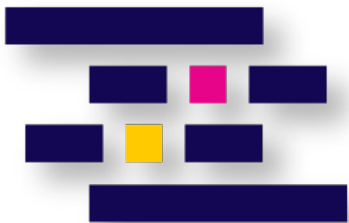




**BRAINALYST**  
A Data Driven Company

pandas



**BRAINALYST'S**

# Learning the **PANDAS** *Library*

## ABOUT BRAINALYST

**Brainalyst** is a pioneering data-driven company dedicated to transforming data into actionable insights and innovative solutions. Founded on the principles of leveraging cutting-edge technology and advanced analytics, Brainalyst has become a beacon of excellence in the realms of data science, artificial intelligence, and machine learning.

### OUR MISSION

At Brainalyst, our mission is to empower businesses and individuals by providing comprehensive data solutions that drive informed decision-making and foster innovation. We strive to bridge the gap between complex data and meaningful insights, enabling our clients to navigate the digital landscape with confidence and clarity.

### WHAT WE OFFER

#### 1. Data Analytics and Consulting

Brainalyst offers a suite of data analytics services designed to help organizations harness the power of their data. Our consulting services include:

- **Data Strategy Development:** Crafting customized data strategies aligned with your business objectives.
- **Advanced Analytics Solutions:** Implementing predictive analytics, data mining, and statistical analysis to uncover valuable insights.
- **Business Intelligence:** Developing intuitive dashboards and reports to visualize key metrics and performance indicators.

#### 2. Artificial Intelligence and Machine Learning

We specialize in deploying AI and ML solutions that enhance operational efficiency and drive innovation. Our offerings include:

- **Machine Learning Models:** Building and deploying ML models for classification, regression, clustering, and more.
- **Natural Language Processing:** Implementing NLP techniques for text analysis, sentiment analysis, and conversational AI.
- **Computer Vision:** Developing computer vision applications for image recognition, object detection, and video analysis.

#### 3. Training and Development

Brainalyst is committed to fostering a culture of continuous learning and professional growth. We provide:

- **Workshops and Seminars:** Hands-on training sessions on the latest trends and technologies in data science and AI.
- **Online Courses:** Comprehensive courses covering fundamental to advanced topics in data analytics, machine learning, and AI.
- **Customized Training Programs:** Tailored training solutions to meet the specific needs of organizations and individuals.



#### 4. Generative AI Solutions

As a leader in the field of Generative AI, Brainalyst offers innovative solutions that create new content and enhance creativity. Our services include:

- **Content Generation:** Developing AI models for generating text, images, and audio.
- **Creative AI Tools:** Building applications that support creative processes in writing, design, and media production.
- **Generative Design:** Implementing AI-driven design tools for product development and optimization.

#### OUR JOURNEY

Brainalyst's journey began with a vision to revolutionize how data is utilized and understood. Founded by Nitin Sharma, a visionary in the field of data science, Brainalyst has grown from a small startup into a renowned company recognized for its expertise and innovation.

#### KEY MILESTONES:

- **Inception:** Brainalyst was founded with a mission to democratize access to advanced data analytics and AI technologies.
- **Expansion:** Our team expanded to include experts in various domains of data science, leading to the development of a diverse portfolio of services.
- **Innovation:** Brainalyst pioneered the integration of Generative AI into practical applications, setting new standards in the industry.
- **Recognition:** We have been acknowledged for our contributions to the field, earning accolades and partnerships with leading organizations.

Throughout our journey, we have remained committed to excellence, integrity, and customer satisfaction. Our growth is a testament to the trust and support of our clients and the relentless dedication of our team.

#### WHY CHOOSE BRAINALYST?

Choosing Brainalyst means partnering with a company that is at the forefront of data-driven innovation. Our strengths lie in:

- **Expertise:** A team of seasoned professionals with deep knowledge and experience in data science and AI.
- **Innovation:** A commitment to exploring and implementing the latest advancements in technology.
- **Customer Focus:** A dedication to understanding and meeting the unique needs of each client.
- **Results:** Proven success in delivering impactful solutions that drive measurable outcomes.

