



## Operations

There are lots of operations with pandas that will be really useful to you, but don't fall into any distinct category. Let's show them here in this lecture.

```
In [252]: import pandas as pd
import numpy as np
```

### 1) Viewing data

```
In [253]: titanic = pd.read_csv("train.csv")
```

**head()**

Return the first **n** rows.

- Syntax:** DataFrame.head(**n**: 'int' = 5)

```
In [254]: titanic.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cummings, Mrs. John Bradley (Florence Briggs Th.	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

**tail()**

Return the last **n** rows.

- Syntax:** DataFrame.tail(**n**: 'int' = 5)

```
In [255]: titanic.tail()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.00	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.00	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.45	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	113869	30.00	C148	C
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.75	NaN	Q

**sample()**

Return a random sample of items from an axis of object.

- Syntax:** DataFrame.sample(**n**=None, **frac**=None, **replace**=False, **weights**=None, **random\_state**=None, **axis**=None, **ignore\_index**=False)

```
In [256]: titanic.sample()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
834	835	0	3	Allum, Mr. Owen George	male	18.0	0	0	2223	8.3	NaN	S

```
In [257]: # We can choose a random sample size
titanic.sample(3)
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
443	444	1	2	Reynaldo, Ms. Encarnacion	female	28.00	0	0	230434	13.0000	NaN	S
803	804	1	3	Thomas, Master. Assad Alexander	male	0.42	0	1	2625	8.5167	NaN	C
324	325	0	3	Sage, Mr. George John Jr	male	NaN	8	2	CA 2343	69.5500	NaN	S

### 2) DataFrame Attributes

**a) shape**

Return the dimensionality of the DataFrame in the form of a tuple.

```
In [258]: titanic.shape # 891 rows and 12 columns
```

```
Out[258]: (891, 12)
```

```
In [259]: titanic.shape[0] # rows
```

```
Out[259]: 891
```

**b) size**

Return the number of elements in the DataFrame.

```
In [260]: titanic.size
```

```
Out[260]: 10692
```

**c) dtypes**

Return the dtypes of each column in the DataFrame

```
In [261]: titanic.dtypes
```

```
Out[261]: PassengerId    int64
Survived              int64
Name                  object
Sex                   object
Age                   float64
SibSp                 int64
Parch                 int64
Ticket                float64
Fare                  float64
Cabin                 object
Embarked              object
dtype: object
```

**d) columns**

This attribute is used to fetch the label values for columns present in a particular data frame.

```
In [262]: titanic.columns
```

```
Out[262]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],
dtype='object')
```

**e) index**

the index attribute is used to display the row labels of a data frame object.

```
In [263]: titanic.index
```

```
Out[263]: RangeIndex(start=0, stop=891, step=1)
```

**f) empty**

check whether the data frame is empty or not. This attribute returns true if the data frame is empty and false if the DataFrame is not empty.

```
In [264]: titanic.empty
```

```
Out[264]: False
```

**g) values**

This attribute is used to represent the values/data of dataframe in NumPy array form.

```
In [265]: titanic.values # The return value is a 2-dimensional array with one array for each row.
```

```
Out[265]: array([[0, 0, 3, ..., 7.25, nan, 'S'],
[2, 1, 1, ..., 71.2833, 'C85', 'C'],
[3, 1, 3, ..., 7.925, nan, 'S'],
...,
[889, 0, 3, ..., 23.45, nan, 'S'],
[890, 1, 1, ..., 30.6, 'C148', 'C'],
[891, 0, 3, ..., 7.75, nan, 'Q']], dtype=object)
```

### 3) DataFrame methods

**info()**

Print a concise summary of a DataFrame.

This method prints information about a DataFrame including the **index dtype** and **columns**, **non-null values** and **memory usage**.

```
In [266]: titanic.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   PassengerId           891 non-null    int64
1   Survived              891 non-null    int64
2   Pclass                891 non-null    int64
3   Name                  891 non-null    object
4   Sex                   891 non-null    object
5   Age                   714 non-null    float64
6   SibSp                 891 non-null    int64
7   Parch                891 non-null    int64
8   Ticket                891 non-null    object
9   Fare                  891 non-null    float64
10  Cabin                 204 non-null    object
11  Embarked              889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

**describe()**

Generate descriptive statistics.

Descriptive statistics include those that summarize the central tendency, dispersion and shape of a dataset's distribution, excluding **NaN** values.

