

Pandas Assignment

```
import pandas and numpy with their aliases

In [6]: import pandas as pd
import numpy as np

Create a variable a = pd.Series[ 100, 200, 300, 400]

In [7]: a = pd.Series([100, 200, 300, 400])
print(a)

0    100
1    200
2    300
3    400
dtype: int64

Print a, and data type

In [8]: print(a)
print("ndata type of Series 'a':", type(a))

0    100
1    200
2    300
3    400
dtype: int64

Data type of Series 'a': <class 'pandas.core.series.Series'>

Using indexing access the element 300 from the series a

In [9]: element_300 = a[2]
print("Element 300 from the Series 'a':", element_300)

Element 300 from the Series 'a': 300

What are the values of index for series a?

In [10]: a = pd.Series([100, 200, 300, 400])
print("Index values for Series 'a':", a.index)

Index values for Series 'a': RangeIndex(start=0, stop=4, step=1)

Change the index to ['a', 'b', 'c', 'd']

In [11]: a.index = ['a', 'b', 'c', 'd']
print("Series 'a' with new index:")
print(a)

Series 'a' with new index:
a    100
b    200
c    300
d    400
dtype: int64

Access the value in the series with index 'd'

In [12]: value_d = a['d']
print("Value in the Series with index 'd':", value_d)

Value in the Series with index 'd': 400

Sort the values wrt to the index and print it

In [13]: sorted_a = a.sort_index()
print("Sorted Series 'a' with respect to the index:")
print(sorted_a)

Sorted Series 'a' with respect to the index:
a    200
b    300
c    100
d    400
dtype: int64

Create a new Pandas Series b having index as 'v', 'r', and 'g' and value 800,450,100 and print it

In [14]: b = pd.Series([800, 450, 100], index=['v', 'r', 'g'])

# Printing the Series 'b'
print(b)

v    800
r    450
g    100
dtype: int64

Append a series at the end of a series

In [15]: a_appended = pd.concat([a, b])

# Printing the appended Series
print("Appended Series:")
print(a_appended)

Appended Series:
c    100
a    200
b    300
d    400
e    800
f    450
g    100
dtype: int64

In [16]: #Print a again after appending b into it
print(a_appended)

c    100
a    200
b    300
d    400
e    800
f    450
g    100
dtype: int64

Sort the values in descending order of a and print the index of the sorted series

In [17]: sorted_a_appended = a_appended.sort_values(ascending=False)

# Print the index of the sorted series
print(sorted_a_appended.index)

Index(['e', 'f', 'd', 'b', 'a', 'c', 'g'], dtype='object')
```

Pandas DataFrame

Part 1

Create a pandas dataframe df from the series 'a' that we used in the last section, print the dataframe

```
In [18]: df = pd.DataFrame(a_appended, columns=['values'])
print("DataFrame 'df' from Series 'a':")
print(df)

DataFrame 'df' from Series 'a':
   values
c    100
a    200
b    300
d    400
e    800
f    450
g    100

What is the shape of the dataframe
(also, what does it imply?)

In [19]: print("Shape of the DataFrame 'df':", df.shape)

Shape of the DataFrame 'df': (7, 1)

Hey! remember shape (7,1) implies dataframe has 7 rows and 1 column.

What is the index of the dataframe, is it same as the series 'a'

In [20]: # yes its same as the series.
print("Index of the DataFrame 'df':", df.index)

Index of the DataFrame 'df': Index(['c', 'a', 'b', 'd', 'e', 'f', 'g'], dtype='object')

print the head and tail of the dataframe.
Additional- (what does head and tail represent?)

In [21]: print("Head of the DataFrame 'df':")
print(df.head())

Head of the DataFrame 'df':
   values
c    100
a    200
b    300
d    400
e    800

In [22]: print("\nTail of the DataFrame 'df':")
print(df.tail())

Tail of the DataFrame 'df':
   values
a    200
d    400
e    800
f    450
g    100

Rename the column of the dataframe as 'points'

In [23]: df = df.rename(columns={'values': 'points'})
print("DataFrame 'df' with column 'points':")
print(df)

DataFrame 'df' with column 'points':
   points
c    100
a    200
b    300
d    400
e    800
f    450
g    100

Create another Series 'fruits', which contains random names of fruits from ['orange','mango','apple']. The series should contain 7 elements, randomly selected from ['orange','mango','apple']

In [24]: #Create fruits array
import random
fruit_list = ['orange', 'mango', 'apple']
fruits = pd.Series(random.choice(fruit_list, n=7))
print("Series 'fruits' with random fruit names:")
print(fruits)

Series 'fruits' with random fruit names:
0    apple
1    apple
2    apple
3    apple
4    orange
5    apple
6    orange
dtype: object

In [25]: #Create series fruits out of fruits array

fruits_array = ['orange', 'mango', 'apple', 'orange', 'mango', 'apple', 'orange']
fruits = pd.Series(fruits_array)
print("Series 'fruits' from the fruits array:")
print(fruits)

Series 'fruits' from the fruits array:
0    orange
1    mango
2    apple
3    orange
4    mango
5    apple
6    orange
dtype: object

Change the index of fruits to the index of dataframe df

In [26]: fruits.index = df.index
print("Series 'fruits' with the index matching DataFrame 'df':")
print(fruits)

Series 'fruits' with the index matching DataFrame 'df':
c    orange
a    mango
b    apple
d    orange
e    mango
f    apple
g    orange
dtype: object

Add this fruits series as a new column to the dataframe df with its column name as 'fruits'
print the head of the dataframe to verify

In [27]: df['fruits'] = fruits

print("Head of the DataFrame 'df' with the new 'fruits' column:")
print(df.head())

Head of the DataFrame 'df' with the new 'fruits' column:
   points  fruits
c    100  orange
a    200  mango
b    300  apple
d    400  orange
e    800  mango
```

