

定义一个字典

```
1 # 定义一个 DataFrame 对象
2 # 数据使用字典
3 # 也要列出列的名称以便于一一匹配
4 df = pd.DataFrame(
5     {
6         "id": [1001, 1002, 1003, 1004, 1005, 1006],
7         "date": pd.date_range("20130102", periods=6),
8         "city": ["Beijing", "SH", "guangzhou",
9         "Shenzhen", "shanghai", "BEIJING"],
10        "age": [23, 44, 54, 32, 34, 32],
11        "category": ["100-A", "100-B", "110-A", "110-C",
12        "210-A", "130-F"],
13        "price": [1200, np.nan, 2133, 5433, np.nan,
14        4432],
15    },
16    columns=["id", "date", "city", "category", "age",
17    "price"],
18 )
19 df
```

	id	date	city	category	age	price
0	1001	2013-01-02	Beijing	100-A	23	1200.0
1	1002	2013-01-03	SH	100-B	44	NaN
2	1003	2013-01-04	guangzhou	110-A	54	2133.0
3	1004	2013-01-05	Shenzhen	110-C	32	5433.0
4	1005	2013-01-06	shanghai	210-A	34	NaN
5	1006	2013-01-07	BEIJING	130-F	32	4432.0

值得注意的是，`date` 数据可以按照时间戳的方式输出。

输出形状

```
1 # 输出形状
2 df.shape
```

```
1 (6, 6)
```

输出信息

```
1 # 输出信息
2 df.info
```

```
1 <bound method DataFrame.info of          id date          city
   category age  price
2 0  1001  NaN  Beijing  100-A  23  1200.0
3 1  1002  NaN         SH  100-B  44      NaN
4 2  1003  NaN  guangzhou  110-A  54  2133.0
5 3  1004  NaN  Shenzhen  110-C  32  5433.0
6 4  1005  NaN  shanghai  210-A  34      NaN
7 5  1006  NaN   BEIJING  130-F  32  4432.0>
```

返回每行的下标和相应下标的值

```
1 # 行索引的下标以及相应下标的值
2 row_index_name_1 = df.index
3 row_index_name_2 = df.index.values
4 print(row_index_name_1)
5 print(row_index_name_2)
```

```
1 RangeIndex(start=0, stop=6, step=1)
2 [0 1 2 3 4 5]
```

输出相应变量的值

```
1 # 输出相应变量的值
2 df["city"].values
```

```
1 array(['Beijing', 'SH', 'guangzhou', 'Shenzhen',
        'shanghai', 'BEIJING'],
2        dtype=object)
```

输出变量的数据类型

```
1 df.dtypes
```

```
1 id           int64
2 date         object
3 city         object
4 category     object
5 age          int64
6 price        float64
7 dtype: object
```

判断是否为空

```
1 # 判断是否为空
2 df["price"].isnull()
```

```
1 0    False
2 1     True
3 2    False
4 3    False
5 4     True
6 5    False
7 Name: price, dtype: bool
```

去除重复值

```
1 # 去除重复值
2 df["age"].unique()
```

```
1 array([23, 44, 54, 32, 34], dtype=int64)
```

筛选

```
1 dataframe = pd.DataFrame({"a": [1, 2, 3], "b": ["aaa",
2         "bbb", "ccc"]})
3 # 选取属性 "a" 中大于 1 的行:
4 dataframe[dataframe.a > 1]
```

	a	b
1	2	bbb
2	3	ccc

```
1 # 选取属性 "a" 中大于 1 的 "b" 列:
2 dataframe["b"][dataframe.a > 1]
```

```
1 1    bbb
2 2    ccc
3 Name: b, dtype: object
```

```
1 # 选取属性 "b" 在 list 中的行
2 array = np.array(["aaa", "ccc"])
3 dataframe[dataframe.b.isin(array)]
```

	a	b
0	1	aaa
2	3	ccc

```
1 dataframe_1 = df.loc[df["city"] == "Beijing", ["age",
  "id"]]
2 dataframe_1
```

	age	id
0	23	1001

```
1 dataframe_2 = df.loc[df["age"] < 34, ["city", "id"]]
2 dataframe_2
```

	city	id
0	Beijing	1001
3	Shenzhen	1004
5	BEIJING	1006

```
1 dataframe_3 = df.loc[(df["age"] < 50) & (df["price"] <
  6000), ["id", "city"]]
2 dataframe_3
```

	id	city
0	1001	Beijing
3	1004	Shenzhen
5	1006	BEIJING

填充空缺值

```
1 # 填充空缺值，如果在原对象修改，则添加 inplace 参数，并置为 True，
  即 fillna(value=10, inplace=True)
2 new_df = df.fillna(value=10)
3 new_df
```

	id	date	city	category	age	price
0	1001	2013-01-02	Beijing	100-A	23	1200.0

1	id	date	city	category	age	price
2	1003	2013-01-04	guangzhou	110-A	54	2133.0
3	1004	2013-01-05	Shenzhen	110-C	32	5433.0
4	1005	2013-01-06	shanghai	210-A	34	10.0
5	1006	2013-01-07	BEIJING	130-F	32	4432.0

```
1 # 均值填充
2 df["price"] =
  df["price"].fillna(value=df["price"].mean())
3 # 线性填充
4 df["price"] = df["price"].interpolate()
5 # 三次样条填充
6 df["price"] = df["price"].interpolate()
```