**JOIN US ON THIS JOURNEY TO HARNESS THE POWER OF DATA AND AI. WITH BRAINALYST, THE FUTURE IS DATA-DRIVEN AND LIMITLESS.**



## Basic Python Data Structures:

In basic Python, when working with panel data evaluation, not unusual statistics systems include tuples, lists, dictionaries, and sets.

**One-Dimensional:** These simple facts structures are one-dimensional, which means they constitute a **single collection of elements**.

**Heterogeneous:** They can maintain different varieties of facts consisting of integers, floats, strings, and Booleans inside the identical shape.

**No Broadcasting or Vectorization:** These simple facts systems do no longer inherently support broadcasting or vectorization, making it much less green for numerical computations on big datasets.

## NumPy-ndarray:

NumPy introduces the ndarray, that are a extential statistics shape for numerical computations in Python.

**N-Dimensional:** NumPy arrays are n-dimensional, making an allowance for green storage and manipulation of multi-dimensional statistics.

**Homogeneous:** NumPy arrays are homogeneous, meaning they kept factors of the identical facts kind, leading to green reminiscence usage and operations!!!

**Broadcasting and Vectorization:** NumPy arrays support broadcasting and vectorization, allowing efficient element-smart operations on entire arrays.

```
import pandas as pd
```

## Pandas Series:

Pandas introduces the Series data shape, which is much like NumPy arrays however with extra functions like labelled indices!

**ID Homogeneous:** Series are one-dimensional and homogeneous, making an allowance for efficient garage and operations on categorized facts!

**DataFrame:** Pandas also introduces the DataFrame, a -dimensional tabular data structure.

**Heterogeneous:** DataFrames are heterogeneous, that means they are able to maintain different kinds of facts in exclusive columns!!!

**Broadcasting and Vectorization:** Like NumPy arrays, DataFrames assist broadcasting and vectorization, making them green for data manipulation and evaluation!

When operating with Pandas, it is common to import it the usage of the alias pd, following the convention used within the documentation.



**Creating a Series from a Dictionary:** A Series can be created via passing in a dictionary, in which the keys come to be the index labels and the values emerge as the fact's factors!

**Accessing Elements:** Elements of a Series can be accessed through index or using cutting, much like lists in Python. Additionally, a Series may be accessed like a dictionary, wherein the index values act as keys!

**Handling Missing Keys:** Accessing a non-present key will enhance an exception until you use the '.get()' technique, which returns None if the important thing does not exist.

**Operating on Series:** Series guides mathematical operations like NumPy arrays, taking into consideration detail-smart operations!

**Defining Series with Lists or Arrays:** Series can be described the usage of lists or arrays, in which Pandas robotically assigns index labels.

**Defining Custom Index Labels:** Custom index labels may be described through passing a list to the index parameter when developing a Series.

These examples illustrate the flexibility and capability of Pandas Series, making them a flexible device for records manipulation and evaluation in Python!!

The dir(pd) characteristic in Python returns a looked after listing of all attributes and methods available within the pd module, that is the alias for the Pandas library. Here's a generalized assessment of what you would possibly assume to look!

```
print(dir(pd))
```

