# Pandas
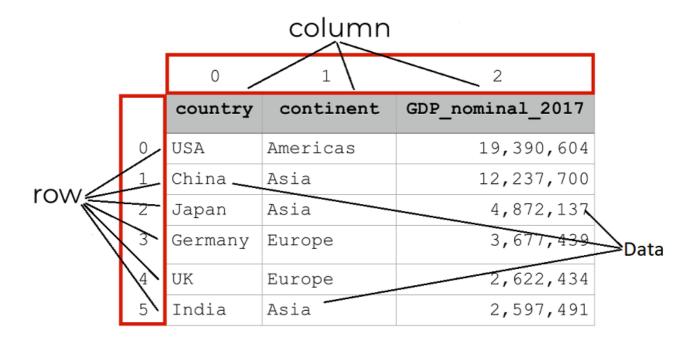
# DataFrame

**Pandas DataFrame** is a two-dimensional data structure with labeled axes (rows and columns).

It consists of three principal components:
the **data**, **rows** and **columns**.

Topics:

– Create a Pandas DataFrame

– Rows and Columns Handling

– Subset the DataFrame

# Create a Pandas DataFrame

In the real world, a Pandas DataFrame will be created by loading the datasets from existing storage. The storage can be CSV file, Excel file and SQL Database etc.

Pandas DataFrame can be created from the lists, dictionary, and from a list of dictionary etc.

Dataframe can be created in different ways and here, we have discussed some ways
by which we create a dataframe:

# Creating a dataframe from List

DataFrame can be created using a single list or a list of lists.

```
1   # import pandas as pd
2   import pandas as pd
3
4   # list of strings
5   lst = ['Code', 'Data', 'AI', 'is',
6               'welcoming', 'all', 'Data Enthusiasts']
7
8   # Calling DataFrame constructor on list
9   df = pd.DataFrame(lst)
10  print(df)
```

Output

```
1                    0
2   0            Code
3   1            Data
4   2              AI
5   3              is
6   4       welcoming
7   5             all
8   6  Data Enthusiasts
```

# Creating DataFrame from dictionary of lists

Suppose you want to create a dataframe out of the below dictionary **dict**:

dict = {'name':["Sunil", "pankaj", "sudhir", "Geeku"],

'degree': ["MBA", "BCA", "M.Tech", "MBA"],

'score':[90, 40, 80, 98]}

Here, the dictionary **keys** 'name' ,'degree' & 'score' will become Columns and the
**values (lists)** will become data.

Therefore all the values ["Sunil", "pankaj", "sudhir", "Geeku"], ["MBA", "BCA", "M.Tech", "MBA"] and [90, 40, 80, 98] need to be of **same length**.(here, it is 4).

Therefore, the basic condition to convert a dictionary of lists to a dataframe is that all the lists should be of same length.

```
# importing pandas as pd
import pandas as pd

# dictionary of lists
dict = {'name':["Sunil", "pankaj", "sudhir", "Geeku"],
        'degree': ["MBA", "BCA", "M.Tech", "MBA"],
        'score':[90, 40, 80, 98]}

df = pd.DataFrame(dict)

print(df)
```

Output

```
     name  degree  score
0   Sunil     MBA     90
1  pankaj     BCA     40
2  sudhir  M.Tech     80
3   Geeku     MBA     98
```

# Rows and Columns Handling

A Data frame is a two-dimensional data structure where data are stored in rows and columns. Each row is called observation and each coulmn is termed as feature.

Here, you will perform basic operations on rows/columns like **selecting**, **adding**, **deleting** and **renaming**.

## Select Rows

Pandas provide a unique method to retrieve rows from a Data frame. `DataFrame.loc[]` method is used to retrieve rows from Pandas DataFrame.

Rows can also be selected by passing integer location to an **iloc[]** function.

**loc** gets rows (**or** columns) with particular labels from the index.

**iloc** gets rows (**or** columns) at particular positions in the index (it only takes integers)

```
1   # Import pandas package
2   import pandas as pd
3
4   # Define a dictionary containing Students data
5   data = {'Height': [5.1, 6.2, 5.1, 5.2],
6           'Qualification': ['Msc', 'MA', 'Msc', 'Msc'],
7           'Hobby': ['Poetry', 'Travelling', 'Biking', 'Sports']}
8
9   df = pd.DataFrame(data)
10  df.index = ['Sarah', 'Princi', 'Gaurav', 'Anuj']
11
12  print(df)
13
14  print ('\n Select Height, Qualification & Hobby for Sarah \n')
15  df.loc['Sarah']
```