字符处理方法

```
1 df["city"] = df["city"].str.upper()
2 df
```

数据处理方法

改变类型

```
1 df["price"] = df["price"].astype(int)
2 df
```

去重

```
1 df["city"] = df["city"].drop_duplicates()
2 df
```

合并

```
1 df_inner = pd.merge(df, df1, how="inner")
2 df_left = pd.merge(df, df1, how="left")
3 df_right = pd.merge(df, df1, how="right")
4 df_outer = pd.merge(df, df1, how="outer")
```

left 和 right 代表着左合并与右合并，一个以左边数据集为基础，一个以右边数据集为基础。

inner 只对两个所共用的进行保存，outer 是并集。

设置下标

```
1 df_inner.set_index("id")
```

将数据中的某一列设置为下标。

	date	city	category	age	price	gender	pay	m-point
id								
1001	2013-01-02	Beijing	100-A	23	1200	male	Y	10
1002	2013-01-03	Sh	100-B	44	3299	female	N	12
1003	2013-01-04	Guangzhou	110-A	54	2133	male	Y	20
1004	2013-01-05	Shenzhen	110-C	32	5433	female	Y	40
1005	2013-01-06	Shanghai	210-A	34	3299	male	N	40
1006	2013-01-07	NaN	130-F	32	4432	female	Y	40

排序

根据某一列进行排序

```
1 df_inner.sort_values(by=["price"])
```

根据下标进行排序

```
1 df_inner = df_inner.sort_index()  
2 df_inner
```

分组

```
1 data = {  
2     'Product': ['A', 'B', 'A', 'B', 'A', 'B'],  
3     'Sales': [100, 200, 150, 300, 120, 250],  
4     'Quantity': [10, 15, 12, 20, 8, 18]  
5 }  
6  
7 df = pd.DataFrame(data)  
8  
9 grouped = df.groupby("Product")  
10 total_quantities = grouped["Quantity"].sum()  
11 total_quantities
```

```
1 Product  
2 A      30  
3 B      53  
4 Name: Quantity, dtype: int64
```

抽样

在 *Pandas* 中，`sample()` 方法用于从 `DataFrame` 中随机抽取指定数量或比例的样本。下面是该方法的基本语法和一些示例用法：

1. 基本语法：


```
1 df.sample(n=None, frac=None, replace=False,
            random_state=None)
```

- `n`：要抽取的样本数量。
- `frac`：要抽取的样本比例，取值范围为 `[0, 1]`。
- `replace`：是否允许重复抽样，默认为 `False`。
- `random_state`：随机种子，用于复现随机抽样结果。

参数 `n` 和 `frac` 二选一，如果同时指定了，`frac` 会被忽略。

2. 示例：

假设有一个包含学生信息的 `DataFrame`：

```
1 import pandas as pd
2
3 data = {
4     'StudentID': [1, 2, 3, 4, 5],
5     'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eva'],
6     'Age': [20, 21, 22, 23, 24]
7 }
8
9 df = pd.DataFrame(data)
10 print(df)
```

输出：

	StudentID	Name	Age
0	1	Alice	20
1	2	Bob	21
2	3	Charlie	22
3	4	David	23
4	5	Eva	24

现在，我们可以使用 `sample()` 方法随机抽取一定数量的样本：

```
1 # 抽取 2 个样本
2 sample1 = df.sample(n=2)
3 print(sample1)
```

输出：

```
1      StudentID      Name  Age
2  3              4    David   23
3  2              3  Charlie   22
```

你也可以指定抽样的比例：

```
1 # 抽取 50% 的样本
2 sample2 = df.sample(frac=0.5)
3 print(sample2)
```

输出：

```
1      StudentID      Name  Age
2  0              1   Alice   20
3  4              5     Eva   24
```

如果你希望允许重复抽样，可以设置 `replace=True`：

```
1 # 允许重复抽样，抽取 3 个样本
2 sample3 = df.sample(n=3, replace=True)
3 print(sample3)
```

输出：

```
1      StudentID      Name  Age
2  2              3  Charlie   22
3  3              4    David   23
4  2              3  Charlie   22
```

生成统计摘要

```
1 # 生成统计摘要
2 df.describe()
```

	Sales	Quantity
count	6.000000	6.000000
mean	186.666667	13.833333
std	77.888810	4.665476

	Sales	Quantity
min	100.000000	8.000000
25%	127.500000	10.500000
50%	175.000000	13.500000
75%	237.500000	17.250000
max	300.000000	20.000000