than 2 layer nested array
18 print("matrix/(row: col) is :",arr4.shape)
19 print("total number of ele is:",arr4.size)
```

```
arr3 is : [[[1 2 3]
 [4 5 6]]]
```

```
[[[1 2 3]
 [4 5 6]]]
```

```
arr4 is: [['1' '2' '4' '3']
 ['d' 'e' 'f' 'G']
 ['ten' 'ele' 'twe' '13th']]
dimentsion: 2
matrix/(row: col) is : (3, 4)
total number of ele is: 12
```

In [30]:

```
1 """4.change to multi-list to array
2 numpy.reshape(a, newshape, order='C')
3 order:'C','K','A'
4 """
5 a1=np.arange(8) #(start,end,step)
6 print(a1)
7 a1=np.arange(1,17).reshape(4,4) #(row,col)
8 print(a1)
```

```
[0 1 2 3 4 5 6 7]
[[ 1  2  3  4]
 [ 5  6  7  8]
 [ 9 10 11 12]
 [13 14 15 16]]
```

In [37]:

```
1 #5 access array values
2 #a1[row,col]=[<start:end:step:>,<:::>] start from 0
3 print(a1[1,1:3]) #1st row, 1,2 ele
4 print(a1[1:3,2]) # 1&2 row, 2nd col
5 print(a1[:,::2])
6 print(a1[:,2])
7 print(a1[1,:])
```

```
[6 7]
[ 7 11]
[[ 1  3]
 [ 9 11]]
[ 3  7 11 15]
[5 6 7 8]
```

chapter 13.3 initialize diff types of arrays

In [43]:

```
1 #1.all 1s/0s matrix
2 #numpy.ones(shape, dtype = None, order = 'C')
3 #shape =(row,col)
4 print(np.zeros(3,))
```

```
5 print(np.zeros((2,3)))
6 print(np.ones((3,3)))
7
```

```
[0. 0. 0.]
[[0. 0. 0.]
 [0. 0. 0.]]
[[1. 1. 1.]
 [1. 1. 1.]
 [1. 1. 1.]]
```

In [44]:

```
1 #2.any other number martrix
2 #numpy.full(shape, fill_value, dtype = None, order = 'C')
3 np.full((2,3),12)
```

Out[44]: array([[12, 12, 12],
 [12, 12, 12]])

In [45]:

```
1 #full_like, change copy a matrix shape
2 arr2 = np.arange(16, dtype=float).reshape(4, 4)
3 print("\n\narr2 before full_like : \n", arr2)
4
5 # using full_like
6 print("\narr2 after full_like : \n", np.full_like(arr2, -3))
```

```
arr2 before full_like :
[[ 0.  1.  2.  3.]
 [ 4.  5.  6.  7.]
 [ 8.  9. 10. 11.]
 [12. 13. 14. 15.]]
```

```
arr2 after full_like :
[[-3. -3. -3. -3.]
 [-3. -3. -3. -3.]
 [-3. -3. -3. -3.]
 [-3. -3. -3. -3.]]
```

chapter 14 numpy generate random data

In [47]:

```
1 """
2 1.random decimal number
3 np.random.rand(d0,d2,d3) :d means dimension
4 #value[0,1]
5 """
6 print(np.random.rand(2))
7 print(np.random.rand(2,3)) #(row,cols)
```

```
[0.18144125 0.18824898]
[[0.55449679 0.15667761 0.77931262]
 [0.7523837  0.38728601 0.72929654]]
```

In [51]:

```
1 """
2 np.random.randint(min,max,size,dtype)
3 random.randint(low, high=None, size=None, dtype=int)
4 #size(row,col) d-array
5 """
6 print("1st arr: \n",
7      np.random.randint(7,size=(4))) #0-7,7 exclude
8 print("2nd arr: \n",
9      np.random.randint(low=2,high=8,size=(2,3))
10     )
```

```
1st arr:
[6 4 0 2]
2nd arr:
[[3 3 2]
 [5 2 3]]
```

In [58]:

```
1 """most practical use to gen list data
2 choice(a, size=None, replace=True, p=None)
3 NOTE:replace=False :not repeat value
4 a=array
5 """
6 print("1st arr: \n",
7      np.random.choice(range(2,10),size=(5),replace=False)
8      )
9 print("2nd arr: \n",
10     np.random.choice(range(3,50),size=(2,4),replace=False)
11     )
```

```
1st arr:
[7 9 8 3 6]
2nd arr:
[[13 33 24 40]
 [46 8 16 30]]
```

real case on plot

In [67]:

```

1 boolean=[True, False]
2 gender=["男", "女"]
3 color=["white", "black", "yellow"]
4 data=pd.DataFrame({
5     "height":np.random.randint(150,190,100),
6     "weight":np.random.randint(40,90,100),
7     "smoker":[boolean[x] for x in np.random.randint(0,2,100)],
8     "gender":[gender[x] for x in np.random.randint(0,2,100)],
9     "age":np.random.randint(15,90,100),
10    "color":[color[x] for x in np.random.randint(0,len(color),100) ]
11 }
12 )

```

In []:

```

1 """gender列的男替换为1 · 女替换为0 ·
2 Series.map() vs apply()
3 """
4 data["gender"] = data["gender"].map({"男":1, "女":0})
5 data["gender"] = data["gender"].map(gender_map)
6 #
7 data["age"] = data["age"].apply(apply_age,args=(-3,))

```

chapter 14.2 numpy math

given random arrays

In [59]:

```

1 a=np.random.choice(range(2,10),size=(5),replace=False)
2 b=a.copy()
3 b

```

Out[59]: array([4, 2, 7, 6, 9])

In [67]:

```
1 #1given
2 #random draw [1,10] to build 2*3
3 a=np.random.choice(range(1,10),size=(2,3),replace=False)
4 """
5 #[2,13] all element to build 2*3
6 #error:b=np.arange(2,13).reshape(2,3)
7 #need to arange=2*3=6
8 """
9
10 b=np.arange(5,11).reshape(2,3) #range total must=6
11
12 print("a array is:\n",a)
13 print("b array is:\n",b)
```

a array is:

```
[[8 4 3]
 [7 2 5]]
```

b array is:

```
[[ 5  6  7]
 [ 8  9 10]]
```

In [73]:

```
1 #1 add
2 c= a+1
3 print("c is :\n",c)
4 print()
5
6 #expoent 次方
7 d=a**2
8 print("d is : \n",d)
9 print()
10
11 #extract specific element
12 e=(a[0,:]+1) # fist row,all cols
13 print("e is :\n",e)
```

c is :
[[9 5 4]
 [8 3 6]]

d is :
[[64 16 9]
 [49 4 25]]

e is :
[9 5 4]

In [77]:

```
1 """Agg of arrays
2 """
3 #ex1:
4 print(a)
5 print()
6 f=np.sum(a)
7 print(f)
8 print()
9
10 #ex2:expod
11 #所有相乘
12 g=np.prod(a)
13 print("prod is :",g)
14
15 #min/max
16 #(0=y)min:each col
17 #(1=x)max : each row
18 h=[np.min(a, axis=0), np.max(a, axis=1)]
19 print("h is \n",h)
20
21 """others stat
22 1.cumsum 累計加
23 2.cumprod 累計乘
24 3.ptp :diff=max-min
25 4.percentile(arr,75) 第75% value
26 """
```

