

# Pandas Tutorial

## Import Libraries

```
In [ ]: import pandas as pd
import numpy as np
```

## Object Series

```
In [ ]: series1 = pd.Series([1,2,3,4,5,6,7,8,9,10])
series1

Out [ ]:
0    1
1    2
2    3
3    4
4    5
5    6
6    7
7    8
8    9
9   10
dtype: int64

In [ ]: dates = pd.date_range("2022-01-01", periods=20)
dates

Out [ ]: DatetimeIndex(['2022-01-01', '2022-01-02', '2022-01-03', '2022-01-04',
                  '2022-01-05', '2022-01-06', '2022-01-07', '2022-01-08',
                  '2022-01-09', '2022-01-10', '2022-01-11', '2022-01-12',
                  '2022-01-13', '2022-01-14', '2022-01-15', '2022-01-16',
                  '2022-01-17', '2022-01-18', '2022-01-19', '2022-01-20'],
                  dtype='datetime64[ns]', freq='D')
```

## Creating Data-Frame

```
In [ ]: df = pd.DataFrame(np.random.randn(20,4), index=dates, columns=list("ABCD"))
df

Out [ ]:
           A         B         C         D
2022-01-01 -0.243597 -0.807238  0.566472  0.407585
2022-01-02 -0.453285  0.501168  0.098791  0.060577
2022-01-03 -0.542059 -1.559297 -1.319488 -1.042029
2022-01-04  0.603553 -1.062545 -1.369933 -0.771599
2022-01-05  1.078564  1.351742  0.813419 -0.246128
2022-01-06  0.031488 -1.328991 -0.602976 -1.109899
2022-01-07  0.836313  1.767735 -0.452541 -0.938514
2022-01-08 -0.964943 -0.665613 -0.248710  0.096500
2022-01-09 -0.347374 -0.275315 -0.564968  0.811769
2022-01-10  2.146198  0.227445  0.366495  1.291886
2022-01-11  0.045809  0.226590 -0.445373 -0.786563
2022-01-12  1.017386 -1.550262  0.308897 -0.369246
2022-01-13  0.358586  0.417658  0.886548  0.847326
2022-01-14 -0.621551  0.848705 -0.348097 -0.671176
2022-01-15  1.061288 -0.765717  1.715390  0.468456
2022-01-16  1.119669  1.251402  0.203521  1.244519
2022-01-17 -0.363791 -0.470904  0.171403 -0.472825
2022-01-18  0.340325 -1.526882  0.285516  0.144267
2022-01-19 -0.410259  0.400678  1.087062 -1.059641
2022-01-20 -0.530977  2.066283  0.250032  0.172869
```

```
In [ ]: df2 = pd.DataFrame({
    "A": 1.0,
    "B": pd.Timestamp("2022-01-01"),
    "C": pd.Series(1, index=list(range(4)), dtype='float32'),
    "D": np.array([1,3,4], dtype='int32'),
    "E": pd.Categorical(["cat", "dog", "dog", "cat"]),
    "F": "female",
})
df2

Out [ ]:
           A         B         C         E         F
0  1.0 2022-01-01  1.0  3    cat  female
1  1.0 2022-01-01  1.0  3  women  female
2  1.0 2022-01-01  1.0  3    cat  female
3  1.0 2022-01-01  1.0  3  women  female

In [ ]: df2.dtypes

Out [ ]:
A          float64
B    datetime64[ns]
C              int32
D              int32
E      Categorical
F      object
dtype: object

head() & tail() Functions
```

```
In [ ]: df.head(4)

Out [ ]:
           A         B         C         D
2022-01-01 -0.243597 -0.807238  0.566472  0.407585
2022-01-02 -0.453285  0.501168  0.098791  0.060577
2022-01-03 -0.542059 -1.559297 -1.319488 -1.042029
2022-01-04  0.603553 -1.062545 -1.369933 -0.771599
```