Pandas Concatenation

Create a dataframe d1 where the cols are 'city' ['Chandigarh', 'Delhi', 'Kanpur', 'Chennai', 'Mumbai'] and 'Temperature' [15, 22, 20, 26, -2]

```
In [28]: data = {
    'city': ['Chandigarh', 'Delhi', 'Kanpur', 'Chennai', 'Mumbai'],
    'Temperature': [15, 22, 20, 26, -2]
}

d1 = pd.DataFrame(data)

Print(d1)

In [29]: print("DataFrame 'd1' with city and Temperature:")
print(d1)

DataFrame 'd1' with city and Temperature:
   city  Temperature
0  Chandigarh      15
1    Delhi        22
2   Kanpur        20
3   Chennai        26
4    Mumbai        -2

What is the shape of d1.

In [30]: d1.shape

(5, 2)

Out[30]: (5, 2)

Set city = d1['city']

In [32]: city = d1['city']

print city

What is the type of city.

In [33]: # Print the variable 'city'
print("Variable 'city':")
print(city)

Variable 'city':
0  Chandigarh
1    Delhi
2   Kanpur
3   Chennai
4    Mumbai
Name: city, dtype: object

Create another dataframe 'd2' where the columns are
'city' ['Bengaluru','Coimbatore','Srirangan','Pondicherry']
'Temperature' [24,25,36,39]

In [34]: # Creating the data for DataFrame 'd2'
data_d2 = {
    'city': ['Bengaluru', 'Coimbatore', 'Srirangan', 'Pondicherry'],
    'Temperature': [24, 35, 36, 39]
}

# Creating the DataFrame 'd2'
d2 = pd.DataFrame(data_d2)

# Printing the DataFrame 'd2'
print("DataFrame 'd2' with city and Temperature:")
print(d2)

DataFrame 'd2' with city and Temperature:
   city  Temperature
0  Bengaluru      24
1  Coimbatore      36
2  Srirangan      39
3  Pondicherry     39

print the shape of this dataframe

In [35]: d2.shape

(4, 2)

Out[35]: (4, 2)

merge the two dataframes together, save it in a new dataframe named 'd3'

In [36]: # d3 = pd.concat([d1, d2], ignore_index=True)
print("Merged DataFrame 'd3':")
print(d3)

Merged DataFrame 'd3':
   city  Temperature
0  Chandigarh      15
1    Delhi        22
2   Kanpur        20
3   Chennai        26
4    Mumbai        -2
5  Bengaluru      24
6  Coimbatore      36
7  Srirangan      39
8  Pondicherry     39

Select the part of the dataframe such that it contains cities where temp is less than or equal to 20
How many cities are there?

In [37]: selected_cities = d3[d3['Temperature'] <= 20]
print("Cities where temperature is less than or equal to 20:")
print(selected_cities)

num_cities = len(selected_cities)
print("Number of cities where temperature is less than or equal to 20:", num_cities)

Cities where temperature is less than or equal to 20:
   city  Temperature
0  Chandigarh      15
2   Kanpur        20
4    Mumbai        -2
Number of cities where temperature is less than or equal to 20: 3

Select the part of the dataframe such that it contains the cities where temperature greater than or equal to 35

In [38]: selected_cities_high_temp = d3[d3['Temperature'] >= 35]
print("Cities where temperature is greater than or equal to 35:")
print(selected_cities_high_temp)

num_cities_high_temp = len(selected_cities_high_temp)
print("Number of cities where temperature is greater than or equal to 35:", num_cities_high_temp)

Cities where temperature is greater than or equal to 35:
   city  Temperature
6  Coimbatore      36
7  Srirangan      39
8  Pondicherry     39
Number of cities where temperature is greater than or equal to 35: 3
```