```
[[8 4 3]
 [7 2 5]]
```

29

```
prod is : 6720
h is
[array([7, 2, 3]), array([8, 7])]
```

14.3 sorting array

In [81]:

```
1 #randint has not replace agr
2 #elements may repeat
3 x=np.random.randint(2,15,size=(3,5))
4 x
```

Out[81]: array([[3, 13, 11, 9, 9],
 [7, 6, 10, 6, 4],
 [7, 3, 11, 9, 3]])

In [78]:

```
1
```

Out[78]: array([[10, 4, 13, 9],
 [6, 11, 3, 14],
 [2, 12, 8, 7]])

In [88]:

```
1 #replace=False(not repeat)
2 a=np.random.choice(range(2,15),size=(3,4),replace=False)
3 print("given arrays is:\n",a)
4 print()
5
6 #sort by row
7 b=np.sort(a,axis=1)
8 print("sort by row ascending:\n",b)
9 print()
10
11 b1=-np.sort(-a,axis=1)
12 print("sort by row descending:\n",b1)
13 print()
14
15 c=np.sort(a,axis=0)
16 print("sort by cols:\n", c)
```

given arrays is:

```
[[ 8  7 12  2]
 [14  5 11  9]
 [ 6 13 10  3]]
```

sort by row ascending:

```
[[ 2  7  8 12]
 [ 5  9 11 14]
 [ 3  6 10 13]]
```

sort by row descending:

```
[[12  8  7  2]
 [14 11  9  5]
 [13 10  6  3]]
```

sort by cols:

```
[[ 6  5 10  2]
 [ 8  7 11  3]
 [14 13 12  9]]
```

In []:

1

In []:

1

chapter 15. outport DF to excel

In []:

1

Part 2:

chapter 16: OS module

In [1]:

```
1 """get the jupyter files info:  
2 eg:path ,file names  
3 """  
4 import os
```

In [3]:

```
1 #cwd: current working directly  
2 os.getcwd()
```

Out[3]: 'C:\\\\Users\\\\fengs'

In [7]:

```
1 """
2 default:list cwd files
3 dir=directory
4 """
5 #os.listdir()
6 os.listdir('C:/Users/fengs/Desktop/pd_real')
```

Out[7]:

```
['data.xlsx',
 'forget',
 'gd',
 'gd_2',
 'gd_3',
 'GeeksforGeeks',
 'mathplot',
 'mathplot.pdf',
 'matplotlib_tutorial-master.zip',
 'Misc',
 'py_loc.pdf',
 'py_road.pdf',
 'README.md',
 'SalesAnalysis']
```

In [9]:

```
1 import glob
2 # All files and directories ending with .pdf and that don't begin with a dot:
3 print(glob.glob("C:/Users/fengs/Desktop/pd_real/*.pdf"))
```

```
['C:/Users/fengs/Desktop/pd_real\\mathplot.pdf', 'C:/Users/fengs/Desktop/pd_real\\py_loc.pdf', 'C:/Users/fengs/Desktop/pd_real\\py_road.pdf']
```

In [11]:

```
1 # All files and directories ending with .pdf with depth of 2 folders, ignoring names beginning with a dot:
2 print(glob.glob("C:/Users/fengs/Desktop/pd_real/*.*pdf"))
```

```
['C:/Users/fengs/Desktop/pd_real\\mathplot.pdf', 'C:/Users/fengs/Desktop/pd_real\\py_loc.pdf', 'C:/Users/fengs/Desktop/pd_real\\py_road.pdf']
```

In [16]:

```
1 #print paths of files as well
2 pdf_files = []
3 for file in glob.glob("C:/Users/fengs/Desktop/pd_real/*.pdf"):
4     #pdf_files.append(file)
5     print(file)
```

```
C:/Users/fengs/Desktop/pd_real\mathplot.pdf
C:/Users/fengs/Desktop/pd_real\py_loc.pdf
C:/Users/fengs/Desktop/pd_real\py_road.pdf
```

chapter 16.2 formal ways

In [17]:

```
1 path = "C:/Users/fengs/Desktop/pd_real"
2 dir_list = os.listdir(path)
3
4 print("Files and directories in '", path, "' :")
5
6 # prints all files
7 print(dir_list)
```

```
Files and directories in ' C:/Users/fengs/Desktop/pd_real ' :
['data.xlsx', 'forget', 'gd', 'gd_2', 'gd_3', 'GeeksforGeeks', 'mathplot', 'mathplot.pdf', 'matplotlib_tutorial-master.zip', 'Misc', 'py_loc.pdf', 'py_road.pdf', 'README.md', 'SalesAnalysis']
```

In [19]:

```
1 #print files names only
2 path='C:/Users/fengs/Desktop/pd_real/mathplot'
3 for x in os.listdir(path):
4     if x.endswith(".csv"):
5         # Prints only csv file present in My Folder
6         print(x)
```

```
fifa_data.csv
gas_prices.csv
```

chapter 17:pandas plot

```
1 there are 3 commonly used plots:  
2 1.pandas plot :simple,easy to use  
3 2.mathplotlib plot: complicated, but more flexible  
4 3.seaborn plot :optional to know
```

```
1 Although they may appear similar, these modules have unique purposes and functionalities.  
2  
3 The Pandas module is used for working with tabular data. It allows us to work with data in table form,  
4 such as in CSV or SQL database formats. We can also create tables of our own,  
5 and edit or add columns or rows to tables.  
6 Pandas provides us with some powerful objects like DataFrames and Series which  
7 are very useful for working with and analyzing data.  
8  
9 The Numpy module is mainly used for working with numerical data.  
10 It provides us with a powerful object known as an Array.  
11 With Arrays, we can perform mathematical operations on multiple values  
12 in the Arrays at the same time, and also perform operations between different Arrays,  
13 similar to matrix operations.  
14  
15 Last, but not least, the Matplotlib module is used for data visualization.  
16 It provides functionality for us to draw charts and graphs,  
17 so that we can better understand and present the data visually.  
18  
19 These modules have different purposes and functionality they excel at,  
20 and together they allow us to analyze,  
21 manipulate and visualize data in very useful ways.
```

In [21]:

```
1 import matplotlib.pyplot as plt  
2 import pandas as pd  
3 import numpy as np
```

```
1 df.plot(kind = '<graph_type>', x = 'x_tilte', y = 'y_axis',color='',ls='',marker='<o/>',xlim,xticks)  
2 1.kind='line','bar','hist','bar','barh(橫條)','scatter','pie'  
3 2.x/y=names of axis  
4 3.ls=linestyle,<->實線,<-->虛線,<:>點線,<-*>虛點線  
5 4.lw=linewidth,  
6 5.figsize=[10,5],[8,4]: x,y axis length ratio  
7  
8 ~~~~
```