```
In [267]: titanic.describe()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare
count	891.000000	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
mean	446.000000	0.383838	2.308642	29.699118	0.523008	0.381594	32.004208
std	257.353842	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
min	1.000000	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
25%	223.500000	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
50%	446.000000	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
75%	668.500000	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
max	891.000000	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

**unique()**

Return unique values

```
In [268]: titanic["Embarked"].unique()
```

```
Out[268]: array(['S', 'C', 'Q', nan], dtype=object)
```

```
In [269]: # Example
titanic["Sex"].unique()
```

```
Out[269]: array(['male', 'female'], dtype=object)
```

**nunique()**

Count number of distinct elements in specified axis

```
In [270]: titanic.nunique() # on whole dataframe
```

```
Out[270]: PassengerId    891
Survived                2
Pclass                  3
Name                    891
Sex                     2
Age                     88
SibSp                   7
Parch                   7
Ticket                 681
Fare                   248
Cabin                  147
Embarked                3
dtype: int64
```

**nlargest()**

Return the first **n** rows with the largest values in **columns**, in descending order

```
In [271]: titanic["Fare"].nlargest()
```

```
Out[271]: 258    512.3292
679    512.3292
737    512.3292
27     263.0000
88     263.0000
Name: Fare, dtype: float64
```

```
In [272]: # Example
titanic["Age"].nlargest(6) # How many largest values we want can be specified.
```

```
Out[272]: 630    80.0
851    74.0
96     71.0
493    71.0
116    70.5
672    70.0
Name: Age, dtype: float64
```

**nsmlallest()**

Return the first **n** rows with the smallest values in **columns**, in ascending order.

```
In [273]: titanic["Age"].nsmlallest()
```

```
Out[273]: 803    0.42
755    0.67
469    0.75
644    0.75
78     0.83
Name: Age, dtype: float64
```

**value\_counts()**

Return a Series containing counts of unique values.

```
In [274]: titanic["Sex"].value_counts()
```

```
Out[274]: male      577
female    314
Name: Sex, dtype: int64
```

```
In [275]: # Example
titanic["Pclass"].value_counts()
```

```
Out[275]: 1      491
2      216
1      184
Name: Pclass, dtype: int64
```

```
In [276]: # Example
titanic[["Sex", "Pclass"].value_counts()
```

```
Out[276]: Sex    Pclass
male      3      347
female    3      144
male      1      122
male      2      108
female    1      94
female    2      76
dtype: int64
```

```
In [277]: titanic["Age"].value_counts(bins=4)
# Bins can be useful for going from a continuous variable to a categorical variable; instead of counting unique
```

```
Out[277]: (20.315, 40.21]    385
(0.339, 20.315]    179
(40.21, 60.105]    128
(60.105, 80.01]     22
Name: Age, dtype: int64
```

**sort\_values()**

The **sort\_values()** method sorts the DataFrame by the specified label.

```
In [278]: titanic.sort_values(by="Sex") # inplace=False by default
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
383	384	1	1	Holvenson, Mrs. Alexander Oskar (Mary Aline To..	female	35.0	1	0	11813	52.0000	NaN	S
218	219	1	1	Bazzani, Miss. Albina	female	40.0	0	0	11813	76.2917	D15	S
609	610	1	1	Shutes, Miss. Elizabeth W	female	32.0	0	0	PC 17582	153.4625	C125	C
216	217	1	3	Honkanen, Miss. Eliina	female	27.0	0	0	STON/O2. 3101283	7.9250	NaN	S
215	216	1	1	Newell, Miss. Madeleine	female	31.0	1	0	35273	113.2750	D36	C
...	...	...	...	...	...	...	...	...	...	...	...	...
371	372	0	3	Wiklund, Mr. Jakob Alfred	male	18.0	1	0	3101267	6.4958	NaN	S
372	373	0	3	Beavan, Mr. William Thomas	male	19.0	0	0	323951	8.0500	NaN	S
373	374	0	1	Ringhini, Mr. Sante	male	22.0	0	0	PC 17760	135.6333	NaN	C
360	361	0	3	Skoog, Mr. Wilhelm	male	40.0	1	4	347088	27.9000	NaN	S
890	891	0	3	Dooley, Mr. Patrick	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows x 12 columns

