



TAIJI LABORATORY
FOR GRAVITATIONAL WAVE UNIVERSE



ICTP-AP
International Centre
for Theoretical Physics Asia-Pacific
国际理论物理中心-亚太地区



中国科学院大学
University of Chinese Academy of Sciences

引力波数据探索：编程与分析实战训练营

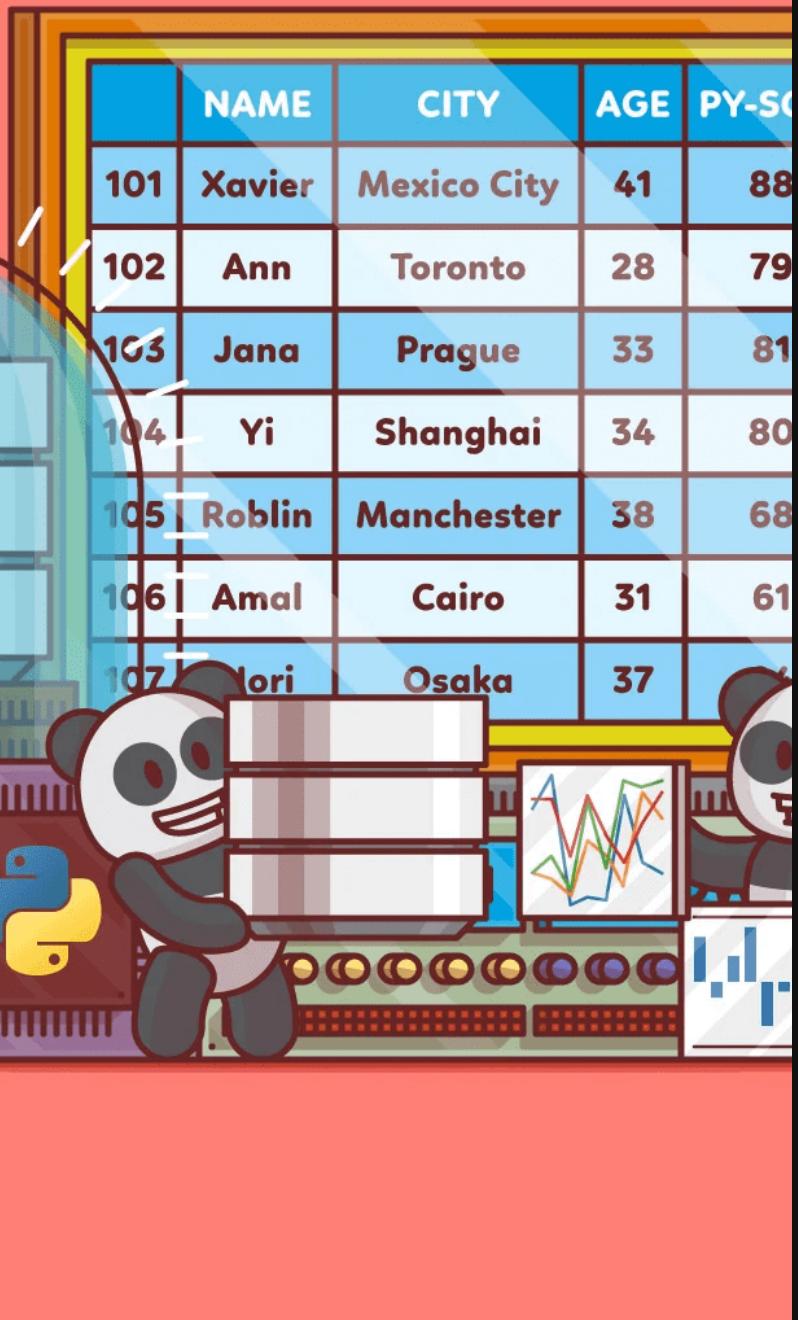
第 2 部分 基于 Python 的数据分析基础 数据分析实训之 Pandas

主讲老师：王赫

ICTP-AP, UCAS

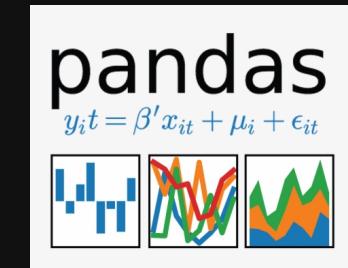
2023/12/03





数据分析实训之 Pandas

- Pandas 中的数据结构: Series 和 DataFrame
- DataFrame 的增删改查
- DataFrame 的组合与聚合&透视表
- 实例: GW Catalog 数据分析案例
- 实例: 股票数据分析案例 (optional)



“Talk is
cheap. Show
me the code.”

Linus Torvalds



Ndarray → Series/DataFrame (Pandas)

- Pandas 是基于 Numpy 的强大的分析结构化数据的工具集
 - Pandas 官方网址: <https://pandas.pydata.org>
 - Pandas 官方中文文档: <https://www.pypandas.cn>
- 从 Numpy 的 Ndarray 到 Pandas 的 Series / DataFrame
 - (Numpy) 1-dimensional array \Leftrightarrow Series (Pandas)
 - (Numpy) 2-dimensional array \Leftrightarrow DataFrame (Pandas)



<https://numpy.org>

index timestamp	col2	col3	col4	col5
200 (a datetime obj)	a	b	c	d