## Index()

```
In [ ]: df.index

Out [ ]: DatetimeIndex(['2022-01-01', '2022-01-02', '2022-01-03', '2022-01-04',
                  '2022-01-05', '2022-01-06', '2022-01-07', '2022-01-08',
                  '2022-01-09', '2022-01-10', '2022-01-11', '2022-01-12',
                  '2022-01-13', '2022-01-14', '2022-01-15', '2022-01-16',
                  '2022-01-17', '2022-01-18', '2022-01-19', '2022-01-20'],
                  dtype='datetime64[ns]', freq='D')
```

## Converting Data Frame into Numpy Array using numpy()

```
In [ ]: df.to_numpy()

Out [ ]: array([[ -0.24359652, -0.80723815,  0.56647183,  0.40758465],
       [ -0.4532853 ,  0.50116843,  0.09878993,  0.06057657],
       [ -0.54205886, -1.55929732, -1.31948775, -1.0420289 ],
       [  0.60355305, -1.06254535, -1.36993386, -0.77159945],
       [  1.07856413,  1.35174352,  0.81341914, -0.24612845],
       [  0.03148845, -1.32898995, -0.60297452, -1.10988962],
       [  0.83631343,  1.767735  , -0.45254133, -0.93851375],
       [  0.96494395, -0.66561329, -0.24871002,  0.09650052],
       [  0.37347277, -0.27531453, -0.56496882,  0.81776875],
       [  2.14619839,  0.22744577,  0.36649585,  1.29188602],
       [  0.04580897,  0.22659879,  0.44537331, -0.78656264],
       [  1.01738646, -1.55926223,  0.30889725, -0.36924591],
       [  0.35858593,  0.41765896,  0.88654817,  0.84732543],
       [ -0.62155861,  0.84878551, -0.34809732, -0.67176851],
       [  1.0612881 , -0.76571719,  1.71539395,  0.46845067],
       [  1.11966891,  1.25146439,  0.203521  ,  1.24451943],
       [  0.36379149, -0.47898382,  0.17148343, -0.4782478 ],
       [  0.34032525, -1.52688237,  0.28551651,  0.14426651],
       [ -0.41025986,  0.40067775,  1.08706189, -1.05964131],
       [ -0.53097709,  2.06628323,  0.25003195,  0.17286856]])

describe() this function describe about data Mean-Median-Mode
```

```
In [ ]: df.describe()

Out [ ]:
           A         B         C         D
count  20.000000  20.000000  20.000000  20.000000
mean   0.008807  0.012107  0.009070  0.001043
std    0.901647  1.097864  0.764982  0.694549
min    -0.964943  -1.559297  -1.369933  -1.059641
25%    -0.421016  -0.770097  0.187462  -0.528413
50%    0.038459  0.046048  0.187462  -0.021836
75%    0.881582  0.588053  0.411489  0.422903
max     2.146198  2.066283  0.715390  1.291886
```