Applying functions to columns and creating new columns

We need to create another column in d3, which contains a boolean value for each city to indicate whether it's a union territory or not.

- HINT: Chandigarh, Pondicherry and Delhi are only 3 union territories here.

```
In [54]: def is_union_territory(city):
    union_territories = ['Chandigarh', 'Pondicherry', 'Delhi']
    if city in union_territories:
        return True
    else:
        return False
d3['is_ut'] = d3['city'].apply(is_ut)

In [55]: # Print d3
print("DataFrame 'd3' with is_ut column:")
print(d3)

DataFrame 'd3' with is_ut column:
   city  Temperature  is_ut
0  Chandigarh      15   True
1    Delhi        22   True
2   Kanpur        20   False
3   Chennai        26   False
4    Mumbai        -2   False
5  Bengaluru      24   False
6  Coimbatore      36   False
7  Srirangan      39   False
8  Pondicherry     39   True

The temperatures mentioned in 'Temperature' column are mentioned in Celsius, we need another column which contains the same in Fahrenheit.

HINT -
• Define a function c_to_f which takes input temp in celsius and returns a value with temperature in Fahrenheit.
• To check c_to_f(10) should return 50.

In [56]: # Write function here

def c_to_f(celsius_temp):
    fahrenheit_temp = (celsius_temp * 9/5) + 32
    return fahrenheit_temp
d3['fahrenheit'] = d3['temperature'].apply(c_to_f)
print(d3)

DataFrame 'd3' with Temperature_F column (in Fahrenheit):
   city  Temperature  fahrenheit
0  Chandigarh      15          59.0
1    Delhi        22          71.6
2   Kanpur        20          68.0
3   Chennai        26          78.8
4    Mumbai        -2          28.4
5  Bengaluru      24          75.2
6  Coimbatore      36          96.8
7  Srirangan      39          102.2
8  Pondicherry     39          102.2

In [57]: # Check function c_to_f(30)
def c_to_f(celsius_temp):
    fahrenheit_temp = (celsius_temp * 9/5) + 32
    return fahrenheit_temp

print("Temperature in Fahrenheit for 10 Celsius:", c_to_f(10))

Temperature in Fahrenheit for 10 Celsius: 50.0

In [58]: # Apply function c_to_f to d3 to create a column 'temp_fahrenheit'
d3['temp_fahrenheit'] = d3['Temperature'].apply(c_to_f)
print("DataFrame 'd3' with temp_fahrenheit column (in Fahrenheit):")
print(d3)

DataFrame 'd3' with temp_fahrenheit column (in Fahrenheit):
   city  Temperature  fahrenheit  temp_fahrenheit
0  Chandigarh      15          59.0          59.0
1    Delhi        22          71.6          71.6
2   Kanpur        20          68.0          68.0
3   Chennai        26          78.8          78.8
4    Mumbai        -2          28.4          28.4
5  Bengaluru      24          75.2          75.2
6  Coimbatore      36          96.8          96.8
7  Srirangan      39          102.2          102.2
8  Pondicherry     39          102.2          102.2
```