3 (a datetime obj)	e	f	g	h
58 (a datetime obj)	i	j	k	l
666 (a datetime obj)	m	n	o	p
14 (a datetime obj)	q	r	s	t
1 (a datetime obj)	u	v	w	x
...

pandas dataframe

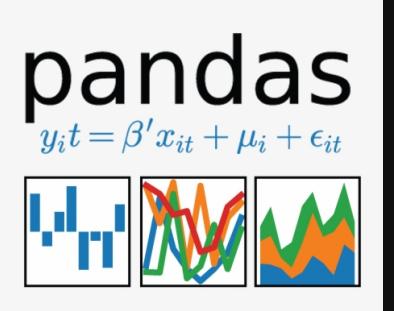
```
df.loc[[58, 14],  
       ['col3', 'col4', 'col5']]
```

0	1	2
0 b'	c'	d'
1 f'	g'	h'
2 j'	k'	l'
3 n'	o'	p'
4 r'	s'	t'
5 v'	w'	x'
...

numpy array

```
a[[2,4], :3]
```

Source





Series vs DataFrame

- **Series** 是一种类似于一维数组的对象，是由一组数据(各种 NumPy 数据类型)以及一组与之相关的数据标签(即索引)组成。仅由一组数据也可产生简单的 Series 对象。
- **DataFrame** 是 Pandas 中的一个表格型的数据结构，包含有一组有序的列，每列可以是不同的值类型(数值、字符串、布尔型等)，DataFrame 即有行索引也有列索引，可以被看做是由 Series 组成的字典。

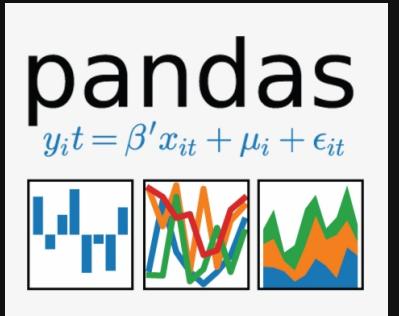


<https://numpy.org>

Series		Series		DataFrame	
	apples		oranges		
0	3	0	0	0	0
1	2	1	3	2	3
2	0	2	7	0	7
3	1	3	2	1	2

+ =

Source





Excel vs DataFrame

- **Series** 是一种类似于一维数组的对象，是由一组数据(各种 NumPy 数据类型)以及一组与之相关的数据标签(即索引)组成。仅由一组数据也可产生简单的 Series 对象。
- **DataFrame** 是 Pandas 中的一个表格型的数据结构，包含有一组有序的列，每列可以是不同的值类型(数值、字符串、布尔型等)，DataFrame 即有行索引也有列索引，可以被看做是由 Series 组成的字典。