['ArrowDtype', 'BooleanDtype', 'Categorical', 'CategoricalDtype', 'CategoricalIndex', 'DataFrame', 'DateOffset', 'DatetimeIndex', 'DatetimeTZDtype', 'ExcelFile', 'ExcelWriter', 'Flags', 'Float32Dtype', 'Float64Dtype', 'Grouper', 'HDFStore', 'Index', 'IndexSlice', 'Int16Dtype', 'Int32Dtype', 'Int64Dtype', 'Int8Dtype', 'Interval', 'IntervalDtype', 'IntervalIndex', 'MultiIndex', 'NA', 'NaT', 'NamedAgg', 'Period', 'PeriodDtype', 'PeriodIndex', 'RangeIndex', 'Series', 'SparseDtype', 'StringDtype', 'Timedelta', 'TimedeltaIndex', 'Timestamp', 'UInt16Dtype', 'UInt32Dtype', 'UInt64Dtype', 'UInt8Dtype', '\_\_all\_\_', '\_\_builtins\_\_', '\_\_cache\_\_', '\_\_doc\_\_', '\_\_docformat\_\_', '\_\_file\_\_', '\_\_git\_version\_\_', '\_\_loader\_\_', '\_\_name\_\_', '\_\_package\_\_', '\_\_path\_\_', '\_\_spec\_\_', '\_\_version\_\_', '\_\_built\_with\_meson', '\_\_config\_\_', '\_\_is\_numpy\_dev\_\_', '\_\_libs\_\_', '\_\_pandas\_datetime\_CAPI', '\_\_pandas\_parser\_CAPI', '\_\_testing\_\_', '\_\_typing\_\_', '\_\_version\_meson\_\_', 'annotations', 'api', 'array', 'arrays', 'bdate\_range', 'concat', 'concat', 'core', 'crosstab', 'cut', 'date\_range', 'describe\_option', 'errors', 'eval', 'factorize', 'from\_dummies', 'get\_dummies', 'get\_option', 'infer\_freq', 'interval\_range', 'io', 'isna', 'isnull', 'json\_normalize', 'lreshape', 'melt', 'merge', 'merge\_asof', 'merge\_ordered', 'notna', 'notnull', 'offsets', 'option\_context', 'options', 'pandas', 'period\_range', 'pivot', 'pivot\_table', 'plotting', 'qcut', 'read\_clipboard', 'read\_csv', 'read\_excel', 'read\_feather', 'read\_fwf', 'read\_gbq', 'read\_hdf', 'read\_html', 'read\_json', 'read\_orc', 'read\_parquet', 'read\_pickle', 'read\_sas', 'read\_spss', 'read\_sql', 'read\_sql\_query', 'read\_sql\_table', 'read\_stata', 'read\_table', 'read\_xml', 'reset\_option', 'set\_eng\_float\_format', 'set\_option', 'show\_versions', 'test', 'testing', 'timedelta\_range', 'to\_datetime', 'to\_numeric', 'to\_pickle', 'to\_timedelta', 'tseries', 'unique', 'util', 'value\_counts', 'wide\_to\_long']

Pandas is primarily designed to work with based facts, which includes tabular facts with categorised rows and columns (like a spreadsheet or database desk!).



In comparison, NumPy is greater perfect for homogeneous numerical facts! Its number one records structure, the ndarray, is homogeneous, which means it could best contain factors of the same statistics kind! This makes NumPy efficient for numerical computations and array operations.

Pandas builds on pinnacle of NumPy and affords extra statistics systems like DataFrame and Series, which are more bendy and appropriate for managing dependent information with specific varieties of values in each column.

This flexibility allows Pandas to address various information types inside a single record shape, making it nicely desirable for data manipulation and analysis tasks normally encountered in information science and analytics!!

- **Default Index:** When you create a Series in Pandas without specifying an index, a default index is mechanically generated by means of the gadget. This default index can't be updated or modified by using the consumer.

```
import pandas as pd

data = [10, 20, 30, 40]
series_default_index = pd.Series(data)
series_default_index

0    10
1    20
2    30
3    40
dtype: int64
```

- **User-Defined Index:** You can also create a Series with a user-described index, in which you specify the index labels yourself. This lets in for personalization of the index labels. If the consumer does not offer any User-Defined Index (UDI), then the UDI might be like the Direct Index (DI)!!!
- **Displayed Index:** When you print or view a Series in Python output, you notice the User-Defined Index (UDI). This is the index that Pandas shows for readability!

```
data = [10, 20, 30, 40]
index_labels = ['A', 'B', 'C', 'D']
series_custom_index = pd.Series(data, index=index_labels)
series_custom_index

A    10
B    20
C    30
D    40
dtype: int64
```