```
In [279]: titanic.sort_values(by="Fare", ascending=False) # Now, It will sort in descending order by fare.
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
258	259	1	1	Ward, Miss. Anna	female	35.0	0	0	PC 17755	512.3292	NaN	C
737	738	1	1	Lesurer, Mr. Gustave J	male	35.0	0	0	PC 17755	512.3292	8101	C
679	680	1	1	Cardeza, Mr. Thomas Drake Martinez	male	36.0	0	1	PC 17755	512.3292	B51 B53 B55	C
88	89	1	1	Fortune, Miss. Mabel Helen	female	23.0	3	2	19950	263.0000	C23 C25 C27	S
27	28	0	1	Fortune, Mr. Charles Alexander	male	19.0	3	2	19950	263.0000	C23 C25 C27	S
...	...	...	...	...	...	...	...	...	...	...	...	...
633	634	0	1	Parr, Mr. William Henry Marsh	male	NaN	0	0	112052	0.0000	NaN	S
413	414	0	2	Cunningham, Mr. Alfred Fleming	male	NaN	0	0	239853	0.0000	NaN	S
822	823	0	1	Reuchlin, Jonkheer, John George	male	38.0	0	0	19972	0.0000	NaN	S
732	733	0	2	Knight, Mr. Robert J	male	NaN	0	0	239855	0.0000	NaN	S
674	675	0	2	Watson, Mr. Ennis Hastings	male	NaN	0	0	239856	0.0000	NaN	S

891 rows x 12 columns

**Sort by multiple columns**

```
In [280]: titanic.sort_values(by=["Sex", "Fare"], ascending=False)
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
679	680	1	1	Cardeza, Mr. Thomas Drake Martinez	male	36.0	0	1	PC 17755	512.3292	B51 B53 B55	C
737	738	1	1	Lesurer, Mr. Gustave J	male	35.0	0	0	PC 17755	512.3292	8101	C
27	28	0	1	Fortune, Mr. Charles Alexander	male	19.0	3	2	19950	263.0000	C23 C25 C27	S
438	439	0	1	Fortune, Mr. Mark	male	64.0	1	4	19950	263.0000	C23 C25 C27	S
118	119	0	1	Baxter, Mr. Quigg Edmond	male	24.0	0	1	PC 17558	427.5028	B58 B60	C
...	...	...	...	...	...	...	...	...	...	...	...	...
367	368	1	3	Moussa, Mrs. (Mantoura Boulos)	female	NaN	0	0	2626	7.2292	NaN	C
780	781	1	3	Ayoub, Miss. Bancoura	female	13.0	0	0	2687	7.2292	NaN	C
19	20	1	3	Masselemani, Mrs. Fatima	female	NaN	0	0	2649	7.2250	NaN	C
875	876	1	3	Najib, Miss. Adele Kiamie "Jane"	female	15.0	0	0	2667	7.2250	NaN	C
654	655	0	3	Hegarty, Miss. Hanora "Nora"	female	18.0	0	0	365226	6.7500	NaN	Q

891 rows x 12 columns

**drop\_duplicates()**

Return boolean Series denoting duplicate rows.

```
In [281]: # Create a DataFrame
dupd_data = pd.DataFrame({
    "name": ['Alice Smith', 'Jones', "John", "Mary"],
    "age": [50, 40, 30, 40]
})
```

```
In [282]: dupd_data.duplicated()
```

```
Out[282]: 0    False
1    False
2    False
3     True
dtype: bool
```

```
In [283]: # select rows
dupd_data[dupd_data.duplicated()] # There are one duplicate row in this table.
```