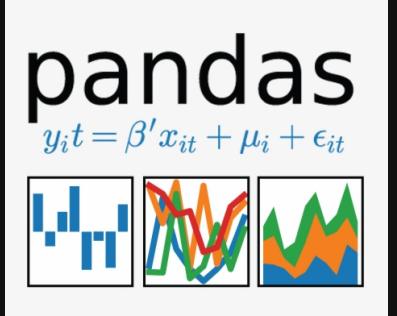


<https://numpy.org>

music.csv

pandas.read_csv('music.csv')

	Artist	Genre	Listeners	Plays
0	Billie Holiday	Jazz	1,300,000	27,000,000
1	Jimi Hendrix	Rock	2,700,000	70,000,000
2	Miles Davis	Jazz	1,500,000	48,000,000
3	SIA	Pop	2,000,000	74,000,000





DataFrame 的“增删改查”



- DataFrame 构建与初始化: 1. 按列构建 2. 按行构建

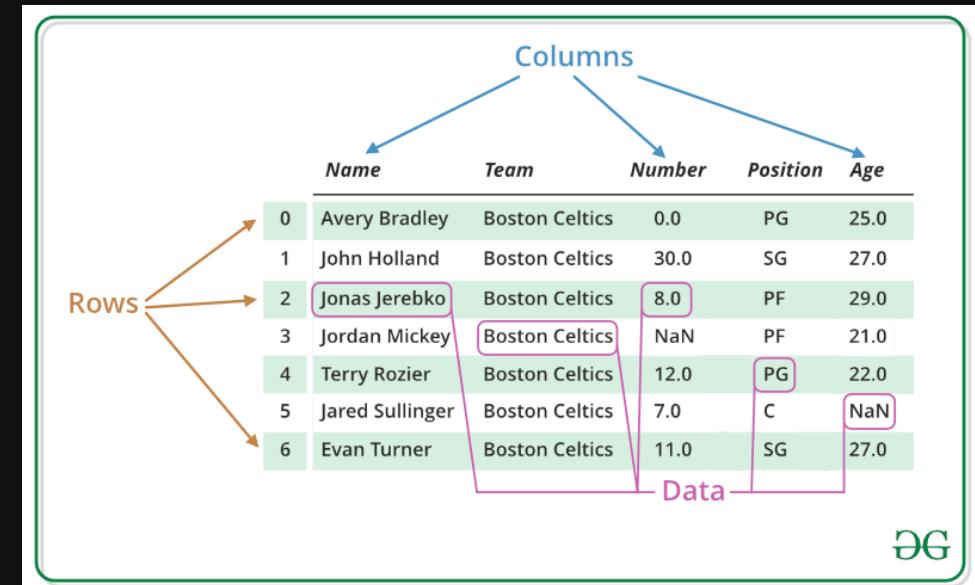
	a	b	c
1	4	7	10
2	5	8	11
3	6	9	12

```

df = pd.DataFrame(
    {"a" : [4 ,5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.

df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.

```



Source



DataFrame 的“增删改查”

- DataFrame 的构建与初始化: 1. 按列构建 2. 按行构建

		a	b	c
n	v			
d	1	4	7	10
	2	5	8	11
e	2	6	9	12

```

df = pd.DataFrame(
    {"a" : [4, 5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
        names=['n','v']))
Create DataFrame with a MultiIndex
  
```

Make New Variables

```

df=df.assign(Area=lambda df: df.Length*df.Height)
Compute and append one or more new columns.

df['Volume'] = df.Length*df.Height*df.Depth
Add single column.

pd.qcut(df.col, n, labels=False)
Bin column into n buckets.
  