- **Accessing Data:** You can get right of entry to the records in a Series the usage of both the Direct Index (DI) or the User-Defined Index (UDI), relying in your desire or requirement. Both methods are legitimate for retrieving data from a Series.

```
print(series_default_index[0]) # Accessing data using Default Index (DI)
# Output: 10

print(series_custom_index['A']) # Accessing data using User-Defined Index (UDI)
# Output: 10

10
10
```

## Series – Pandas

| Function     | Description   |
|--------------|---|
| Create       | <ul style="list-style-type: none"> <li>- <code>pd.Series()</code> can be used to convert any 1D data structure into a Pandas Series.</li> <li>- Series can also be created from DataFrames.</li> </ul>  |
| Access       | <ul style="list-style-type: none"> <li>- Using square brackets (<code>[]</code>):</li> <li>- Positional access: Series uses UDI for indexing if it's an integer, DI if it's a string.</li> <li>- Slicing: Retrieves data from start till end +-1.</li> <li>- Using <code>.iloc[]</code>:</li> <li>- Always uses DI for accessing data.</li> <li>- Slicing retrieves data from start till end +-1.</li> <li>- Indexing is possible.</li> <li>- Using <code>.loc[]</code>:</li> <li>- Always uses UDI for accessing data.</li> <li>- Slicing retrieves data from start till end.</li> <li>- Indexing is not possible.</li> <li>- Using <code>.ix[]</code>:</li> </ul> |
| Update       | <ul style="list-style-type: none"> <li>- To update values in a Series, use indexing or methods like <code>.at[]</code> or <code>.iat[]</code>.</li> </ul>   |
| Mathematical | <ul style="list-style-type: none"> <li>- Pandas Series support various mathematical operations such as addition, subtraction, etc.</li> </ul>   |
| Indexing     | <ul style="list-style-type: none"> <li>- Accessing data by position or label using square brackets, <code>.iloc[]</code>, or <code>.loc[]</code>.</li> </ul>  |

### Create

- **pd.Series():** You can create a Series in Pandas by means of passing in any individual-dimensional facts shape, like a list or a NumPy array. Additionally, you may create a Series from DataFrames!

```
l1 = [10,20,30,40,50]
```

```
type(l1)
```

```
list
```

```
s1 = pd.Series(l1)
```

```
s2 = pd.Series(l1, index = range(1,6))
```

```
s3 = pd.Series(l1, index = ["a","b","c","d","e"])
```

- **Access[]:** For a Series, the use of square brackets allows you to get right of entry to values by using role. When reducing, if the index is integers, it operates on positional integer based reducing (DI - Direct Indexing). If the index is strings, reducing is primarily based on the label based totally indexing (DI - Direct Indexing)!

.iIoc[]: This attribute always

accesses data the use of Direct Indexing (DI). Slicing operates similarly to square brackets, from the begin until the end plus/minus one!

.loc[]: It constantly accesses statistics the use of Unique Direct Indexing (UDI). Slicing operates from the begin until the quit.

.ix[]: Deprecated and now not advocated to be used!

```
s1
0    10
1    20
2    30
3    40
4    50
dtype: int64
```

```
s1[2]
30
```

```
s1[1:4]
1    20
2    30
3    40
dtype: int64
```

```
s2
1    10
2    20
3    30
4    40
5    50
dtype: int64
```





```
s2[3]
```

```
30
```

```
s3
```

```
a    10
b    20
c    30
d    40
e    50
dtype: int64
```

```
s3['c']
```

```
30
```

```
s3[1:4]
```

```
b    20
c    30
d    40
dtype: int64
```

```
s3['b':'d']
```

```
b    20
c    30
d    40
dtype: int64
```

```
s1
```

```
0    10
1    20
2    30
3    40
4    50
dtype: int64
```

```
s1.iloc[2]
```

```
30
```

```
s1.loc[2]
```

```
30
```

```
s1[1:3]
```

```
1    20
2    30
dtype: int64
```

```
s1.iloc[1:3]
```

```
1    20
2    30
dtype: int64
```

```
s1.loc[1:3]
```

```
1    20
2    30
```