## How to Transpose Data (Row to Column and Column to Row)

```
In [ ]: df.T

Out [ ]:
           0         1         2         3
A  2022-01-01 00:00:00 2022-01-01 00:00:00 2022-01-01 00:00:00 2022-01-01 00:00:00
B           1.0         1.0         1.0         1.0
C           1.0         1.0         1.0         1.0
D           3         3         3         3
E           cat         women         cat         women
F           female        female        female        female

In [ ]: # using Axis=0 data rows values will be in ascending form
df.sort_index(axis=0,ascending=True)

Out [ ]:
           A         B         C         D
2022-01-01 -0.243597 -0.807238  0.566472  0.407585
2022-01-02 -0.453285  0.501168  0.098791  0.060577
2022-01-03 -0.542059 -1.559297 -1.319488 -1.042029
2022-01-04  0.603553 -1.062545 -1.369933 -0.771599
2022-01-05  1.078564  1.351742  0.813419 -0.246128
2022-01-06  0.031488 -1.328991 -0.602976 -1.109899
2022-01-07  0.836313  1.767735 -0.452541 -0.938514
2022-01-08 -0.964943 -0.665613 -0.248710  0.096500
2022-01-09 -0.347374 -0.275315 -0.564968  0.811769
2022-01-10  2.146198  0.227445  0.366495  1.291886
2022-01-11  0.045809  0.226590 -0.445373 -0.786563
2022-01-12  1.017386 -1.550262  0.308897 -0.369246
2022-01-13  0.358586  0.417658  0.886548  0.847326
2022-01-14 -0.621551  0.848705 -0.348097 -0.671176
2022-01-15  1.061288 -0.765717  1.715390  0.468456
2022-01-16  1.119669  1.251402  0.203521  1.244519
2022-01-17 -0.363791 -0.470904  0.171403 -0.472825
2022-01-18  0.340325 -1.526882  0.285516  0.144267
2022-01-19 -0.410259  0.400678  1.087062 -1.059641
2022-01-20 -0.530977  2.066283  0.250032  0.172869
```

```
In [ ]: # use using Axis=1 data Columns values will be in ascending form
df.sort_index(axis=1,ascending=True)

Out [ ]:
           A         B         C         D
2022-01-01 -0.243597 -0.807238  0.566472  0.407585
2022-01-02 -0.453285  0.501168  0.098791  0.060577
2022-01-03 -0.542059 -1.559297 -1.319488 -1.042029
2022-01-04  0.603553 -1.062545 -1.369933 -0.771599
2022-01-05  1.078564  1.351742  0.813419 -0.246128
2022-01-06  0.031488 -1.328991 -0.602976 -1.109899
2022-01-07  0.836313  1.767735 -0.452541 -0.938514
2022-01-08 -0.964943 -0.665613 -0.248710  0.096500
2022-01-09 -0.347374 -0.275315 -0.564968  0.811769
2022-01-10  2.146198  0.227445  0.366495  1.291886
2022-01-11  0.045809  0.226590 -0.445373 -0.786563
2022-01-12  1.017386 -1.550262  0.308897 -0.369246
2022-01-13  0.358586  0.417658  0.886548  0.847326
2022-01-14 -0.621551  0.848705 -0.348097 -0.671176
2022-01-15  1.061288 -0.765717  1.715390  0.468456
2022-01-16  1.119669  1.251402  0.203521  1.244519
2022-01-17 -0.363791 -0.470904  0.171403 -0.472825
2022-01-18  0.340325 -1.526882  0.285516  0.144267
2022-01-19 -0.410259  0.400678  1.087062 -1.059641
2022-01-20 -0.530977  2.066283  0.250032  0.172869
```

## Sort by Specific Columns

```
In [ ]: df.sort_values(by="B",)
# df.sort_values(by="B",ascending=False)

Out [ ]:
           A         B         C         D
2022-01-03 -0.542059 -1.559297 -1.319488 -1.042029
2022-01-12  1.017386 -1.550262  0.308897 -0.369246
2022-01-18  0.340325 -1.526882  0.285516  0.144267
2022-01-01 -0.243597 -0.807238  0.566472  0.407585
2022-01-08 -0.964943 -0.665613 -0.248710  0.096500
2022-01-17 -0.363791 -0.470904  0.171403 -0.472825
2022-01-06  0.031488 -1.328991 -0.602976 -1.109899
2022-01-11  0.045809  0.226590 -0.445373 -0.786563
2022-01-19 -0.410259  0.400678  1.087062 -1.059641
2022-01-20 -0.530977  2.066283  0.250032  0.172869
2022-01-02 -0.453285  0.501168  0.098791  0.060577
2022-01-14 -0.621551  0.848705 -0.348097 -0.671176
2022-01-16  1.119669  1.251402  0.203521  1.244519
2022-01-05  1.078564  1.351742  0.813419 -0.246128
2022-01-07  0.836313  1.767735 -0.452541 -0.938514
2022-01-09 -0.347374 -0.275315 -0.564968  0.811769
2022-01-10  2.146198  0.227445  0.366495  1.291886
2022-01-13  0.358586  0.417658  0.886548  0.847326
2022-01-15  1.061288 -0.765717  1.715390  0.468456
2022-01-17 -0.363791 -0.470904  0.171403 -0.472825
2022-01-18  0.340325 -1.526882  0.285516  0.144267
2022-01-19 -0.410259  0.400678  1.087062 -1.059641
2022-01-20 -0.530977  2.066283  0.250032  0.172869
```