```

Creating Pandas DataFrames from Python Lists and Dictionaries

	Dictionary	List																				
Row Oriented	<pre> sales = [{"account": 'Jones LLC', 'Jan': 150, 'Feb': 200, 'Mar': 140}, {"account": 'Alpha Co', 'Jan': 200, 'Feb': 210, 'Mar': 215}, {"account": 'Blue Inc', 'Jan': 50, 'Feb': 90, 'Mar': 95}] df = pd.DataFrame(sales) </pre>	<pre> sales = [('Jones LLC', 150, 200, 50), ('Alpha Co', 200, 210, 90), ('Blue Inc', 140, 215, 95)] labels = ['account', 'Jan', 'Feb', 'Mar'] df = pd.DataFrame.from_records(sales, columns=labels) </pre>																				
default	<table border="1"> <thead> <tr> <th></th> <th>account</th> <th>Jan</th> <th>Feb</th> <th>Mar</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>Jones LLC</td> <td>150</td> <td>200</td> <td>140</td> </tr> <tr> <td>1</td> <td>Alpha Co</td> <td>200</td> <td>210</td> <td>215</td> </tr> <tr> <td>2</td> <td>Blue Inc</td> <td>50</td> <td>90</td> <td>95</td> </tr> </tbody> </table>		account	Jan	Feb	Mar	0	Jones LLC	150	200	140	1	Alpha Co	200	210	215	2	Blue Inc	50	90	95	from_records
	account	Jan	Feb	Mar																		
0	Jones LLC	150	200	140																		
1	Alpha Co	200	210	215																		
2	Blue Inc	50	90	95																		
Column Oriented	<pre> sales = {'account': ['Jones LLC', 'Alpha Co', 'Blue Inc'], 'Jan': [150, 200, 50], 'Feb': [200, 210, 90], 'Mar': [140, 215, 95]} df = pd.DataFrame.from_dict(sales) </pre>	<pre> sales = [['Jones LLC', 'Alpha Co', 'Blue Inc'], [150, 200, 50], [200, 210, 90], [140, 215, 95]] df = pd.DataFrame.from_items(sales) </pre>																				
from_dict	<p>When using a dictionary, column order is not preserved. Explicitly order them:</p> <pre> df = df[['account', 'Jan', 'Feb', 'Mar']] </pre>	from_items																				

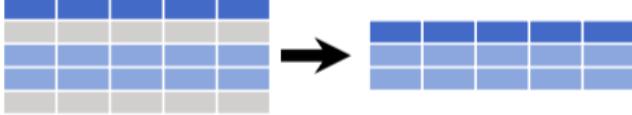
Practical Business Python - pbpython.com



DataFrame 的“增删改查”

- DataFrame 的索引与切片：

Subset Observations (Rows)



`df[df.Length > 7]`
Extract rows that meet logical criteria.

`df.drop_duplicates()`
Remove duplicate rows (only considers columns).

`df.head(n)`
Select first n rows.

`df.tail(n)`
Select last n rows.

`df.sample(frac=0.5)`
Randomly select fraction of rows.

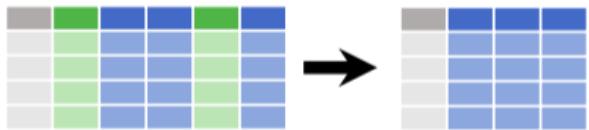
`df.sample(n=10)`
Randomly select n rows.

`df.iloc[10:20]`
Select rows by position.

`df.nlargest(n, 'value')`
Select and order top n entries.

`df.nsmallest(n, 'value')`
Select and order bottom n entries.

Subset Variables (Columns)



`df[['width', 'length', 'species']]`
Select multiple columns with specific names.

`df['width'] or df.width`
Select single column with specific name.

`df.filter(regex='regex')`
Select columns whose name matches regular expression *regex*.

`df.drop(['length', 'species'], axis=1)`

Just like ndarray ...

`df.loc[:, 'x2':'x4']`
Select all columns between x2 and x4 (inclusive).