```
9 df.plot is friendly with tabular data,
10 so CAN combine loc
11 eg:one row one line
12 df.iloc[1].plot()
13 df.iloc[2].plot()
```

example_1: school marks

In [31]:

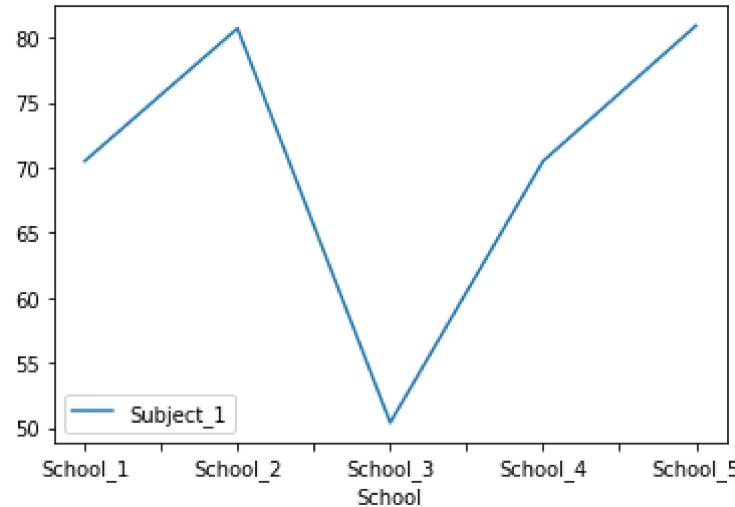
```
1 """ref:example
2 """
3 df = pd.DataFrame({
4     'Subject_1': [70.5, 80.7, 50.4, 70.5, 80.9],
5     'Subject_2': [40.24, 50.9, 70.6, 80.1, 50.9],
6     'Subject_3': [30, 50.5, 70.8, 90.88, 30],
7     'School': ['School_1', 'School_2', 'School_3', 'School_4', 'School_5'],
8     'Date': ['2021/06/01', '2021/06/02', '2021/06/03', '2021/06/04', '2021/06/05']
9 })
10 df
```

Out[31]:

	Subject_1	Subject_2	Subject_3	School	Date
0	70.5	40.24	30.00	School_1	2021/06/01
1	80.7	50.90	50.50	School_2	2021/06/02
2	50.4	70.60	70.80	School_3	2021/06/03
3	70.5	80.10	90.88	School_4	2021/06/04
4	80.9	50.90	30.00	School_5	2021/06/05

In [32]:

```
1 """
2 line plot with x and y values defined
3 抽其中一個object 出來 · 自己define x,y axis values
4 the x,y 100% same as the column name, row value
5 自己識食
6 """
7 #compare schools's sub_1 values
8 df.plot(x='School', y='Subject_1');
```



real excel case for gas_price

In [25]:

```
1 """
2 PATH='/ FOR folder, \file.xxx'
3 """
4 df_gas=pd.read_csv('C:/Users/fengs/Desktop/pd_real/mathplot\gas_prices.csv')
5 df_gas.head()
```

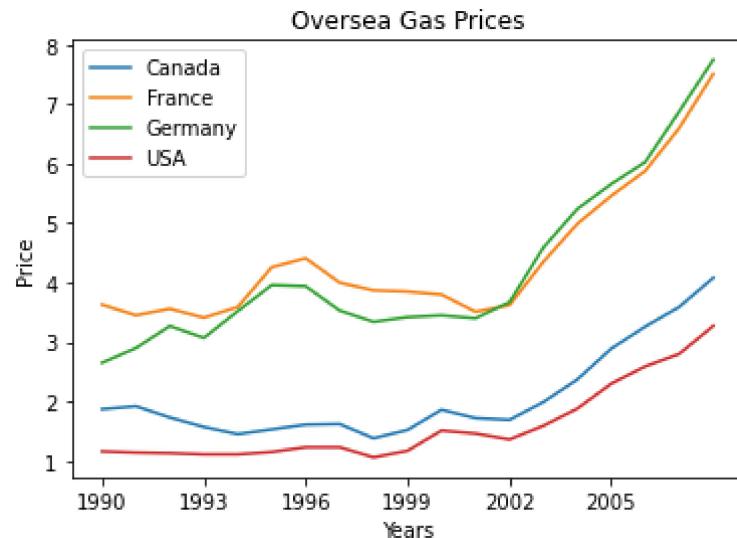
Out[25]:

	Year	Australia	Canada	France	Germany	Italy	Japan	Mexico	South Korea	UK	USA
0	1990	NaN	1.87	3.63	2.65	4.59	3.16	1.00	2.05	2.82	1.16
1	1991	1.96	1.92	3.45	2.90	4.50	3.46	1.30	2.49	3.01	1.14
2	1992	1.89	1.73	3.56	3.27	4.53	3.58	1.50	2.65	3.06	1.13
3	1993	1.73	1.57	3.41	3.07	3.68	4.16	1.56	2.88	2.84	1.11
4	1994	1.84	1.45	3.59	3.52	3.70	4.36	1.48	2.87	2.99	1.11

In [83]:

```
1 #try to use pandas.plot only, no matplotlib
2 #CANNOT g1=df.plot g2=df.plot, it will come out 2 graphs
3 #link:https://realpython.com/pandas-plot-python/
4 df_gas.plot(x='Year',y=['Canada','France','Germany','USA'], #must know:to draw multi-lines
5             xticks=range(df_gas['Year'].min(),df_gas['Year'].max(),3),
6             xlabel='Years',ylabel='Price',title='Oversea Gas Prices'
7             )
8
```

Out[83]: <AxesSubplot:title={'center':'Oversea Gas Prices'}, xlabel='Years', ylabel='Price'>



In [38]:

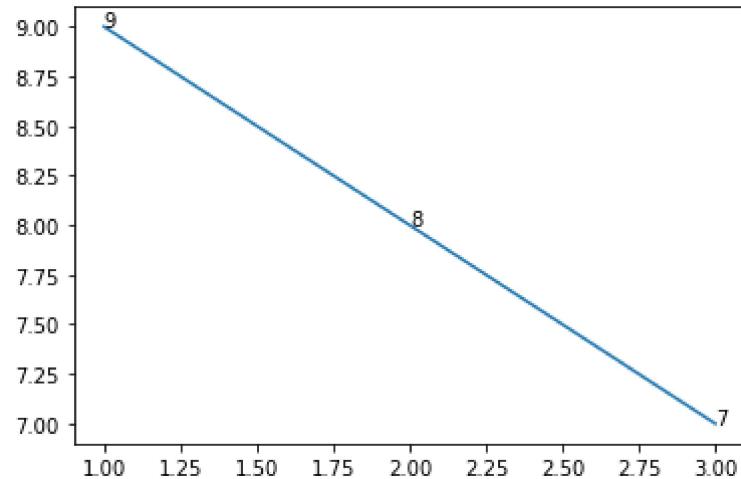
```
1 """prepare the x,y axis values"""
2 print(df_gas['Year'].min(),df_gas['Year'].max())
```

1990 2008

demo show point values

In [52]:

```
1 import matplotlib.pyplot as plt
2
3 x=[1,2,3]
4 y=[9,8,7]
5
6 plt.plot(x,y)
7 for a,b in zip(x, y):
8     plt.text(a, b, str(b))
9 plt.show()
```



chapter 18 :Mathplotlib graphics

```
1 syntax:
2 plt.plot(x,y,color='r,g,b',marker=<o/.>,ls=<--/-/:>,label='')
3
4 1.figsize is size on container
5      # Creating a new figure with width = 5 inches
6      # and height = 4 inches
7      fig = plt.figure(figsize =(5, 4))
```

```
8  
9 2.subplot=container分成幾多個小正方形，用來在一個containers  
10      draw N graphics,如果1個就不用理，default(1,1,1)  
11 3.label=line name  
12 4.ls=line style
```

In [58]:

```
1 import matplotlib.pyplot as plt  
2 df_gas=pd.read_csv('C:/Users/fengs/Desktop/pd_real/mathplot\gas_prices.csv')
```

In [82]:

```
1 df_gas.head(2)
```

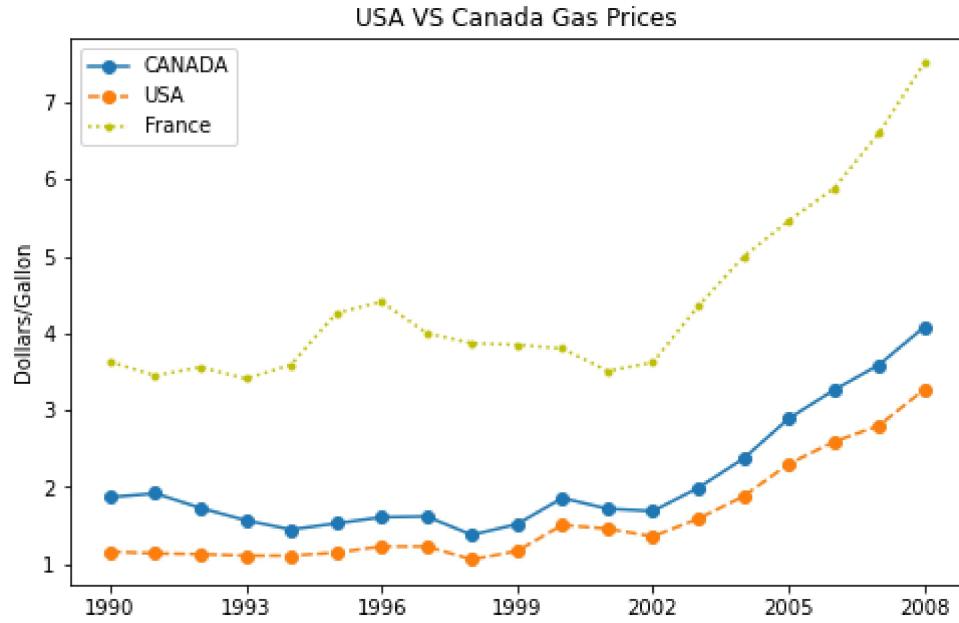
Out[82]:

	Year	Australia	Canada	France	Germany	Italy	Japan	Mexico	South Korea	UK	USA
0	1990	NaN	1.87	3.63	2.65	4.59	3.16	1.0	2.05	2.82	1.16
1	1991	1.96	1.92	3.45	2.90	4.50	3.46	1.3	2.49	3.01	1.14

18.1 :EX1: Multi-line

In [85]:

```
1 #1.commonly shared by Lines
2 plt.figure(figsize=(8,5)) #x=8 inch ,y =5 inch
3 plt.title('USA VS Canada Gas Prices')
4 plt.ylabel('Dollars/Gallon')
5
6 #x-value intervals
7 plt.xticks(df_gas['Year'][::3])
8
9 #2.no :x=[],y=[] arg, directly x,y variable
10 plt.plot(df_gas['Year'],df_gas['Canada'],ls='-',marker='o' ,
11           label='CANADA')
12 plt.plot(df_gas['Year'],df_gas['USA'],ls='--',marker='o' ,
13           label='USA')
14 #shortform yellow line, dot mark,: 虛線
15 plt.plot(df_gas['Year'],df_gas['France'],'y.:' ,
16           label='France')
17
18 #3.Legend=傳說/圖像說明
19 plt.legend() #to display the Line Label
20 plt.show()
21
22 """note
23 1.can y[['1','2'....]],but not sugg, as mark label problem
24 2.type change is not (kind='x'),but pltbar()/plthist()/
25 """
```

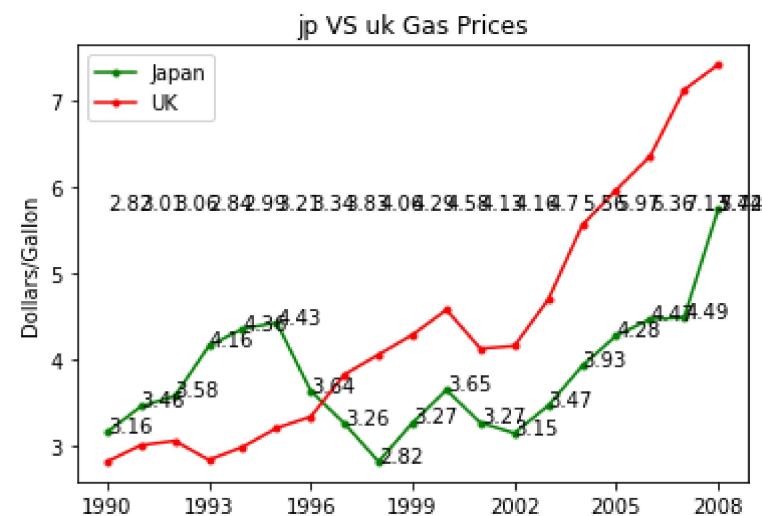


Out[85]: "note\n1.can y[['1','2'....]],but not sugg, as mark label problem\n2.type change is not (kind='x'),but pltbar()/plthist()\n"

18.1.1 print point values of lines

In [105]:

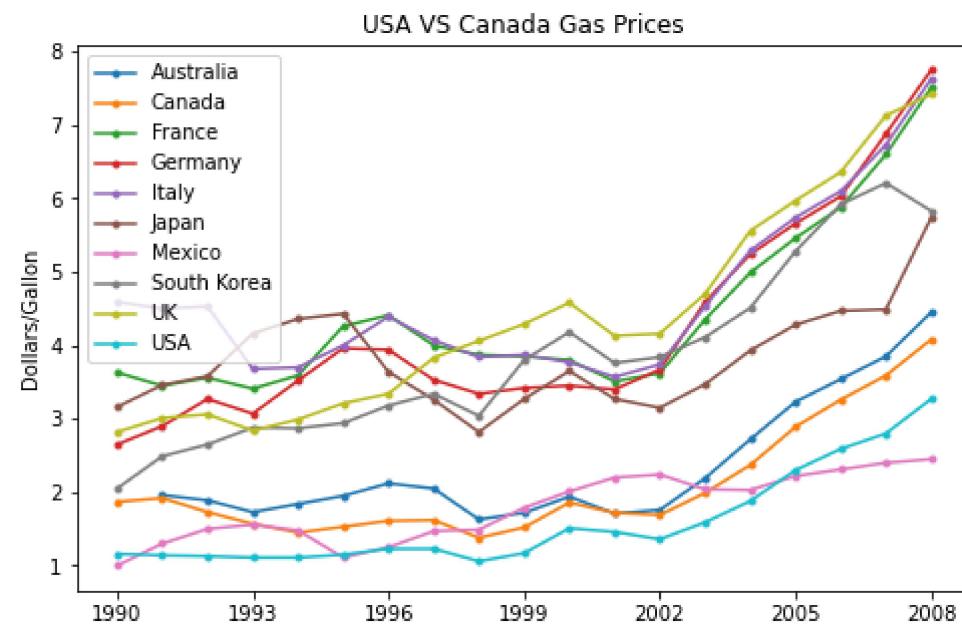
```
1 # plt.figure(figsize=(8,5)) #x=8 inch ,y =5 inch
2 plt.title('jp VS uk Gas Prices')
3 plt.ylabel('Dollars/Gallon')
4 plt.xticks(df_gas['Year'][::3])
5
6 x1=df_gas['Year']
7 y1=df_gas['Japan']
8 y2=df_gas['UK']
9
10 plt.plot(x1,y1,'g.-',label='Japan')
11 plt.plot(x1,y2,'r.-',label='UK')
12
13 for a1,b1 in zip(x1,y1):
14     plt.text(a1,b1,str(b1))
15
16 for a2,b2 in zip(x1,y2):
17     plt.text(a2,b2,str(b2))
18
19 plt.legend() #to display the Line Label
20 plt.show()
```



18.2 quick way to show all lines

In [89]:

```
1 #1. commonly shared by Lines
2 plt.figure(figsize=(8,5)) #x=8 inch ,y =5 inch
3 plt.title('USA VS Canada Gas Prices')
4 plt.ylabel('Dollars/Gallon')
5
6 #x-value intervals
7 plt.xticks(df_gas['Year'][::3])
8
9 #for each Loop print all countries
10 for country in df_gas:
11     if country != 'Year' :
12         plt.plot(df_gas['Year'],df_gas[country],marker='.',label=country)
13
14 #3. Legend=傳說/圖像說明
15 plt.legend() #to display the Line Label
16 plt.show()
```



In [86]:

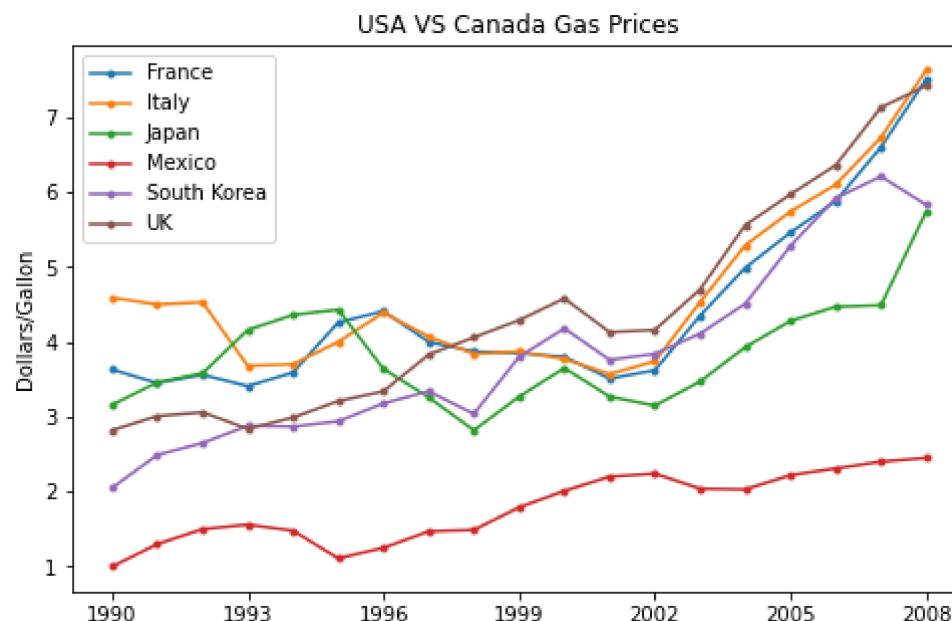
```
1 for country in df_gas:  
2     print(country)
```

Year
Australia
Canada
France
Germany
Italy
Japan
Mexico
South Korea
UK
USA

18.3 only not plot some lines/columns

In [92]:

```
1 #1.commonly shared by lines
2 plt.figure(figsize=(8,5)) #x=8 inch ,y =5 inch
3 plt.title('USA VS Canada Gas Prices')
4 plt.ylabel('Dollars/Gallon')
5
6 #x-value intervals
7 plt.xticks(df_gas['Year'][::3])
8
9
10 #make a not print list
11 country_not_look=['Australia','USA','Canada','Germany','USA']
12
13 #for each Loop print all countries
14 for country in df_gas:
15     #not use isin(); us in/not in
16     if (country!='Year')&(country not in country_not_look):
17         plt.plot(df_gas['Year'],df_gas[country],marker='.',label=country)
18
19 #3.legend=傳說/圖像說明
20 plt.legend() #to display the line label
21 plt.show()
```



chapter 19 pie chart

```
1 1.bar is similar to line, but pie is another thing
2 2.bar need to change value to pct
3 syntax:
4 plt.pie()
5 a)labels
6 b)autopct:% of each parts
7     autopct='%.2f':2 float
8     autopct='%.2f%%': show % as well
9 c)explode(0,0,0.2,0)
10    ~4 parts, 3rd part 凸其0.2格
11    ~突出重點or 不要太靠近
12 d)pctdistance=0.6 : pct 距離圓心0.6
```

In [1]:

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
```

```
In [10]: 1 #ref:df_gas=pd.read_csv('C:/Users/fengs/Desktop/pd_real/mathplot/gas_prices.csv')
2 df_fifa=pd.read_csv(r'C:/Users/fengs/Desktop/pd_real/mathplot/fifa_data.csv')
3 df_fifa.head()
```

Out[10]:

	Unnamed: 0	ID	Name	Age	Photo	Nationality	Flag	Overall	Potential	
0	0	158023	L. Messi	31	https://cdn.sofifa.org/players/4/19/158023.png	Argentina	https://cdn.sofifa.org/flags/52.png	94	94	Ba
1	1	20801	Cristiano Ronaldo	33	https://cdn.sofifa.org/players/4/19/20801.png	Portugal	https://cdn.sofifa.org/flags/38.png	94	94	J
2	2	190871	Neymar Jr	26	https://cdn.sofifa.org/players/4/19/190871.png	Brazil	https://cdn.sofifa.org/flags/54.png	92	93	Par
3	3	193080	De Gea	27	https://cdn.sofifa.org/players/4/19/193080.png	Spain	https://cdn.sofifa.org/flags/45.png	91	93	Mar
4	4	192985	K. De Bruyne	27	https://cdn.sofifa.org/players/4/19/192985.png	Belgium	https://cdn.sofifa.org/flags/7.png	91	92	Mar

5 rows × 89 columns

chapter 19.1 simple pie of football

```
In [3]: 1 #1.df_fifa.columns
2 print(df_fifa['Preferred Foot'].unique())
3 #return the count value only,[0] fist
4 df_fifa.loc[df_fifa['Preferred Foot']=='Left'].count()[0]
```

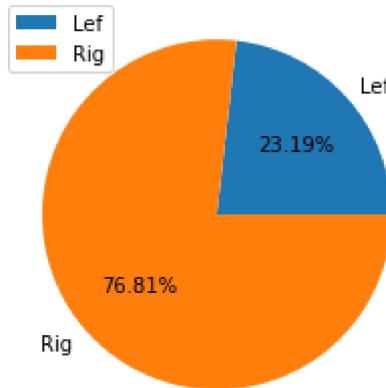
['Left' 'Right' nan]

Out[3]: 4211

```
1 matplotlib.pyplot.pie(x, explode=None, labels=None, colors=None, autopct=None,
2                         pctdistance=0.6, shadow=False, labeldistance=1.1, startangle=0, radius=1,
3                         counterclock=True, wedgeprops=None, textprops=None, center=(0, 0),
4                         rotatelabels=False, *, normalize=True, data=None)
```