Indexing and selecting rows in DataFrame

Select subset of the dataframe d1 such that it contains the cities which are union territories.

```
In [61]: # List of union territories
union_territories = ['Chandigarh', 'Pondicherry', 'Delhi']
subset_d1_ut = d1[d1['city'].isin(union_territories)]
print(subset_d1_ut)

Subset of dataframe d1 with cities that are union territories:
   city  Temperature
0  Chandigarh      15
1    Delhi        22

Select a subset of the dataframe d1 such that it contains the cities which only have temperature above 90 Fahrenheit.

In [62]: subset_d1_above_90F = d1[d1['Temperature'] >= 90]
print("Subset of DataFrame d1 with cities having temperatures above 90 Fahrenheit:")
print(subset_d1_above_90F)

Subset of DataFrame d1 with cities having temperatures above 90 Fahrenheit:
Empty DataFrame
Columns: [city, Temperature]
Index: []

Select only the first three rows of the dataframe d1.

In [63]: subset_d1_first_three_rows = d1.iloc[:3]
print("First three rows of DataFrame d1:")
print(subset_d1_first_three_rows)

First three rows of DataFrame d1:
   city  Temperature
0  Chandigarh      15
1    Delhi        22
2   Kanpur        20

Select all the rows and last two columns in the dataframe.

In [64]: subset_d1_last_two_columns = d1.iloc[:, -2:]
print("All rows and last two columns of DataFrame d1:")
print(subset_d1_last_two_columns)

All rows and last two columns of DataFrame d1:
   city  Temperature
0  Chandigarh      15
1    Delhi        22
2   Kanpur        20
3   Chennai        26
4    Mumbai        -2
```

Groupby

```
In [65]: # Create a dataframe using dictionary of your choice
data = {
    'Name': ['Alice', 'Bob', 'Charlie', 'David'],
    'Age': [25, 30, 35, 40],
    'City': ['New York', 'Los Angeles', 'Chicago', 'Houston']
}

# Create a DataFrame using the dictionary
df = pd.DataFrame(data)

# Print the DataFrame
print("DataFrame created from dictionary:")
print(df)

DataFrame created from dictionary:
   Name  Age  City
0  Alice  25  New York
1   Bob   30  Los Angeles
2  Charlie 35  Chicago
3  David  40  Houston

In [66]: # Use Groupby of single column with aggregate sum()
# Sample DataFrame
data = {
    'Category': ['A', 'B', 'A', 'B', 'A'],
    'Value': [10, 20, 30, 40, 50]
}

df = pd.DataFrame(data)
result = df.groupby('Category').aggregate({'Value': 'sum'})
print(result)

   Category  Value
A         90
B         60

In [67]: # Use Groupby of single column with aggregate count()
result = df.groupby('Category').aggregate({'Value': 'count'})
print(result)

   Category  Value
A         3
B         2

In [68]: # Use Groupby of single column with aggregate min() and max()
result = df.groupby('Category').aggregate({'min': 'min', 'max': 'max'})
print(result)

   Category  min  max
A         X    10  50
B         Y    20  40

In [69]: # Use Groupby of any 2 columns with aggregate mean()
data = {
    'Category': ['A', 'B', 'A', 'B', 'A'],
    'Subcategory': ['X', 'Y', 'X', 'Y', 'X'],
    'Value': [10, 20, 30, 40, 50]
}

df = pd.DataFrame(data)
result = df.groupby(['Category', 'Subcategory']).aggregate({'Value': 'mean'})
# Print the result
print(result)

   Category Subcategory  Value
A         X          30.0
B         Y          30.0

In [70]: # Use Groupby of any 2 columns with aggregate min() and max()
data = {
    'Category': ['A', 'B', 'A', 'B', 'A'],
    'Subcategory': ['X', 'Y', 'X', 'Y', 'X'],
    'Value': [10, 20, 30, 40, 50]
}

df = pd.DataFrame(data)
result = df.groupby(['Category', 'Subcategory']).aggregate({'Value': ['min', 'max']})
print(result)

   Category Subcategory  min  max
A         X          10  50
B         Y          20  40
```

Data Range

Create a pandas datarange where starting date is 1st of January 2020 and end date is 1st of April 2021, store it in a new variable named 'a'

```
In [71]: a = pd.date_range(start='2020-01-01', end='2021-04-01')

print(a)

DatetimeIndex(['2020-01-01', '2020-01-02', '2020-01-03', '2020-01-04',
               '2020-01-05', '2020-01-06', '2020-01-07', '2020-01-08',
               '2020-01-09', '2020-01-10',
               ...,
               '2021-03-23', '2021-03-24', '2021-03-25', '2021-03-26',
               '2021-03-27', '2021-03-28', '2021-03-29', '2021-03-30',
               '2021-03-31', '2021-04-01'],
              dtype='datetime64[ns]', length=467, freq='D')

What is the len of a?

In [73]: length_of_a = len(a)
print("Length of 'a':", length_of_a)

Length of 'a': 467

What is the type of a?

In [74]: type_of_a = type(a)
print("Type of 'a':", type_of_a)

Type of 'a': <class 'pandas.core.indexes.datetimes.DatetimeIndex'>
```