Output

```
1           Height Qualification        Hobby
2   Sarah      5.1           Msc       Poetry
3   Princi     6.2            MA   Travelling
4   Gaurav     5.1           Msc       Biking
5   Anuj       5.2           Msc       Sports
6
7    Select Height, Qualification & Hobby for Sarah
8
9   Height               5.1
10  Qualification        Msc
11  Hobby             Poetry
12  Name: Sarah, dtype: object
```

Selecting rows by **index names**

| | Height | Qualification | Hobby |
|---|---|---|---|
| **Sarah** | 5.1 | Msc | Poetry |
| **Princi** | 6.2 | MA | Travelling |
| **Gaurav** | 5.1 | Msc | Biking |
| **Anuj** | 5.2 | Msc | Sports |

```
1  # select two rows
2  df.loc[['Sarah','Gaurav']]
```

## Output

```
1        Height  Qualification   Hobby
2  Sarah   5.1 Msc         Poetry
3  Gaurav  5.1 Msc         Biking
```

## Selecting rows by slice of index names

| | Height | Qualification | Hobby |
|---|---|---|---|
| **Sarah** | 5.1 | Msc | Poetry |
| **Princi** | 6.2 | MA | Travelling |
| **Gaurav** | 5.1 | Msc | Biking |
| **Anuj** | 5.2 | Msc | Sports |

```
1  # select 1st 3 rows
2  df.loc['Sarah':'Gaurav']
```

## Output

```
1        Height  Qualification   Hobby
2  Sarah   5.1 Msc Poetry
3  Princi  6.2 MA  Travelling
4  Gaurav  5.1 Msc Biking
```

## Add Rows

To add a Row in Pandas DataFrame, you can concat the old dataframe with new one.

```
# Import pandas package
import pandas as pd

# Define a dictionary containing Students data
data = {'Height': [5.1, 6.2, 5.1, 5.2],
        'Qualification': ['Msc', 'MA', 'Msc', 'Msc'],
        'Hobby': ['Poetry', 'Travelling', 'Biking', 'Sports']}

df = pd.DataFrame(data)
df.index = ['Sarah', 'Princi', 'Gaurav', 'Anuj']

print(df)


new_row = pd.DataFrame({'Height': 5.3, 'Qualification': 'Bachelor' , 'Hobby': 'Sports'},
                                                        index =['NewRow'])
df = pd.concat([new_row, df])

print ('\n After adding row NewRow \n')

df
```

## Output

```
        Height Qualification        Hobby
Sarah     5.1           Msc       Poetry
Princi    6.2            MA    Travelling
Gaurav    5.1           Msc       Biking
Anuj      5.2           Msc       Sports

 After adding row NewRow

        Height   Qualification    Hobby
NewRow  5.3 Bachelor     Sports
Sarah   5.1 Msc          Poetry
Princi  6.2 MA           Travelling
Gaurav  5.1 Msc          Biking
Anuj    5.2 Msc          Sports
```

You can also append all the rows of a dataframe to a new dataframe.

```
# Import pandas package
import pandas as pd

# Define a dictionary containing Students data
data = {'Name': ['Sarah', 'Princi', 'Gaurav'],
        'Height': [5.1, 6.2, 5.1],
        'Qualification': ['Msc', 'MA', 'Msc'],
        'Hobby': ['Poetry', 'Travelling', 'Biking']}

# Convert the dictionary into DataFrame
df = pd.DataFrame(data)
```

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 0 | Sarah | 5.1 | Msc | Poetry |
| 1 | Princi | 6.2 | MA | Travelling |
| 2 | Gaurav | 5.1 | Msc | Biking |

*Original DataFrame*

```
1    # Define a dictionary containing New data to append
2    data = {'Name': ['Janvi', 'Rushel'],
3            'Height': [5.1, 6.2],
4            'Qualification': ['Msc', 'MA'],
5            'Hobby': ['Poetry', 'Travelling']}
6
7    # Convert the dictionary into DataFrame
8    df_to_append = pd.DataFrame(data)
```

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 0 | Janvi | 5.1 | Msc | Poetry |
| 1 | Rushel | 6.2 | MA | Travelling |

**DataFrame to Append**

```
1    #append the df_to_append to the original dataframe
2    df.append(df_to_append, ignore_index= True)
```

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 0 | Sarah | 5.1 | Msc | Poetry |
| 1 | Princi | 6.2 | MA | Travelling |
| 2 | Gaurav | 5.1 | Msc | Biking |
| 3 | Janvi | 5.1 | Msc | Poetry |
| 4 | Rushel | 6.2 | MA | Travelling |

*DataFrame after Append*

# Delete Rows

To delete a row in Pandas DataFrame, we can use the drop() method.