In [5]:

```
1 left=df_fifa.loc[df_fifa['Preferred Foot']=='Left'].count()[0]
2 right=df_fifa.loc[df_fifa['Preferred Foot']=='Right'].count()[0]
3
4 #[left,right] as only can full in x, 1-d array
5 #'.2f'=2 float, '%.2f%%'=add % sm
6 plt.pie([left,right],labels=['Lef','Rig'],autopct='%.2f%%')
7
8 plt.title('foot preferences of players')
9 plt.legend()
10 plt.show()
```



chaper 19.2 more real case of footabll

In [21]:

```

1 """background
2 df_fifa['Weight'].head()
3 159lbs is str but not number, need to substr
4 but some are number, some ar nan<null>
5 bad data sameple
6
7 # storing first 3 letters of name as username
8 df['UserName'] = df['Name'].str[:3]
9
10 #strip()=remove space,
11 strip('xx')=remove char
12 """
13 #1:update the Weight column
14 #'folat' object has no attribute 'strip'
15 #list comprehension, only get the true value
16 df_fifa['Weight']=[x.strip('lbs') if type(x)==str else x # 放前面可以令else做額外
17                         for x in df_fifa['Weight']]
18
19 df_fifa['Weight'][0] # get first number
20 #'159' is str but not number

```

Out[21]: '159'

```

1 use for loop+ if in dataframe
2 df['new column name'] = df['column name'].apply(
3     lambda x: 'value if condition is met' if x condition
4     else 'value if condition is not met')
5

```

Convert list of NaNs and strings to int?

In []:

1

In [57]:

```
1 """problem 1: 'list' object has no attribute 'astype'  
2 一定是df=df.astype(int/float)  
3 """  
4 # df_fifa['Weight']=[x.strip('lbs') if type(x)==str else x # 放前面可以令else做額外  
5 #                 for x in df_fifa['Weight']]  
6 #                     .astype[float]  
7 #AttributeError: 'list' object has no attribute 'astype'  
8  
9  
10 """problems 2: some a nan so  
11 failed to convert  
12 ValueError: cannot convert float NaN to integer  
13 #FAILED  
14 df_fifa['Weight']=df_fifa['Weight'].astype(int)  
15 """  
16 print()
```

note:find NaN rows

In [49]:

```

1 """
2 df_fifa['Weight'].isna() only return T/F
3 """
4 #S1:df.loc[<where conditions>,[cols]]
5 #df_fifa.loc[df_fifa['Weight'].isna(),['ID','Name','Age','Weight']].head()
6
7 """
8 S2: no use pandas, but only add where
9 #df_fifa[df_fifa['Weight'].isna()].head()
10 #df_fifa[where conditions]
11
12 df_fifa[(df_fifa['Weight'].isna()) & (df_fifa.Age<30)] #valid
13 """
14 df_fifa.loc[(df_fifa['Weight'].isna()) & (df_fifa.Age<30)
15           ,df_fifa.columns.isin(['Age','ID','Name','Weight'])].head()
16 #isin is df attribute, if(x not in listx) when normal condition
17

```

Out[49]:

	ID	Name	Age	Weight
13237	195380	J. Barrera	29	NaN
13239	240437	A. Semprini	20	NaN
13240	209462	R. Bingham	24	NaN
13241	219702	K. Dankowski	21	NaN
13242	225590	I. Colman	23	NaN

note: convert nan in string cols to empty str

In []:

```
1 df[['column1','column2']] = df[['column1','column2']].fillna('')
```

In [64]:

```

1 #1.fill in nan
2 #df_fifa['Weight']=df_fifa['Weight'].fillna('')
3
4 #2.check whether still nan row
5 #df_fifa[df_fifa['Weight'].isna()]
6
7 #3.change str to int
8 #ref:df['Price'] = df['Price'].astype(int)
9
10 type(df_fifa.dtypes) #pandas.core.series.Series

```

Out[64]: pandas.core.series.Series

note:only not show some columns

In [54]:

```

1 """
2 df.loc[<row condition>, <~df.columns.isin(['a','b'])>]
3
4 1/exclude one clos
5 #df.loc[:, df.columns != 'column1']
6
7 2/exclude some clos
8 df.loc[:, ~df.columns.isin(['a','b'])]
9 """
10 df_fifa.loc[0:3, df_fifa.columns.isin(['Age', 'ID', 'Name', 'Weight'])]

```

Out[54]:

	ID	Name	Age	Weight
0	158023	L. Messi	31	159
1	20801	Cristiano Ronaldo	33	183
2	190871	Neymar Jr	26	150
3	193080	De Gea	27	168

note: drop cols in dataframe

```
In [ ]: 1 newdf = df.drop("age", axis='columns')
```

```
In [66]: 1 #delete dataframe
2 #del [[df_1,df_2]]
3 del [df_fifa]
```

19.2.1 start the football case

```
In [6]: 1 #1.del first
2 #del [df_fifa]
3 df_fifa=pd.read_csv(r'C:/Users/fengs/Desktop/pd_real/mathplot\fifa_data.csv')
4 df_fifa.head()
```

Out[6]:

	Unnamed: 0	ID	Name	Age	Photo	Nationality	Flag	Overall	Potential	Cl
0	0	158023	L. Messi	31	https://cdn.sofifa.org/players/4/19/158023.png	Argentina	https://cdn.sofifa.org/flags/52.png	94	94	Barcelo
1	1	20801	Cristiano Ronaldo	33	https://cdn.sofifa.org/players/4/19/20801.png	Portugal	https://cdn.sofifa.org/flags/38.png	94	94	Juven
2	2	190871	Neymar Jr	26	https://cdn.sofifa.org/players/4/19/190871.png	Brazil	https://cdn.sofifa.org/flags/54.png	92	93	Paris Sa Germ
3	3	193080	De Gea	27	https://cdn.sofifa.org/players/4/19/193080.png	Spain	https://cdn.sofifa.org/flags/45.png	91	93	Manches Uni
4	4	192985	K. De Bruyne	27	https://cdn.sofifa.org/players/4/19/192985.png	Belgium	https://cdn.sofifa.org/flags/7.png	91	92	Manches C

5 rows × 89 columns



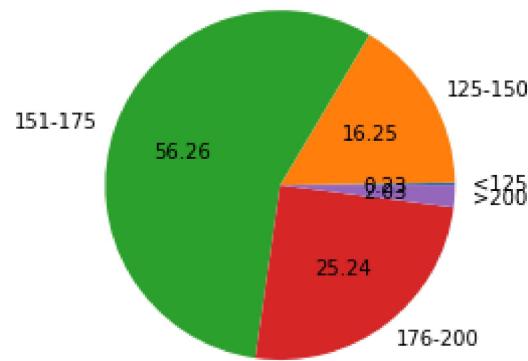
In [4]:

```
1 df_fifa['Weight']=[int(x.strip('lbs')) if type(x)==str else x # 放前面可以令else做額外  
2             for x in df_fifa['Weight']  
3             ]  
4 print(df_fifa['Weight'][0])  
5 print(df_fifa['Weight'].dtype)
```