BRAINALYST  
A Data Company

```
s2.iloc[4]
```

```
50
```

```
s2.loc[4]
```

```
40
```

```
s2[4]
```

```
40
```

```
s2[1:3]
```

```
2    20  
3    30  
dtype: int64
```

```
s2.iloc[1:3]
```

```
2    20  
3    30  
dtype: int64
```

```
s2.loc[2:3]
```

```
2    20  
3    30  
dtype: int64
```

```
s3
```

```
a    10  
b    20  
c    30  
d    40  
e    50  
dtype: int64
```

```
s3["c"]
```

```
30
```



```
s3.iloc[2]
```

```
30
```

```
s3.loc['c']
```

```
30
```

```
s3.loc["c":"d"]
```

```
c    30  
d    40  
dtype: int64
```

```
s3.iloc[2:4]
```

```
c    30  
d    40  
dtype: int64
```

```
s3[2:4]
```

```
c    30  
d    40
```

```
s3["c":"d"]
```

```
c    30  
d    40  
dtype: int64
```

**Update:** You can replace values in a Series using indexing or label-based indexing?

## Mathematical Operations:

Pandas Series supports various mathematical operations, consisting of addition, subtraction, multiplication, and department!!! These operations can be achieved on whole Series or between Series.



```
s1/10
```

```
0    1.0
1    2.0
2    3.0
3    4.0
4    5.0
dtype: float64
```

```
s1 + s2
```

```
0    NaN
1    30.0
2    50.0
3    70.0
4    90.0
5    NaN
dtype: float64
```

```
s1[s1 > 30]
```

```
3    40
4    50
dtype: int64
```

```
s1.loc[s1>30]
```

```
3    40
4    50
dtype: int64
```

```
s1.iloc[s1>30]
```

```
-----
NotImplementedError                                Traceback (most recent c
ast)
Cell In[56], line 1
----> 1 s1.iloc[s1>30]

File ~/anaconda3/lib/python3.11/site-packages/pandas/core/indexing
153, in _iLocIndexer._getitem__(self, key)
    1150 axis = self.axis or 0
    1152 maybe_callable = com.apply_if_callable(key, self.obj)
-> 1153 return self._getitem_axis(maybe_callable, axis=axis)

File ~/anaconda3/lib/python3.11/site-packages/pandas/core/indexing
700, in _iLocIndexer._getitem_axis(self, key, axis)
    1697 key = np.asarray(key)
    1699 if com.is_bool_indexer(key):
-> 1700     self._validate_key(key, axis)
```

## Indexing:

Indexing in Pandas Series permits for gaining access to records based totally on labels or positional integers, depending on the indexing method used ([], .iloc[], .loc[])!

```
import numpy as np
s4 = pd.Series(np.random.randint(10,100,20))
s4
```

```
0    11
1    56
2    21
3    43
4    41
5    22
6    51
7    33
8    41
9    68
10   35
11   96
12   14
13   66
14   99
15   58
16   31
17   52
18   67
```

```
s4[s4%2==0]
```

```
1    56
5    22
9    68
11   96
12   14
13   66
15   58
17   52
19   90
dtype: int64
```

```
s4[(s4>50) & (s4 <= 70)]
```

```
1    56
6    51
9    68
13   66
15   58
17   52
18   67
dtype: int64
```

## DataFrame

DataFrame is an essential aspect of Pandas and are widely used for statistics manipulation and evaluation obligations. They constitute tabular data with rows and columns, just like spreadsheets or database tables!!!

