

Lecture 4-3

Introduction to Pandas

Week 4 Friday

Miles Chen, PhD

Pandas

NumPy creates ndarrays that must contain values that are of the same data type.

Pandas creates dataframes. Each column in a dataframe is an ndarray. This allows us to have traditional tables of data where each column can be a different data type.

Important References:

<https://pandas.pydata.org/pandas-docs/stable/reference/series.html>

<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.html>

```
In [1]: import numpy as np  
import pandas as pd
```

The basic data structure in pandas is the *series*. You can construct it in a similar fashion to making a numpy array.

The command to make a Series object is

```
pd.Series(data, index=index)
```

the `index` argument is optional

```
In [2]: data = pd.Series([0.25, 0.5, 0.75, 1.0])
        print(data)
        print(type(data))
```

```
0    0.25
1    0.50
2    0.75
3    1.00
dtype: float64
<class 'pandas.core.series.Series'>
```

```
In [3]: data
```

```
Out[3]: 0    0.25
        1    0.50
        2    0.75
        3    1.00
        dtype: float64
```

The series is printed out in a table form. The type is a Pandas Series

```
In [4]: print(data.values)
```

```
[0.25 0.5  0.75 1.  ]
```

```
In [5]: print(type(data.values))
```

```
<class 'numpy.ndarray'>
```

The values attribute of the series is a numpy array.

```
In [6]: print(data.index)
```

```
RangeIndex(start=0, stop=4, step=1)
```

```
In [7]: print(type(data.index)) # the row names are known as the index
```

```
<class 'pandas.core.indexes.range.RangeIndex'>
```

You can subset a pandas series like other python objects

```
In [8]: print(data[1])
```

```
0.5
```

```
In [9]: print(type(data[1])) # when you select only one value, it simplifies the object
```

```
<class 'numpy.float64'>
```

```
In [10]: print(data[1:3])
```

```
1    0.50  
2    0.75  
dtype: float64
```

```
In [11]: print(type(data[1:3])) # slicing / selecting multiple values returns a series
```

```
<class 'pandas.core.series.Series'>
```

```
In [12]: print(data[np.array([1, 0, 1, 2])]) # You can also do fancy indexing by subsetting w/a numpy array
```

```
1    0.50  
0    0.25  
1    0.50  
2    0.75  
dtype: float64
```

```
In [13]: # Pandas uses a 0-based index by default. You may also specify the index values  
data = pd.Series([0.25, 0.5, 0.75, 1.0],  
                 index = ['a', 'b', 'c', 'd'])  
print(data)
```

```
a    0.25  
b    0.50  
c    0.75  
d    1.00  
dtype: float64
```

```
In [14]: data.values
```

```
Out[14]: array([0.25, 0.5 , 0.75, 1.  ])
```

```
In [15]: data.index
```

```
Out[15]: Index(['a', 'b', 'c', 'd'], dtype='object')
```

```
In [16]: data[1] # subset with index position
```