159.0
float64

In [13]:

```
1 light=df_fifa.loc[df_fifa.Weight<125].count()[0]  
2 medium=df_fifa.loc[df_fifa.Weight.between(125,150)].count()[0]  
3 heavy=df_fifa.loc[df_fifa.Weight.between(151,175)].count()[0]  
4 giant=df_fifa.loc[df_fifa.Weight.between(175,200)].count()[0]  
5 sg=df_fifa.loc[df_fifa.Weight>200].count()[0]  
6  
7 #number List but not string name List,as put as x in plt.pie(x)  
8 weights=[light,medium,heavy,giant,sg]  
9 #Labels names list  
10 labels=['<125','125-150','151-175','176-200','>200']  
11  
12 plt.pie(weights,labels=labels,autopct='%0.2f')  
13 plt.show()
```

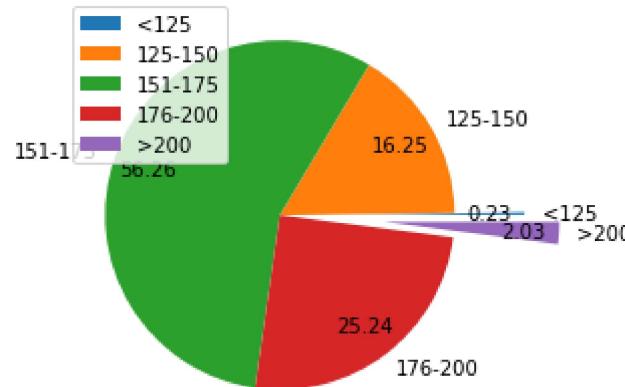


In [24]:

```

1 """fine-tune pie above
2 .explode()
3 .pctdistance=<0.6-0.8>
4 """
5 plt.style=('ggplot') # one color/display template
6 plt.title='Weight of Player'
7
8 plt.pie(weights,labels=labels,autopct='%0.2f',pctdistance=0.8,explode=(0.4,0,0,0,0.6))
9 #as only <125 ,>200 are tight, 數番上面weights[]排位
10
11 plt.legend() # give the color description
12 plt.show()

```



chapter 19.3 boxplot

```

1 boxplot(箱形圖)compare 大家的 5 個數字：
2 1.頭min · max,top 25%/75% · median,
3
4 display the five-number summary of a dataset
5 min/max/1st quartile/median/3rd quartile

```

19.3.1 basic example

In []:

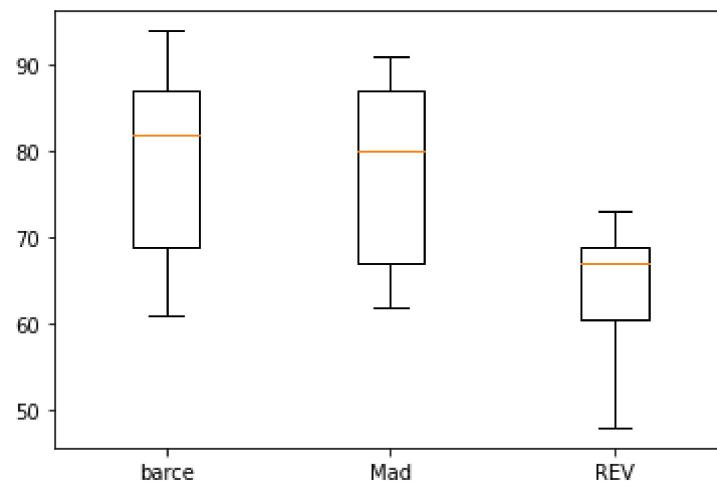
```
1 df=pf.DataFrame{'A':[1,1,2,2],  
2                 'B':[ '' ],  
3  
4 }
```

In []:

```
1
```

In [30]:

```
1 """syntax
2 plt.boxplot([d1,d2],....) [x,y] is data,if single just x,
3 """
4 barcelona=df_fifa.loc[df_fifa.Club=='FC Barcelona']['Overall']
5 madrid=df_fifa.loc[df_fifa.Club=='Real Madrid']['Overall']
6 revs=df_fifa.loc[df_fifa.Club=='New England Revolution']['Overall']
7
8 plt.boxplot([barcelona,madrid,revs],labels=['barce','Mad','REV'])
9
10 plt.title('Team Soccer Comparsion')
11 plt.ylabel('Overall Rate')
12
13 plt.show()
```



chapter 19.x:create nested pie with subplot

In []:

1

chpater 20:statistic Pandas

In []:

```

1 1.video link:
2     https://www.youtube.com/watch?v=txMdrV1Ut64
3 2.ref link
4     https://www.geeksforgeeks.org/python-pandas-dataframe-groupby/

```

chapter 20.1 groupby

```

1 Syntax: df.groupby(by=None, axis=0, level=None,
2                         as_index=True, sort=True, group_keys=True, squeeze=False, **kwargs)

```

EX1: NBA worker info

In [7]:

```

1 #ref:df_fifa=pd.read_csv(r'C:/Users/fengs/Desktop/pd_real/mathplotlib/fifa_data.csv')
2 df_nba=pd.read_csv(r'C:\Users\fengs\Desktop\pd_real\nba.csv')
3 df_nba.head()

```

Out[7]:

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
0	Avery Bradley	Boston Celtics	0.0	PG	25.0	6-2	180.0	Texas	7730337.0
1	Jae Crowder	Boston Celtics	99.0	SF	25.0	6-6	235.0	Marquette	6796117.0
2	John Holland	Boston Celtics	30.0	SG	27.0	6-5	205.0	Boston University	NaN
3	R.J. Hunter	Boston Celtics	28.0	SG	22.0	6-5	185.0	Georgia State	1148640.0
4	Jonas Jerebko	Boston Celtics	8.0	PF	29.0	6-10	231.0	NaN	5000000.0

In [8]:

```

1 #show unique
2 print(df_nba['Team'].unique())

```

Out[8]: array(['Boston Celtics', 'Brooklyn Nets', 'New York Knicks',
 'Philadelphia 76ers', 'Toronto Raptors', 'Golden State Warriors',
 'Los Angeles Clippers', 'Los Angeles Lakers', 'Phoenix Suns',
 'Sacramento Kings', 'Chicago Bulls', 'Cleveland Cavaliers',
 'Detroit Pistons', 'Indiana Pacers', 'Milwaukee Bucks',
 'Dallas Mavericks', 'Houston Rockets', 'Memphis Grizzlies',
 'New Orleans Pelicans', 'San Antonio Spurs', 'Atlanta Hawks',
 'Charlotte Hornets', 'Miami Heat', 'Orlando Magic',
 'Washington Wizards', 'Denver Nuggets', 'Minnesota Timberwolves',
 'Oklahoma City Thunder', 'Portland Trail Blazers', 'Utah Jazz',
 nan], dtype=object)

In [14]:

```

1 """
2 df.groupby(['<col_name>']).<arg>()
3 #<arg>.():mean(),count()
4 """
5 team=df_nba.groupby('Team') #return a ref_no
6 team.first().head() #first() return 1st record of each group

```

Out[14]:

	Name	Number	Position	Age	Height	Weight	College	Salary
--	------	--------	----------	-----	--------	--------	---------	--------

Team

Atlanta Hawks	Kent Bazemore	24.0	SF	26.0	6-5	201.0	Old Dominion	2000000.0
Boston Celtics	Avery Bradley	0.0	PG	25.0	6-2	180.0	Texas	7730337.0
Brooklyn Nets	Bojan Bogdanovic	44.0	SG	27.0	6-8	216.0	Oklahoma State	3425510.0
Charlotte Hornets	Nicolas Batum	5.0	SG	27.0	6-8	200.0	Virginia Commonwealth	13125306.0
Chicago Bulls	Cameron Bairstow	41.0	PF	25.0	6-9	250.0	New Mexico	845059.0