Creating a DataFrame can be executed in numerous approaches. Let's explore one approach using lists and dictionaries!

```
import pandas as pd

# Define lists containing data
cities = ["London", "Paris", "Berlin", "Madrid", "Rome"]
population = [8908081, 2148327, 3769495, 3266126, 2870493]

# Create a dictionary with keys "City" and "Population"
city_dict = {"City": cities, "Population": population}

# Create a DataFrame from the dictionary
city_df = pd.DataFrame(city_dict)
city_df
```

|   | City   | Population |
|---|--------|------------|
| 0 | London | 8908081    |
| 1 | Paris  | 2148327    |
| 2 | Berlin | 3769495    |
| 3 | Madrid | 3266126    |
| 4 | Rome   | 2870493    |

## Here's what befell:

Two lists, cities and populace, have been defined containing the names of towns and their respective populations.

These lists were then mixed into a dictionary called city\_dict, in which the keys are "City" and "Population".

Finally, the dictionary turned into handed into the pd.DataFrame() function to create a DataFrame object, city\_df, with columns labeled "City" and "Population". Pandas routinely assigned index labels from 0 to 4 to correspond with the quantity of elements in every list.

DataFrames are one of the fundamental data structures provided by the Pandas library in Python, widely used for data manipulation and analysis tasks. They are designed to handle two-dimensional data in a tabular format, similar to a spreadsheet or database table.

## Creating DataFrame:

DataFrame can be created using various methods:

- **From dictionaries:** Each key-value pair in the dictionary corresponds to a column in the DataFrame.
- **From lists:** Lists can represent either rows or columns of the DataFrame.
- **From external data sources:** Data can be read from CSV, Excel, SQL databases, or other file formats.

```
pd.read_csv('path_data/file_name.csv')
```



```
import pandas as pd

# Example 1: Creating DataFrame from dictionaries
data_dict = {
    'Name': ['John', 'Alice', 'Bob', 'Emily'],
    'Age': [25, 30, 35, 28],
    'City': ['New York', 'Paris', 'London', 'Berlin']
}
df_dict = pd.DataFrame(data_dict)
print("DataFrame created from dictionary:")
print(df_dict)
print()
```

DataFrame created from dictionary:

|   | Name  | Age | City     |
|---|-------|-----|----------|
| 0 | John  | 25  | New York |
| 1 | Alice | 30  | Paris    |
| 2 | Bob   | 35  | London   |
| 3 | Emily | 28  | Berlin   |

```
# Example 2: Creating DataFrame from lists
names = ['John', 'Alice', 'Bob', 'Emily']
ages = [25, 30, 35, 28]
cities = ['New York', 'Paris', 'London', 'Berlin']
data_list = list(zip(names, ages, cities))
df_list = pd.DataFrame(data_list, columns=['Name', 'Age', 'City'])
print("DataFrame created from lists:")
print(df_list)
print()
```

DataFrame created from lists:

|   | Name  | Age | City     |
|---|-------|-----|----------|
| 0 | John  | 25  | New York |
| 1 | Alice | 30  | Paris    |
| 2 | Bob   | 35  | London   |
| 3 | Emily | 28  | Berlin   |

## Accessing Data:

Once a DataFrame is created, data can be accessed using various methods:

- **Accessing columns:** Columns can be accessed using square brackets [] or dot notation.
- **Accessing rows:** Rows can be accessed using methods like loc[] and iloc[], which allow for label-based and integer-based indexing, respectively.

```
# Example 3: Creating DataFrame from external data source (CSV file)
df_csv = pd.read_csv('data.csv')
print("DataFrame created from CSV file:")
print(df_csv)
print()

# Accessing data
print("Accessing data:")
print("First row of 'Name' column:", df_dict['Name'][0])
print("Second row of 'Age' column:", df_list['Age'][1])
print("Data from row 1 to row 2 of 'City' column:", df_csv['City'][1:3])
print()
```

## Manipulating Data:

DataFrame offer numerous methods for manipulating data:

- **Adding and deleting columns:** Columns can be added using the assignment or the insert() method and deleted using the del keyword or the pop() method.
- **Adding and deleting rows:** Rows can be added using the append() method and deleted using the drop() method or Boolean indexing.
- **Modifying values:** Values in a DataFrame can be modified using assignment or various transformation functions.

```
# Manipulating data
# Adding a new column
df_dict['Gender'] = ['Male', 'Female', 'Male', 'Female']
print("After adding 'Gender' column:")
print(df_dict)
print()
```