C:\Users\miles\AppData\Local\Temp\ipykernel_35940\2055430528.py:1: FutureWarning: Series.__getitem__ treating keys as positions is deprecated. In a future version, integer keys will always be treated as labels (consistent with DataFrame behavior). To access a value by position, use `ser.iloc[pos]`

```
data[1] # subset with index position
```

```
Out[16]: np.float64(0.5)
```

```
In [17]: data["a"] # subset with index names
```

```
Out[17]: np.float64(0.25)
```

```
In [18]: data[0:2] # slicing behavior is unchanged
```

```
Out[18]: a    0.25  
b    0.50  
dtype: float64
```

```
In [19]: data["a":"c"] # slicing using index names includes the last value
```

```
Out[19]: a    0.25  
b    0.50  
c    0.75  
dtype: float64
```

```
In [20]: # creating a series from a python dictionary  
# remember, dictionary construction uses curly braces {}  
samp_dict = {'Tony Stark': "Robert Downey Jr.",  
             'Steve Rogers': "Chris Evans",  
             'Natasha Romanoff': "Scarlett Johansson",  
             'Bruce Banner': "Mark Ruffalo",  
             'Thor': "Chris Hemsworth",  
             'Clint Barton': "Jeremy Renner"}  
samp_series = pd.Series(samp_dict)  
samp_series
```

```
Out[20]: Tony Stark           Robert Downey Jr.  
         Steve Rogers         Chris Evans  
         Natasha Romanoff     Scarlett Johansson  
         Bruce Banner         Mark Ruffalo  
         Thor                 Chris Hemsworth  
         Clint Barton         Jeremy Renner  
         dtype: object
```

```
In [21]: print(samp_series.index) # dtype = object is for strings but allows mixed data types.
```

```
Index(['Tony Stark', 'Steve Rogers', 'Natasha Romanoff', 'Bruce Banner',
      'Thor', 'Clint Barton'],
      dtype='object')
```

```
In [22]: samp_series.values
```

```
Out[22]: array(['Robert Downey Jr.', 'Chris Evans', 'Scarlett Johansson',  
                'Mark Ruffalo', 'Chris Hemsworth', 'Jeremy Renner'], dtype=object)
```

```
In [23]: # ages during the First Avengers film (2012)
# I don't have an exact source, don't get mad at me.
age_dict = {'Thor': 1493,
            'Steve Rogers': 104,
            'Natasha Romanoff': 28,
            'Clint Barton': 41,
            'Tony Stark': 42,
            'Bruce Banner': 42} # note that the dictionary order is not same here
ages = pd.Series(age_dict)
print(ages)
```

```
Thor          1493
Steve Rogers   104
Natasha Romanoff  28
Clint Barton   41
Tony Stark     42
Bruce Banner   42
dtype: int64
```

```
In [24]: # Super Hero Names
hero_dict = {'Thor': np.nan,
             'Steve Rogers': 'Captain America',
             'Natasha Romanoff': 'Black Widow'}
```

```

        'Clint Barton': 'Hawkeye',
        'Tony Stark': 'Iron Man',
        'Bruce Banner': 'Hulk'}
hero_names = pd.Series(hero_dict)
print(hero_names)

```

```

Thor                NaN
Steve Rogers        Captain America
Natasha Romanoff    Black Widow
Clint Barton        Hawkeye
Tony Stark           Iron Man
Bruce Banner        Hulk
dtype: object

```

Creating a DataFrame

There are multiple ways of creating a DataFrame in Pandas. The next few slides show a few.

We can create a dataframe by providing a dictionary of series objects. The dictionary key becomes the column name. The dictionary values become values. The keys within the dictionaries become the index.

```

In [25]: avengers = pd.DataFrame({'actor': samp_series,
                                'hero name': hero_names,
                                'age': ages})

# the DataFrame will match the indices and sort them
print(avengers)

```

	actor	hero name	age
Bruce Banner	Mark Ruffalo	Hulk	42
Clint Barton	Jeremy Renner	Hawkeye	41
Natasha Romanoff	Scarlett Johansson	Black Widow	28
Steve Rogers	Chris Evans	Captain America	104
Thor	Chris Hemsworth	NaN	1493
Tony Stark	Robert Downey Jr.	Iron Man	42

```

In [26]: print(type(avengers)) # this is a DataFrame object

```

```
<class 'pandas.core.frame.DataFrame'>
```

The data is a list of dictionaries. Each dictionary needs to have the same set of keys, otherwise, NaNs will appear.

```
In [27]: data = [{'a': 0, 'b': 0},
                 {'a': 1, 'b': 2},
                 {'a': 2, 'b': 5}]
data
```

```
Out[27]: [{'a': 0, 'b': 0}, {'a': 1, 'b': 2}, {'a': 2, 'b': 5}]
```

```
In [28]: print(pd.DataFrame(data, index = [1, 2, 3]))
```

```
   a  b
1  0  0
2  1  2
3  2  5
```

```
In [29]: data2 = [{'a': 0, 'b': 0},
                  {'a': 1, 'b': 2},
                  {'a': 2, 'c': 5}] # mismatch of keys. NAs will appear
data2
```

```
Out[29]: [{'a': 0, 'b': 0}, {'a': 1, 'b': 2}, {'a': 2, 'c': 5}]
```

```
In [30]: pd.DataFrame(data2) # if the index argument is not supplied, it defaults to integer index start at 0
```

```
Out[30]:
```

