

Pytorch

Pandas

Create dataframe

```
1 import pandas as pd
2 import numpy as np
3
4 s = pd.Series([1,3,6,np.nan,44,1])
5
6 dates = pd.date_range('20200101',periods=6)
7
8 # Create a 6x4 dataframe
9 df = pd.DataFrame(np.random.randn(6,4),index=dates,column=['a','b','c','d'])
10
11 # Use dict to create DataFrame
12 df2 = pd.DataFrame({'A' : 1.,
13                    'B' : pd.Timestamp('20130102'),
14                    'C' : pd.Series(1,index=list(range(4)),dtype='float32'),
15                    'D' : np.array([3] * 4,dtype='int32'),
16                    'E' : pd.Categorical(["test","train","test","train"]),
17                    'F' : 'foo'})
```

Attributes of DataFrame

```
1 df2.dtypes # return the type of DataFrame
2
3 df2.index # return the row index of the DataFrame
4
5 df2.columns #return the column index of the DataFrame
6
7 df2.values # return the elements of the DataFrame
8
9 df2.describe() # return the info of the DataFrame
10
11 df2.T # return the transpose DataFrame
12
13 # axis=1 sort the DataFrame based on column
14 # ascending=False, Reverse order
15 df2.sort_index(axis=1,ascending=False))
16
17 # sort the DataFrame based on specific column
18 df2.sort_values(by='E')
```

DataFrame slicing

```
1 import pandas as pd
2 import numpy as np
3
4 dates = pd.date_range('20200101',periods=6)
5 df = pd.DataFrame(np.arange(24).reshape((6,4)),index=dates,columns=['A','B','C','D'])
6
7 # print column A
8 print(['A'])
9 print(df.A)
10
11 # print 0-3 lines
12 print(df[0:3])
13 print(df['20200101':'20200103'])
```

Select by label: loc

```
1 # print specific element by label
2 print(df.loc['20200101'])
3
4 # print all the element in A to C columns
5 print(df.loc[:, 'A':'C'])
6
7 # print the element with '20200103' label and in A/B columns
8 print(df.loc['20200103', ['A', 'B']])
```

Select by position: iloc

```
1 # print the element which in the 3 line and 1 column
2 print(df.iloc[3,1])
3
4 print(df.iloc[3:5,1:3])
```

Boolean indexing

```
1 # print the element when it is larger than 8
2 print(df[df.A>8])
```

Assign values to elements

```
1 df.iloc[2,2] = 111
2
3 df.loc['20200102','A'] = 222
4
5 # When the element in column A is less than 4,
6 # change the corresponding element in column B to 0
7 df.B[df.A<4] = 0
8
9 # Add a new column all the elements are NaN
10 df['F'] = np.nan
11 df['E'] = [1,2,3,4,5,6]
```

Error value handling

```
1 df.iloc[0,1] = np.nan
2 df.iloc[1,2] = np.nan
3
4 # Delete line with NaN
5 df.dropna(axis=0,how='any')
6
7 df.dropna(axis=1)
8
9 # Use specific value to replace the NaN
10 df.fillna(value=0)
11
12 # Use True and False to replace the element in the matrix
13 # If the element is not NaN use False to replace it
14 # otherwise use True
15 df.isnull()
16
17 # Judge whether there is any NaN in the matrix
18 df.any(df.isnull()) ==True
```

Read and save

```
1 data = pd.read_csv('student.csv')
2
3 data.to_pickle('student.pickle')
```

Concatenate

```
1 df1 = pd.DataFrame(np.ones((3,4))*0,columns=['a','b','c','d'])
2 df2 = pd.DataFrame(np.ones((3,4))*1,columns=['a','b','c','d'])
3 df3 = pd.DataFrame(np.ones((3,4))*2,columns=['a','b','c','d'])
4
5 # axis=0 change the lines, ignore_index=True means use new index
6 res = pd.concat([df1,df2,df3],axis=0,ignore_index = True)
7
8
9
10 df1 = pd.DataFrame(np.ones((3,4))*0,columns=['a','b','c','d'],index=[1,2,3])
11 df2 = pd.DataFrame(np.ones((3,4))*1,columns=['b','c','d','e'],index=[2,3,4])
12
13 # The default setting for join is outer, which will use NaN to fill the lo
14 res = pd.concat([df1,df2])
15
16 # inner: only keep the same part of the two DataFrame
17 res = pd.concat([df1,df2],axis=0, join='inner',ignore_index=False)
```