`df.iloc[:, [1, 2, 5]]`
Select columns in positions 1, 2 and 5 (first column is 0).

`df.loc[df['a'] > 10, ['a', 'c']]`
Select rows meeting logical condition, and only the specific columns .



DataFrame 的“增删改查”

- DataFrame 的改造与合并:

Reshaping Data – Change the layout of a data set

<p><code>pd.melt(df)</code> Gather columns into rows.</p>	<p><code>df.pivot(columns='var', values='val')</code> Spread rows into columns.</p>	<p><code>df=df.sort_values('mpg')</code> Order rows by values of a column (low to high).</p> <p><code>df=df.sort_values('mpg', ascending=False)</code> Order rows by values of a column (high to low).</p> <p><code>df=df.rename(columns = {'y':'year'})</code> Rename the columns of a DataFrame</p>
<p><code>pd.concat([df1, df2])</code> Append rows of DataFrames</p>	<p><code>pd.concat([df1, df2], axis=1)</code> Append columns of DataFrames</p>	<p><code>df=df.sort_index()</code> Sort the index of a DataFrame</p> <p><code>df=df.reset_index()</code> Reset index of DataFrame to row numbers, moving index to columns.</p> <p><code>df=df.drop(['Length','Height'], axis=1)</code> Drop columns from DataFrame</p>



DataFrame 的“增删改查”

- DataFrame 的改造与合并:
 - merge vs concat ? ([Ref](#))
 - concat ⇒ "append"
 - concat < merge

ydf + **zdf** =

x1	x2
A	1
B	2
C	3

Set-like Operations

`pd.merge(ydf, zdf)`
Rows that appear in both ydf and zdf
(Intersection).

x1	x2
B	2
C	3

`pd.merge(ydf, zdf, how='outer')`
Rows that appear in either or both ydf and zdf
(Union).

x1	x2
A	1
B	2
C	3
D	4

`pd.merge(ydf, zdf, how='outer', indicator=True).query('_merge == "left_only") .drop(['_merge'], axis=1)`
Rows that appear in ydf but not zdf (Setdiff).

x1	x2
A	1

Combine Data Sets

adf		bdf	
x1	x2	x1	x3
A	1	A	T
B	2	B	F
C	3	D	T

Standard Joins

`pd.merge(adf, bdf, how='left', on='x1')`
Join matching rows from bdf to adf.

x1	x2	x3
A	1	T
B	2	F
C	3	NaN

`pd.merge(adf, bdf, how='right', on='x1')`
Join matching rows from adf to bdf.

x1	x2	x3
A	1.0	T
B	2.0	F
D	NaN	T

`pd.merge(adf, bdf, how='inner', on='x1')`
Join data. Retain only rows in both sets.

x1	x2	x3
A	1	T
B	2	F

`pd.merge(adf, bdf, how='outer', on='x1')`
Join data. Retain all values, all rows.

x1	x2	x3
A	1	T
B	2	F
C	3	NaN
D	NaN	T

Filtering Joins

`adf[adf.x1.isin(bdf.x1)]`
All rows in adf that have a match in bdf.

x1	x2
A	1
B	2

`adf[~adf.x1.isin(bdf.x1)]`
All rows in adf that do not have a match in bdf.

x1	x2
C	3



DataFrame 的“增删改查”



- Series / DataFrame 常用的属性与方法

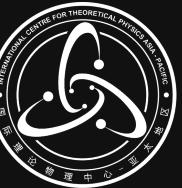
- shape / size / index / dtype / astype() / ...
- head() / tail() / describe() / values \Rightarrow ndarray (Numpy)
- max() / min() / mean() / median() / ...
- to_*() / sort_index() / sort_values() / ...
- apply() / drop() / drop_duplicates() / ...
- isin() / isna() / isnull() /fillna() / ...

- Series 的常用属性与方法

- unique()
- tolist()
- value_counts()
- map()
- is_unique() / is_monotonic()

- DataFrame 的常用属性与方法

- columns
- info()
- stack()
- insert()
- ...