Rows is deleted by dropping Rows by index label.

```
1    # Import pandas package
2    import pandas as pd
3
4    # Define a dictionary containing Students data
```

```
5    data = {'Name': ['Sarah', 'Princi', 'Gaurav'],
6            'Height': [5.1, 6.2, 5.1],
7            'Qualification': ['Msc', 'MA', 'Msc'],
8            'Hobby': ['Poetry', 'Travelling', 'Biking']}
9
10   # Convert the dictionary into DataFrame
11   df = pd.DataFrame(data)
12
13   print(df)
14
15   print('\n After dropping \n')
16
17   df.drop([0,1])
```

## Output

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 0 | Sarah | 5.1 | Msc | Poetry |
| 1 | Princi | 6.2 | MA | Travelling |
| 2 | Gaurav | 5.1 | Msc | Biking |

*Before Dropping Rows*

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 2 | Gaurav | 5.1 | Msc | Biking |

*After dropping 1st 2 rows*

# Rename Index

You can rename index names **rename()** function.

```
1    import pandas as pd
2
3    # making data frame from csv file
4    # Define a dictionary containing Students data
5    data = {'Name': ['Sarah', 'Princi', 'Gaurav'],
6            'Height': [5.1, 6.2, 5.1],
7            'Qualification': ['Msc', 'MA', 'Msc'],
8            'Hobby': ['Poetry', 'Travelling', 'Biking']}
9
10   # Convert the dictionary into DataFrame
11   df = pd.DataFrame(data)
```

```
1    # changing index cols with rename() to 1st, 2nd & 3rd Record respectively
2    df.rename(index = {  0: "1st Record",
3                         1:"2nd Record",
4                         2: "3rd Record"},
5                               inplace = True)
6    # display
7    df
```

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 0 | Sarah | 5.1 | Msc | Poetry |
| 1 | Princi | 6.2 | MA | Travelling |
| 2 | Gaurav | 5.1 | Msc | Biking |

*Before Renaming Index*

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 1st Record | Sarah | 5.1 | Msc | Poetry |
| 2nd Record | Princi | 6.2 | MA | Travelling |
| 3rd Record | Gaurav | 5.1 | Msc | Biking |

*After Renaming Index*

# Select Column

To select a column in Pandas DataFrame, you can either access the columns by calling them by their columns name or column number.

Select columns by **column name**

```
# Import pandas package
import pandas as pd

# Define a dictionary containing employee data
data = {'Name':['Sarah', 'Princi', 'Gaurav', 'Anuj'],
        'Age':[27, 24, 22, 32],
        'Address':['Kolkata', 'Kanpur', 'Allahabad', 'Delhi'],
        'Qualification':['Msc', 'MA', 'MCA', 'Phd']}

# Convert the dictionary into DataFrame
df = pd.DataFrame(data)

# select two columns
df[['Name', 'Qualification']]
```

## Output

```
1          Name    Qualification
2     0    Sarah   Msc
3     1    Princi  MA
4     2    Gaurav  MCA
5     3    Anuj    Phd
```

## Select columns by **column number**

```
1   # Import pandas package
2   import pandas as pd
3
4   # Define a dictionary containing employee data
5   data = {'Name':['Sarah', 'Princi', 'Gaurav', 'Anuj'],
6           'Age':[27, 24, 22, 32],
7           'Address':['Kolkata', 'Kanpur', 'Allahabad', 'Delhi'],
8           'Qualification':['Msc', 'MA', 'MCA', 'Phd']}
9
10  # Convert the dictionary into DataFrame
11  df = pd.DataFrame(data)
12
13
14  # select all rows by ':'
15  # select two columns Name and Qualification by their position
16
17  df.iloc[:, [0,3]]
```

## Output

```
1          Name    Qualification
2     0    Sarah   Msc
3     1    Princi  MA
4     2    Gaurav  MCA
5     3    Anuj    Phd
```

## Select **1st 3 columns**

```
1   df[df.columns[0:3]]
```

## Output

```
1          Name    Age  Address
2     0    Sarah   27   Kolkata
3     1    Princi  24   Kanpur
4     2    Gaurav  22   Allahabad
5     3    Anuj    32   Delhi
```

Select columns from "**Name**" to "**Address**"

```
1   # select two rows and
2   # column "name" to "Address"
3   # Means total three columns
4   df.loc[ :,'Name':'Address']
```

## Output

```
1        Name    Age Address
2   0    Sarah   27  Kolkata
3   1    Princi  24  Kanpur
4   2    Gaurav  22  Allahabad
5   3    Anuj    32  Delhi
```