	a	b	c
0	0	0.0	NaN
1	1	2.0	NaN
2	2	NaN	5.0

You can convert a dictionary to a DataFrame. The keys form column names, and the values are lists/arrays of values. The arrays need to be of the same length.

```
In [31]: data3 = {'a': [1, 2, 3],
                  'b': ['x', 'y', 'z']}
data3
```

```
Out[31]: {'a': [1, 2, 3], 'b': ['x', 'y', 'z']}
```

```
In [32]: pd.DataFrame(data3)
```

```
Out[32]:
```

	a	b
0	1	x
1	2	y
2	3	z

```
In [33]: data4 = {'a': [1, 2, 3, 4],  
                  'b': ['x', 'y', 'z']} # arrays are not of the same length  
pd.DataFrame(data4)
```

ValueError

Traceback (most recent call last)

Cell In[33], line 3

```
1 data4 = {'a': [1, 2, 3, 4],
2         'b': ['x', 'y', 'z']} # arrays are not of the same length
----> 3 pd.DataFrame(data4)
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\frame.py:778, in DataFrame.__init__(self, data, index, columns, dtype, copy)

```
772 mgr = self._init_mgr(
773     data, axes={"index": index, "columns": columns}, dtype=dtype, copy=copy
774 )
776 elif isinstance(data, dict):
777     # GH#38939 de facto copy defaults to False only in non-dict cases
--> 778 mgr = dict_to_mgr(data, index, columns, dtype=dtype, copy=copy, typ=manager)
779 elif isinstance(data, ma.MaskedArray):
780     from numpy.ma import mrecords
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\internals\construction.py:503, in dict_to_mgr(data, index, columns, dtype, typ, copy)

```
499 else:
500     # dtype check to exclude e.g. range objects, scalars
501     arrays = [x.copy() if hasattr(x, "dtype") else x for x in arrays]
--> 503 return arrays_to_mgr(arrays, columns, index, dtype=dtype, typ=typ, consolidate=copy)
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\internals\construction.py:114, in arrays_to_mgr(arrays, columns, index, dtype, verify_integrity, typ, consolidate)

```
111 if verify_integrity:
112     # figure out the index, if necessary
113     if index is None:
--> 114         index = _extract_index(arrays)
115 else:
116     index = ensure_index(index)
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\internals\construction.py:677, in _extract_index(data)

```
675 lengths = list(set(raw_lengths))
676 if len(lengths) > 1:
--> 677     raise ValueError("All arrays must be of the same length")
679 if have_dicts:
680     raise ValueError(
681         "Mixing dicts with non-Series may lead to ambiguous ordering."
```

```
682 )
```

ValueError: All arrays must be of the same length

Turn a 2D Numpy array (matrix) into a DataFrame by adding column names and optionally index values.

```
In [34]: data = np.random.randint(10, size = 10).reshape((5,2))
print(data)
```

```
[[8 6]
 [5 7]
 [9 0]
 [8 5]
 [3 2]]
```

```
In [35]: print(pd.DataFrame(data, columns = ["x","y"], index = ['a','b','c','d','e']))
```

```
   x  y
a  8  6
b  5  7
c  9  0
d  8  5
e  3  2
```

Subsetting the DataFrame

In a DataFrame, the `.columns` attribute show the column names and the `.index` attribute show the row names.

```
In [36]: print(avengers)
```

	actor	hero name	age
Bruce Banner	Mark Ruffalo	Hulk	42
Clint Barton	Jeremy Renner	Hawkeye	41
Natasha Romanoff	Scarlett Johansson	Black Widow	28
Steve Rogers	Chris Evans	Captain America	104
Thor	Chris Hemsworth	NaN	1493
Tony Stark	Robert Downey Jr.	Iron Man	42

```
In [37]: print(avengers.columns)
```

```
Index(['actor', 'hero name', 'age'], dtype='object')
```

```
In [38]: print(avengers.index)
```

```
Index(['Bruce Banner', 'Clint Barton', 'Natasha Romanoff', 'Steve Rogers',  
      'Thor', 'Tony Stark'],  
      dtype='object')
```