DataFrame 的组合与聚合&透视表

— 数据高级统计分析

- Pandas 中常用的组合和聚合方法 [Ref](#)
 - groupby() / apply() / agg()
 - stack() / unstack()
 - melted()
 - pivot_table()

```
1 pivoted = pd.pivot_table(s4g, index=['Symbol', 'Year'],
2                           values=['Open', 'Close'], aggfunc='mean',
3                           columns=['Month'], fill_value = 0)
```



```
1 table = s4g.groupby(['Symbol', 'Year', 'Month'])['Open', 'Close'].mean()
2 table = table.unstack('Month')
3 table = table.fillna(0)
```

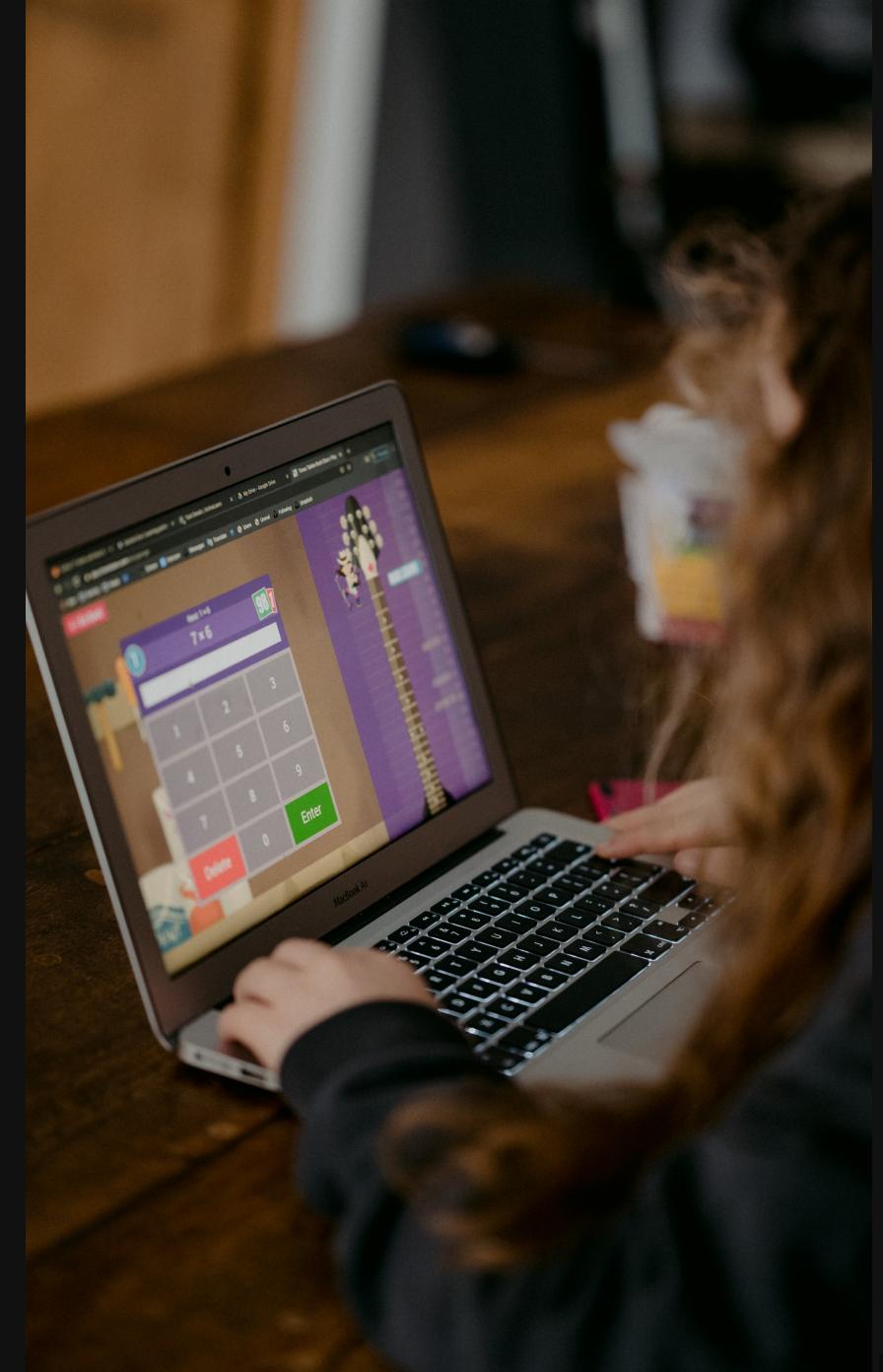
pd.pivot_table(df,
index=["Manager", "Status"],
columns=["Product"],
aggfunc=[np.sum], values=["Price"],
fill_value=0,
margins=True,
dropna=True)