# Add Column

To add a column in Pandas DataFrame, you can declare a new list as a column and add to a existing Dataframe.

```
1    # Import pandas package
2    import pandas as pd
3
4    # Define a dictionary containing Students data
5    data = {'Name': ['Sarah', 'Princi', 'Gaurav', 'Anuj'],
6            'Height': [5.1, 6.2, 5.1, 5.2],
7            'Qualification': ['Msc', 'MA', 'Msc', 'Msc']}
8
9    # Convert the dictionary into DataFrame
10   df = pd.DataFrame(data)
11   df
```

|   | Name | Height | Qualification |
|---|------|--------|---------------|
| 0 | Sarah | 5.1 | Msc |
| 1 | Princi | 6.2 | MA |
| 2 | Gaurav | 5.1 | Msc |
| 3 | Anuj | 5.2 | Msc |

*Before Adding Column*

```
1    # Declare a list that is to be converted into a column
2    address = ['Kolkata', 'Bangalore', 'Chennai', 'Mumbai']
3
4    # Using 'Address' as the column name
5    # and equating it to the list
6    df['Address'] = address
7    df
```

| | Name | Height | Qualification | Address |
|---|---|---|---|---|
| 0 | Sarah | 5.1 | Msc | Kolkata |
| 1 | Princi | 6.2 | MA | Bangalore |
| 2 | Gaurav | 5.1 | Msc | Chennai |
| 3 | Anuj | 5.2 | Msc | Mumbai |

*After Adding Column Address*

# Delete Column

To delete a column in Pandas DataFrame, you can use the `drop()` method. Columns are deleted by **dropping** columns with columnnames.

```
# Import pandas package
import pandas as pd

# Define a dictionary containing Students data
data = {'Name': ['Sarah', 'Princi', 'Gaurav', 'Anuj'],
        'Height': [5.1, 6.2, 5.1, 5.2],
        'Qualification': ['Msc', 'MA', 'Msc', 'Msc'],
        'Hobby': ['Poetry', 'Travelling', 'Biking', 'Sports']}

# Convert the dictionary into DataFrame
df = pd.DataFrame(data)
df
```

| | Name | Height | Qualification | Hobby |
|---|---|---|---|---|
| 0 | Sarah | 5.1 | Msc | Poetry |
| 1 | Princi | 6.2 | MA | Travelling |
| 2 | Gaurav | 5.1 | Msc | Biking |
| 3 | Anuj | 5.2 | Msc | Sports |

*Before Deleting Columns*

```
df.drop(["Height", "Hobby"], axis = 1, inplace = True)
```

|   | Name | Qualification |
|---|------|---------------|
| 0 | Sarah | Msc |
| 1 | Princi | MA |
| 2 | Gaurav | Msc |
| 3 | Anuj | Msc |

*After Deleting Columns*

# Rename Column

To Rename columns, you can use **df.columns = new list of column names**.

```
1   # Import pandas package
2   import pandas as pd
3
4   # Define a dictionary containing Students data
5   data = {'Name': ['Sarah', 'Princi', 'Gaurav', 'Anuj'],
6           'Height': [5.1, 6.2, 5.1, 5.2],
7           'Qualification': ['Msc', 'MA', 'Msc', 'Msc'],
8           'Hobby': ['Poetry', 'Travelling', 'Biking', 'Sports']}
9
10  # Convert the dictionary into DataFrame
11  df = pd.DataFrame(data)
12  df
13  # change the column name Qualification & Hobby to Degree and Leisure respectively
14  df.columns = ['Name','Height', 'Degree', 'Leisure']
15  df
```

## Output

```
1         Name   Height  Degree  Leisure
2   0    Sarah    5.1 Msc Poetry
3   1    Princi   6.2 MA  Travelling
4   2    Gaurav   5.1 Msc Biking
5   3    Anuj     5.2 Msc Sports
```

You can choose to **rename a particular column** as well.

```
1   # renaming 'Degree' back to 'Qualification'
2   df.rename(columns={'Degree':'Qualification'}, inplace=True)
3   df
```

```
1      Name   Height Qualification   Leisure
2   0  Sarah   5.1 Msc          Poetry
```

```
3 | 1   Princi  6.2 MA        Travelling
4 | 2   Gaurav  5.1 Msc       Biking
5 | 3   Anuj    5.2 Msc       Sports
```

# Subset the DataFrame

Subsetting a Dataframe is same as slicing a dataframe into smaller
dataframe so that you can focus on a small chunk of a large dataset at a
particular time.

It is one of the main skills for *Exploratory Data Analysis* (EDA).

For example, you need to slice the box out of the dataframe:



*Subsetting the Squared Portion*

To do this, you will subset the two rows (**2nd** and **3rd**) and
1st two coulmns (**Height** & **Qualification**).

```
1 | df.iloc[1:3,0:2]
```