## (Filtering)

- Selecting a Specific Column

```
In [ ]: df["B"]

Out [ ]:
2022-01-01    -0.807238
2022-01-02    0.501168
2022-01-03   -1.559297
2022-01-04   -1.062545
2022-01-05   -1.369933
2022-01-06   -0.602976
2022-01-07   -0.132891
2022-01-08   -0.275315
2022-01-09   -0.665613
2022-01-10   -0.275315
2022-01-11   0.226590
2022-01-12   -1.550262
2022-01-13   0.417658
2022-01-14   0.848705
2022-01-15   -0.765717
2022-01-16   1.251482
2022-01-17   -0.470904
2022-01-18   -1.526882
2022-01-19   0.400678
2022-01-20   2.066283
Freq: D, Name: B, dtype: float64

Row Wise Selection

In [ ]: df[10:]

Out [ ]:
           A         B         C         D
2022-01-11  0.045809  0.226590 -0.445373 -0.786563
2022-01-12  1.017386 -1.550262  0.308897 -0.369246
2022-01-13  0.358586  0.417658  0.886548  0.847326
2022-01-14 -0.621551  0.848705 -0.348097 -0.671176
2022-01-15  1.061288 -0.765717  1.715390  0.468456
2022-01-16  1.119669  1.251402  0.203521  1.244519
2022-01-17 -0.363791 -0.470904  0.171403 -0.472825
2022-01-18  0.340325 -1.526882  0.285516  0.144267
2022-01-19 -0.410259  0.400678  1.087062 -1.059641
2022-01-20 -0.530977  2.066283  0.250032  0.172869
```

## Getting a specific data using index & labels Name

```
In [ ]: df.loc["2022-01-02","2022-01-06",["A","B"]]

Out [ ]:
           A         B
2022-01-02 -0.453285  0.501168
2022-01-03 -0.542059 -1.559297
2022-01-04  0.603553 -1.062545
2022-01-05  1.078564  1.351742
2022-01-06  0.031488 -1.328991
```

## df.at use for

```
In [ ]: # What is a single value for a row/column label pair. Similar to loc, in that both provide label-based lookups. Use at if you only need to get or set a single value in a DataFrame or Series.

Out [ ]:
-1.06254543469570858
```