		sum					
		Price					
		Product	CPU	Maintenance	Monitor	Software	
Manager	Debra Henley	declined	70000	0	0	0	70000
		pending	40000	10000	0	0	50000
		presented	30000	0	0	20000	50000
		won	65000	0	0	0	65000
	Fred Anderson	declined	65000	0	0	0	65000
pending		0	5000	0	0	5000	
presented		30000	0	5000	10000	45000	
won		165000	7000	0	0	172000	
All		465000	22000	5000	30000	522000	

Repo of the course: <https://github.com/iphysresearch/GWData-Bootcamp>

Homework

1. Python、Numpy 和 Pandas 的**基础作业**题目（单选题）位于：
 - `2023/python/homework-python.*`
 - `2023/python/homework-numpy.*`
 - `2023/python/homework-pandas.*`
 - （分别有 `ipynb`, `html`, `md` 三种文档格式，以方便阅读）
2. 将你完成的作业添加到你在上一步中创建的个人作业目录中。根据作业的类型，应将完成的作业分别命名为 `python_submit.txt`、`numpy_submit.txt` 和 `pandas_submit.txt`。其中每个 `.txt` 文档的每行对应于 A,B,C,D,... 等选项当中的一个（注意：行号对应于题号）
3. 在 `homework` 分支上把你完成的作业 `push` 到你自己的关于本课程的远程仓库中，即：`$ git push origin homework`；最后，在 GitHub 上你的远程仓库中，在 `homework` 分支下发起 **Pull Request (PR)** 至本课程远程仓库的 `homework` 分支。
4. GitHub Actions 工作流将自动检查你的提交，并将 `modified` 的 `python_submit.txt`、`numpy_submit.txt` 和 `pandas_submit.txt` 与 `solution` 进行比较。
5. 基础作业的最终成绩，根据 PR 的最新 commit 来定，记得到时候 @ 我记录成绩。
6. 不要修改其他学员的作业目录和作业内容！



Repo of the course: <https://github.com/iophysresearch/GWData-Bootcamp>

Homework

扩展作业

- 完成以下 Leetcode 中的算法题目：
 - 434. Number of Segments in a String (初级)
 - 1869. Longer Contiguous Segments of Ones than Zeros (初级)
 - 1784. Check if Binary String Has at Most One Segment of Ones (初级)
 - 852. Peak Index in a Mountain Array (中级)
 - 162. Find Peak Element (中级)
- 把上面5道算法题目的结果 comment 在你完成的**基础作业**的PR里，要求：
算法的每一行都写好中文说明注释。

```

1 class Solution:
2     def twoSum(self, nums: List[int], target: int) -> List[int]:
3         return [
4             # 取出两个索引
5             [index_i, index_j]
6             # 第一层循环
7             for index_i, i in enumerate(nums)
8             # 第二层循环
9             for index_j, j in enumerate(nums)
10            # 跳过相同的索引
11            if index_i != index_j
12            # 要求两层循环的下的值满足要求
13            if target == i+j][0]
```