You can select a column using dot notation or with single square brackets.

```
In [39]: avengers.actor # extracting the column
```

```
Out[39]: Bruce Banner      Mark Ruffalo  
        Clint Barton      Jeremy Renner  
        Natasha Romanoff  Scarlett Johansson  
        Steve Rogers      Chris Evans  
        Thor              Chris Hemsworth  
        Tony Stark        Robert Downey Jr.  
        Name: actor, dtype: object
```

```
In [40]: avengers["hero name"] # if there's a space in the column name, you'll need to use square brackets
```

```
Out[40]: Bruce Banner      Hulk  
        Clint Barton      Hawkeye  
        Natasha Romanoff  Black Widow  
        Steve Rogers      Captain America  
        Thor              NaN  
        Tony Stark        Iron Man  
        Name: hero name, dtype: object
```

```
In [41]: type(avengers.actor)
```

```
Out[41]: pandas.core.series.Series
```

The selected column is a Pandas Series and can be subset accordingly.

```
In [42]: avengers.actor[1] # 0 based indexing
```

```
C:\Users\miles\AppData\Local\Temp\ipykernel_35940\438569199.py:1: FutureWarning: Series.__getitem__ treating keys as positions is deprecated. In a future version, integer keys will always be treated as labels (consistent with DataFrame behavior). To access a value by position, use `ser.iloc[pos]`  
avengers.actor[1] # 0 based indexing
```

```
Out[42]: 'Jeremy Renner'
```

```
In [43]: avengers.actor[avengers.age == 42]
```

```
Out[43]: Bruce Banner      Mark Ruffalo  
         Tony Stark      Robert Downey Jr.  
         Name: actor, dtype: object
```

```
In [44]: avengers["hero name"]['Steve Rogers']
```

```
Out[44]: 'Captain America'
```

```
In [45]: avengers["hero name"]['Steve Rogers':'Tony Stark']
```

```
Out[45]: Steve Rogers      Captain America  
         Thor              NaN  
         Tony Stark      Iron Man  
         Name: hero name, dtype: object
```

.loc

The `.loc` attribute can be used to subset the DataFrame using the index names.

```
In [46]: avengers.loc['Thor'] # subset based on location to get a row
```

```
Out[46]: actor      Chris Hemsworth  
         hero name      NaN  
         age      1493  
         Name: Thor, dtype: object
```

```
In [47]: print(type(avengers.loc['Thor']))  
         print(type(avengers.loc['Thor'].values)) # the values are of mixed type but is still a numpy array.  
         # this is possible because it is a structured numpy array. (covered in "Python for Data Science" chapter 2)
```

```
<class 'pandas.core.series.Series'>  
<class 'numpy.ndarray'>
```

```
In [48]: print(avengers.loc[:, 'age']) # subset based on location to get a column
```

```
Bruce Banner      42
Clint Barton      41
Natasha Romanoff  28
Steve Rogers      104
Thor              1493
Tony Stark        42
Name: age, dtype: int64
```

```
In [49]: print(type(avengers.loc[:, 'age'])) #the object is a pandas series
        print(type(avengers.loc[:, 'age'].values))
```

```
<class 'pandas.core.series.Series'>
<class 'numpy.ndarray'>
```

```
In [50]: avengers.loc['Steve Rogers', 'age'] # you can provide a pair of 'coordinates' to get a particular value
```

```
Out[50]: np.int64(104)
```