Append

```
1 df1 = pd.DataFrame(np.ones((3,4))*0,columns=['a','b','c','d'],index=[1,2,3])
2 df2 = pd.DataFrame(np.ones((3,4))*1,columns=['a','b','c','d'],index=[2,3,4])
3 df3 = pd.DataFrame(np.ones((3,4))*2,columns=['a','b','c','d'],index=[2,3,4])
4
5 df1.append([df2,df3],ignore_index=True)
6
7 s1=pd.Series([1,2,3,4],index=['a','b','c','d'])
8
9 res=df1.append(s1,ignore_index=True)
```

Merge

```
1 left = pd.DataFrame({'key':['K0','K1','K2','K3'],
2                       'A':['A0','A1','A2','A3'],
3                       'B':['B0','B1','B2','B3']})
4 right = pd.DataFrame({'key':['K0','K1','K2','K3'],
5                       'C':['C0','C1','C2','C3'],
6                       'D':['D0','D1','D2','D3']})
7
8 # merge left and right dataframe based on 'key' column
9
10 res = pd.merge(left,right,on='key')
```

how = ['left','right','inner','outer'], the default value of how = 'inner' : merge the same part

```
1 left = pd.DataFrame({'key1': ['K0', 'K0', 'K1', 'K2'],
2                       'key2': ['K0', 'K1', 'K0', 'K1'],
3                       'A': ['A0', 'A1', 'A2', 'A3'],
4                       'B': ['B0', 'B1', 'B2', 'B3']})
5 right = pd.DataFrame({'key1': ['K0', 'K1', 'K1', 'K2'],
6                       'key2': ['K0', 'K0', 'K0', 'K0'],
7                       'C': ['C0', 'C1', 'C2', 'C3'],
8                       'D': ['D0', 'D1', 'D2', 'D3']})
9 res = pd.merge(left,right,how='inner',on=['key1','key2'])
10
11
12
13 # how = 'outer' merge all part use NaN to fill the missing part
14
15 res = pd.merge(left,right,how='outer',on=['key1','key2'])
16
17 # how = 'right' merge two DataFrame based on the right key,
18 # use NaN to fill the missing part
19
20 res = pd.merge(left,right,how='right',on=['key1','key2'])
```

Indicator: show how to merge

```
1 # Indicator: show how to merge
2
```

```
3 df1 = pd.DataFrame({'col1':[0,1], 'col_left':['a','b']})
4 df2 = pd.DataFrame({'col1':[1,2,2], 'col_right':[2,2,2]})
5
6 res = pd.merge(df1,df2,on='col1',how='outer',indicator=True)
```

Index

```
1 left = pd.DataFrame({'A': ['A0', 'A1', 'A2'],
2                        'B': ['B0', 'B1', 'B2']},
3                      index=['K0', 'K1', 'K2'])
4 right = pd.DataFrame({'C': ['C0', 'C2', 'C3'],
5                          'D': ['D0', 'D2', 'D3']},
6                       index=['K0', 'K2', 'K3'])
7
8 res=pd.merge(left,right,left_index=True,right_index=True,how='outer')
```

solve overlapping problem

```
1 boys = pd.DataFrame({'k': ['K0', 'K1', 'K2'], 'age': [1, 2, 3]})
2 girls = pd.DataFrame({'k': ['K0', 'K0', 'K3'], 'age': [4, 5, 6]})
3
4 print(boys)
5 print(girls)
6 res = pd.merge(boys,girls,on='k',suffixes=['_boy','_girl'],how='outer')
7 print(res)
```