## iLoc[]

```
In [ ]: # implicitly
df.iloc[0,5,:]

Out [ ]:
           A         B         C         D
2022-01-01 -0.243597 -0.807238  0.566472  0.407585
2022-01-02 -0.453285  0.501168  0.098791  0.060577
2022-01-03 -0.542059 -1.559297 -1.319488 -1.042029
2022-01-04  0.603553 -1.062545 -1.369933 -0.771599
2022-01-05  1.078564  1.351742  0.813419 -0.246128

In [ ]: df.iloc[:, 0:2]

Out [ ]:
           A         B
2022-01-01 -0.243597 -0.807238
2022-01-02 -0.453285  0.501168
2022-01-03 -0.542059 -1.559297
2022-01-04  0.603553 -1.062545
2022-01-05  1.078564  1.351742
2022-01-06  0.031488 -1.328991
2022-01-07  0.836313  1.767735
2022-01-08 -0.964943 -0.665613
2022-01-09 -0.347374 -0.275315
2022-01-10  2.146198  0.227445
2022-01-11  0.045809  0.226590
2022-01-12  1.017386 -1.550262
2022-01-13  0.358586  0.417658
2022-01-14 -0.621551  0.848705
2022-01-15  1.061288 -0.765717
2022-01-16  1.119669  1.251402
2022-01-17 -0.363791 -0.470904
2022-01-18  0.340325 -1.526882
2022-01-19 -0.410259  0.400678
2022-01-20 -0.530977  2.066283
```

```
In [ ]: #Selection or Filtration
df[df["A","B"]>1.5]

KeyError                                Traceback (most recent call last)
~\Anaconda3\lib\site-packages\pandas\core\indexes\base.py in get_iloc(self, key, method, tolerance)
    3590         try:
-> 3591             return self._engine.get_iloc(casted_key)
    3592         except KeyError as err:
-> 3593             raise KeyError from err

~\Anaconda3\lib\site-packages\pandas\libs\index.pyx in pandas._libs.index.IndexEngine.get_iloc()

~\Anaconda3\lib\site-packages\pandas\libs\index.pyx in pandas._libs.index.IndexEngine.get_iloc()

pandas._libs\hashtable.class_helper.pxi in pandas._libs.hashtable.pyobjecthashtable.get_iitem()

pandas._libs\hashtable.class_helper.pxi in pandas._libs.hashtable.pyobjecthashtable.get_iitem()

KeyError: ('A', 'B')

The above exception was the direct cause of the following exception:

KeyError                                Traceback (most recent call last)
~\Users\ME295-1-TAM\AppData\Local\Temp\ipykernel_3589\206511266.py in <module>
     1 #Selection or Filtration
     2 i = df[df["A","B"]>1.5]
-> 3         if self.columns.nlevels == 1:
    3459             return self._getitem(self, key)
-> 3459             return self._getitem_multilevel(key)
    3457         indexer = self.columns.get_iloc(key)
-> 3458         if is_integer(indexer):
    3460             indexer = [indexer]

~\Anaconda3\lib\site-packages\pandas\core\indexes\base.py in get_iloc(self, key, method, tolerance)
    3362         except KeyError as err:
-> 3363             raise KeyError from err
    3364
-> 3365         if is_scalar(key) and isna(key) and not self.hasnans:

KeyError: ('A', 'B')
```

```
In [ ]: #Select those values which are greater than 0 in entire data frame
df[df>0]

Out [ ]:
           A         B         C         D         E  new mean
2022-01-01    NaN      NaN      NaN      NaN      NaN
2022-01-02    NaN  0.501168  0.098791      NaN      NaN
2022-01-03    NaN      NaN      NaN      NaN      NaN
2022-01-04  0.603553      NaN      NaN      NaN      NaN
2022-01-05  1.078564  1.351742  0.813419      NaN      NaN
2022-01-06  0.031488      NaN      NaN      NaN      NaN
2022-01-07  0.836313  1.767735 -0.452541 -0.938514      NaN
2022-01-08 -0.964943 -0.665613 -0.248710  0.096500      NaN
2022-01-09 -0.347374 -0.275315 -0.564968  0.811769      NaN
2022-01-10  2.146198  0.227445  0.366495  1.291886      NaN
2022-01-11  0.045809  0.226590      NaN      NaN      NaN
2022-01-12  1.017386 -1.550262  0.308897      NaN      NaN
2022-01-13  0.358586  0.417658  0.886548  0.847326      NaN
2022-01-14 -0.621551  0.848705 -0.348097 -0.671176      NaN
2022-01-15  1.061288 -0.765717  1.715390  0.468456      NaN
2022-01-16  1.119669  1.251402  0.203521  1.244519      NaN
2022-01-17 -0.363791 -0.470904  0.171403 -0.472825      NaN
2022-01-18  0.340325 -1.526882  0.285516  0.144267      NaN
2022-01-19 -0.410259  0.400678  1.087062 -1.059641      NaN
2022-01-20 -0.530977  2.066283  0.250032  0.172869      NaN
```

```
In [ ]: df["mean"] = df[["A","B","C","D","E","new"]].mean(axis=1)
df2

Out [ ]:
           A         B         C         D         E  new mean
2022-01-01 -0.243597 -0.807238  0.566472  0.407585      NaN  1.2
2022-01-02 -0.453285  0.501168  0.098791  0.060577      NaN  2.2
2022-01-03 -0.542059 -1.559297 -1.319488 -1.042029      NaN  Three
2022-01-04  0.603553 -1.062545 -1.369933 -0.771599      NaN  Four
2022-01-05  1.078564  1.351742  0.813419 -0.246128      NaN  Five
2022-01-06  0.031488 -1.328991 -0.602976 -1.109899      NaN  One
2022-01-07  0.836313  1.767735 -0.452541 -0.938514      NaN  Two
2022-01-08 -0.964943 -0.665613 -0.248710  0.096500      NaN  Three
2022-01-09 -0.347374 -0.275315 -0.564968  0.811769      NaN  Four
2022-01-10  2.146198  0.227445  0.366495  1.291886      NaN  Five
2022-01-11  0.045809
```