.iloc

The `.iloc` attribute can be used to subset the DataFrame using the index position (zero-indexed).

```
In [51]: avengers.iloc[3,] # subset based on index location
```

```
Out[51]: actor      Chris Evans
        hero name    Captain America
        age          104
        Name: Steve Rogers, dtype: object
```

```
In [52]: avengers.iloc[0, 1] # pair of coordinates
```

```
Out[52]: 'Hulk'
```

Assignment with `.loc` and `.iloc`

The `.loc` and `.iloc` attributes can be used in conjunction with assignment.

```
In [53]: # set values individually
avengers.loc['Thor', 'age'] = 1500
avengers.loc['Thor', 'hero name'] = 'Thor'
avengers
```

```
Out[53]:
```

	actor	hero name	age
Bruce Banner	Mark Ruffalo	Hulk	42
Clint Barton	Jeremy Renner	Hawkeye	41
Natasha Romanoff	Scarlett Johansson	Black Widow	28
Steve Rogers	Chris Evans	Captain America	104
Thor	Chris Hemsworth	Thor	1500
Tony Stark	Robert Downey Jr.	Iron Man	42

```
In [54]: # assign multiple values at once
avengers.loc['Thor', ['hero name', 'age']] = [np.nan, 1493]
avengers
```

```
Out[54]:
```

	actor	hero name	age
Bruce Banner	Mark Ruffalo	Hulk	42
Clint Barton	Jeremy Renner	Hawkeye	41
Natasha Romanoff	Scarlett Johansson	Black Widow	28
Steve Rogers	Chris Evans	Captain America	104
Thor	Chris Hemsworth	NaN	1493
Tony Stark	Robert Downey Jr.	Iron Man	42

.loc vs **.iloc** with numeric index

The following DataFrame has a numeric index, but it starts at 1 instead of 0.

```
In [55]: data = [{'a': 11, 'b': 2},  
                {'a': 12, 'b': 4},  
                {'a': 13, 'b': 6}]  
df = pd.DataFrame(data, index = [1, 2, 3])  
df
```

```
Out[55]:
```

	a	b
1	11	2
2	12	4
3	13	6

```
In [56]: df.loc[1, :] # .loc always uses the actual index.
```

```
Out[56]: a    11  
         b     2  
         Name: 1, dtype: int64
```

```
In [57]: df.iloc[1, :] # .iloc always uses the position using a 0-based index.
```

```
Out[57]: a    12  
         b     4  
         Name: 2, dtype: int64
```

```
In [58]: df.iloc[3, :] # using a position that doesn't exist results in an exception.
```

IndexError

Traceback (most recent call last)

Cell In[58], line 1

```
----> 1 df.iloc[3, :] # using a position that doesn't exist results in an exception.
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\indexing.py:1184, in [_LocationIndexer.__getitem__](#)(self, key)

```
    1182     if self._is_scalar_access(key):
    1183         return self.obj._get_value(*key, takeable=self._takeable)
-> 1184     return self._getitem_tuple(key)
    1185 else:
    1186     # we by definition only have the 0th axis
    1187     axis = self.axis or 0
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\indexing.py:1690, in [_iLocIndexer._getitem_tuple](#)(self, tup)

```
    1689 def _getitem_tuple(self, tup: tuple):
-> 1690     tup = self._validate_tuple_indexer(tup)
    1691     with suppress(IndexingError):
    1692         return self._getitem_lowerdim(tup)
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\indexing.py:966, in [_LocationIndexer._validate_tuple_indexer](#)(self, key)

```
    964 for i, k in enumerate(key):
    965     try:
--> 966         self._validate_key(k, i)
    967     except ValueError as err:
    968         raise ValueError(
    969             "Location based indexing can only have "
    970             f"[{self._valid_types}] types"
    971         ) from err
```

File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\indexing.py:1592, in [_iLocIndexer._validate_key](#)(self, key, axis)

```
    1590     return
    1591 elif is_integer(key):
-> 1592     self._validate_integer(key, axis)
    1593 elif isinstance(key, tuple):
    1594     # a tuple should already have been caught by this point
    1595     # so don't treat a tuple as a valid indexer
    1596     raise IndexingError("Too many indexers")
```