In [16]:

```
1 #return a specific group value
2 team.get_group('New York Knicks').tail(3)
```

Out[16]:

	Name	Team	Number	Position	Age	Height	Weight	College	Salary
43	Sasha Vujacic	New York Knicks	18.0	SG	32.0	6-7	195.0	NaN	947276.0
44	Derrick Williams	New York Knicks	23.0	PF	25.0	6-8	240.0	Arizona	4000000.0
45	Tony Wroten	New York Knicks	5.0	SG	23.0	6-6	205.0	Washington	167406.0

In [18]:

```
1 """group by 2 cols
2 """
3 #use <df_name>.groupby, but not pd.groupby()
4 duo_group=df_nba.groupby(['Team','Position'])
5 duo_group.first()
```

Out[18]:

		Name	Number	Age	Height	Weight	College	Salary
	Team	Position						
Atlanta Hawks		C	Al Horford	15.0	30.0	6-10	245.0	Florida 12000000.0
		PF	Kris Humphries	43.0	31.0	6-9	235.0	Minnesota 10000000.0
		PG	Dennis Schroder	17.0	22.0	6-1	172.0	Wake Forest 1763400.0
		SF	Kent Bazemore	24.0	26.0	6-5	201.0	Old Dominion 2000000.0
		SG	Tim Hardaway Jr.	10.0	24.0	6-6	205.0	Michigan 1304520.0
...	
Washington Wizards		C	Marcin Gortat	13.0	32.0	6-11	240.0	North Carolina State 11217391.0
		PF	Drew Gooden	90.0	34.0	6-10	250.0	Kansas 3300000.0
		PG	Ramon Sessions	7.0	30.0	6-3	190.0	Nevada 2170465.0
		SF	Jared Dudley	1.0	30.0	6-7	225.0	Boston College 4375000.0
		SG	Alan Anderson	6.0	33.0	6-6	220.0	Michigan State 4000000.0

149 rows × 7 columns

In [27]:

```
1 """2.with agg function
2 groupby(['col1'])['col2'].sum()/mean()
3 yet didnot add any where conditon
4 """
5 df_nba.groupby(['Team'])['Team'].count().head()
```

Out[27]:

```
Team
Atlanta Hawks      15
Boston Celtics     15
Brooklyn Nets       15
Charlotte Hornets   15
Chicago Bulls        15
Name: Team, dtype: int64
```

In [25]:

```
1 #same as above
2 df_nba.value_counts(['Team']).head()
```

Out[25]:

```
Team
New Orleans Pelicans 19
Memphis Grizzlies    18
New York Knicks      16
Milwaukee Bucks      16
Atlanta Hawks         15
dtype: int64
```

Groupby+agg()

In [31]:

```
1 #find the max,min age of each team
2 df_nba.groupby(['Team'])['Age'].agg(['min','max']).head()
```

Out[31]:

min max

Team	min	max
Atlanta Hawks	22.0	35.0
Boston Celtics	20.0	29.0
Brooklyn Nets	21.0	32.0
Charlotte Hornets	21.0	31.0
Chicago Bulls	21.0	35.0

EX2: Phone Info(實用)

In []:

```
1 useful link:
2 https://www.shanelynn.ie/summarising-aggregation-and-grouping-data-in-python-pandas/
```

Note: Steps to insert image:

```
1 1.change cell to 'Markdown' but not code
2 2.Edit>inset image(bottom)
```

"" link: <https://pandas.pydata.org/pandas-docs/version/0.22/generated/> ""

Select these rows

Group by this column

Work on these columns and...

...Perform these operations

```
data[data['item'] == 'call'].groupby('month').agg(
    max_duration=pd.NamedAgg(column='duration', aggfunc=max),
    min_duration=pd.NamedAgg(column='duration', aggfunc=min),
    total_duration=pd.NamedAgg(column='duration', aggfunc=sum),
    num_days=pd.NamedAgg(column="date", aggfunc=lambda x: (max(x) - min(x)).days)
)
```

Rename the outputs

In [32]:

```
1 df_phone=pd.read_csv('C:/Users/fengs/Desktop/pd_real\phone_data.csv')
2 df_phone.head()
```

Out[32]:

	index	date	duration	item	month	network	network_type
0	0	15/10/14 06:58	34.429	data	2014-11	data	data
1	1	15/10/14 06:58	13.000	call	2014-11	Vodafone	mobile
2	2	15/10/14 14:46	23.000	call	2014-11	Meteor	mobile
3	3	15/10/14 14:48	4.000	call	2014-11	Tesco	mobile
4	4	15/10/14 17:27	4.000	call	2014-11	Tesco	mobile

In [34]:

```
1 """
2 1.find the sum of duration of each month
3 #without agg
4 #need to consider sorting as well
5 """
6 df_phone.groupby(['month'])['duration'].sum()
```

Out[34]: month

```
2014-11    26639.441
2014-12    14641.870
2015-01    18223.299
2015-02    15522.299
2015-03    22750.441
Name: duration, dtype: float64
```

In [36]:

```
1 #get num of date each month
2 #in fact it count ref_no
3 df_phone.groupby('month')['date'].count()
```

Out[36]: month

```
2014-11    230
2014-12    157
2015-01    205
2015-02    137
2015-03    101
Name: date, dtype: int64
```

In [43]:

```
1 """
2 Q2.add where clause
3 df[where clause].groupby(['xx'])['yy'].
```

Out[43]:

duration

network	duration
Meteor	7200.0
Tesco	13828.0
Three	36464.0
Vodafone	14621.0
landline	18433.0
voicemail	1775.0

In [49]:

```

1 """
2 Q3:group by N cols,
3 loc[row,col]
4 """
5 #if loc filter out cal cols,
6 #will error:'Column not found: date'
7 #df_phone.loc[:,['month','item']].groupby(['month','item'])['date'].count()
8
9 df_phone.loc[:,['month','item','date']].groupby(['month','item'])['date'].count()

```

Out[49]:

	month	item	
2014-11	call	107	
	data	29	
	sms	94	
2014-12	call	79	
	data	30	
	sms	48	
2015-01	call	88	
	data	31	
	sms	86	
2015-02	call	67	
	data	31	
	sms	39	
2015-03	call	47	
	data	29	
	sms	25	

Name: date, dtype: int64

groupby +agg(Multi-stat)

```

1 # in realiy there will more than agg() cols
2 like the nba example, the min(),max() age of each team

```

```

1 #group by has sort/axis args
2 DataFrame.groupby(by=None, axis=0, level=None, as_index=True,
3                   sort=True, group_keys=True,

```

```
4             squeeze=NoDefault.no_default, observed=False,
5             dropna=True)
6
7 #sort: bool, default True
8 #axis{0 or 'index', 1 or 'columns'}, default 0,
9   group by each row
10 Split along rows (0) or columns (1).
11 #as_index:bool, default True
12 For aggregated output, return object with group labels as the index.
13 Only relevant for DataFrame input. as_index=False
14 is effectively "SQL-style" grouped output.
```

In []:

```
1 """
2 1. same meaning
3 """
4 df_phone.groupby
```

In []:

```
1
```

In []:

chapter 21:pivot table

In []:

In []:

chapter 22:manipulate spreadsheets

In []:

In []:

1