```
File c:\Users\miles\.pyenv\pyenv-win\versions\3.12.5\Lib\site-packages\pandas\core\indexing.py:1685, in _iLocIndexer._validate_integer(self, key, axis)
    1683 len_axis = len(self.obj._get_axis(axis))
    1684 if key >= len_axis or key < -len_axis:
-> 1685     raise IndexError("single positional indexer is out-of-bounds")

IndexError: single positional indexer is out-of-bounds
```

Boolean subsetting examples with `.loc`

```
In [59]: # select avengers whose age is less than 50 and greater than 40
# select the columns 'hero name' and 'age'
avengers.loc[ (avengers.age < 50) & (avengers.age > 40), ['hero name', 'age']]
```

```
Out[59]:
```

	hero name	age
Bruce Banner	Hulk	42
Clint Barton	Hawkeye	41
Tony Stark	Iron Man	42

```
In [60]: # Use the index of the DataFrame, treat it as a string, and select rows that start with B
avengers.loc[ avengers.index.str.startswith('B'), : ]
```

```
Out[60]:
```

	actor	hero name	age
Bruce Banner	Mark Ruffalo	Hulk	42

```
In [61]: # Use the index of the DataFrame, treat it as a string,
# find the character capital R. Find returns -1 if it does not find the letter
# We select rows that did not result in -1, which means it does contain a capital R
avengers.loc[ avengers.index.str.find('R') != -1, : ]
```

```
Out[61]:
```

	actor	hero name	age
Natasha Romanoff	Scarlett Johansson	Black Widow	28
Steve Rogers	Chris Evans	Captain America	104

Other commonly used DataFrame attributes

```
In [62]: avengers.T # the transpose
```

```
Out[62]:
```

	Bruce Banner	Clint Barton	Natasha Romanoff	Steve Rogers	Thor	Tony Stark
actor	Mark Ruffalo	Jeremy Renner	Scarlett Johansson	Chris Evans	Chris Hemsworth	Robert Downey Jr.
hero name	Hulk	Hawkeye	Black Widow	Captain America	NaN	Iron Man
age	42	41	28	104	1493	42

```
In [63]: avengers.dtypes # the data types contained in the DataFrame
```

```
Out[63]: actor      object
hero name  object
age         int64
dtype: object
```

```
In [64]: avengers.shape # shape
```

```
Out[64]: (6, 3)
```

Importing Data with pd.read_csv()

```
In [65]: # Titanic Dataset
url = 'https://assets.datacamp.com/production/course_1607/datasets/titanic_sub.csv'
titanic = pd.read_csv(url)
```

```
In [66]: titanic
```

```
Out[66]:
```

	PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	male	35.0	0	0	373450	8.0500	NaN	S
...
886	887	0	2	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	890	1	1	male	26.0	0	0	111369	30.0000	C148	C
890	891	0	3	male	32.0	0	0	370376	7.7500	NaN	Q

891 rows × 11 columns

```
In [67]: titanic.shape
```

```
Out[67]: (891, 11)
```

```
In [68]: titanic.columns
```

```
Out[68]: Index(['PassengerId', 'Survived', 'Pclass', 'Sex', 'Age', 'SibSp', 'Parch',  
               'Ticket', 'Fare', 'Cabin', 'Embarked'],  
              dtype='object')
```

```
In [69]: titanic.index
```

```
Out[69]: RangeIndex(start=0, stop=891, step=1)
```

```
In [70]: titanic.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 11 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   Sex             891 non-null   object  
4   Age             714 non-null   float64 
5   SibSp          891 non-null   int64  
6   Parch          891 non-null   int64  
7   Ticket         891 non-null   object  
8   Fare           891 non-null   float64 
9   Cabin          204 non-null   object  
10  Embarked       889 non-null   object  
dtypes: float64(2), int64(5), object(4)
memory usage: 76.7+ KB
```

```
In [71]: titanic.describe() # displays summary statistics of the numeric variables
```

```
Out[71]:
```

	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	32.204208
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