After adding 'Gender' column:

|   | Name  | Age | City     | Gender |
|---|-------|-----|----------|--------|
| 0 | John  | 25  | New York | Male   |
| 1 | Alice | 30  | Paris    | Female |
| 2 | Bob   | 35  | London   | Male   |
| 3 | Emily | 28  | Berlin   | Female |

```
# Deleting a column
del df_list['City']
print("After deleting 'City' column:")
print(df_list)
print()
```

After deleting 'City' column:

|   | Name  | Age |
|---|-------|-----|
| 0 | John  | 25  |
| 1 | Alice | 30  |
| 2 | Bob   | 35  |
| 3 | Emily | 28  |

## Filtering and Querying Data:

DataFrame support filtering and querying operations:

- **Filtering rows:** Rows can be filtered based on specific conditions using Boolean indexing.
- **Querying:** DataFrame provide methods like query() for executing SQL-like queries on the DataFrame.

```
# Filtering data
filtered_df = df_csv[df_csv['Age'] > 30]
print("Filtered DataFrame where Age > 30:")
print(filtered_df)
# Filtering data
filtered_df = df_csv[df_csv['Age'] > 30]
print("Filtered DataFrame where Age > 30:")
print(filtered_df)
```

## Indexing and Labelling:

### DataFrame have two kinds of index:

- **Row index:** Each row is assigned a unique integer index by using default.
- **Column index:** Each column is categorised with a column name.

## Displaying Data:

When displaying a DataFrame, Pandas shows the records in a smartly formatted tabular structure, with rows and columns categorised.

```
DataFrame:  
   Name  Age  City  
0  John   25 New York  
1  Alice  30   Paris  
2   Bob   35  London  
3  Emily  28   Berlin
```

## General workflow for data analysis.

### Data Import / Availability:

- Importing the facts from diverse sources along with CSV files, databases, or APIs.
- Ensuring the supply and accessibility of the facts needed for evaluation.

```
# Step 1: Data Import / Availability  
import pandas as pd  
  
# Importing data from a CSV file  
data = pd.read_csv('dataset.csv')
```

### Data Understanding / Exploratory Data Analysis (EDA):

- Examining the primary traits of the dataset, together with the wide variety of rows and columns.
- Gathering metadata, which incorporates information about
- the statistic types of every column, the range of values, and any missing values.
- Gaining an expertise of the records via exploratory information evaluation strategies inclusive of summary information, records visualization, and correlation evaluation.

```
# Step 2: Data Understanding / Exploratory Data Analysis (EDA)
# Checking the number of rows and columns
print("Number of rows and columns:", data.shape)

# Gathering metadata
print("Data types of each column:")
print(data.dtypes)

# Summary statistics
print("Summary statistics:")
print(data.describe())

# Data Visualization
import matplotlib.pyplot as plt

# Histogram of a numerical variable
plt.hist(data['Age'], bins=20, color='skyblue')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.title('Distribution of Age')
plt.show()
```

## Data Cleaning:

- Performing shape-based modifications like sub setting, reordering variables, calculating new variables, renaming, and dropping variables.
- Handling facts type troubles via statistics type casting or conversion.
- Addressing content material or information-primarily based problems together with filtering, sorting, coping with duplicates, outliers, and lacking values.
- Grouping and binning information, as well as applying alterations to the statistics.

```
# Step 3: Data Cleaning
# Handling missing values
data.dropna(inplace=True)
```

## Data Summarization:

- Creating summaries of the statistics through calculations and aggregation features/techniques.
- Generating reviews, dashboards, or charts to summarize key findings and insights.

```
# Step 4: Data Summarization
# Calculating mean and median
mean_age = data['Age'].mean()
median_age = data['Age'].median()

print("Mean Age:", mean_age)
print("Median Age:", median_age)
```

## Data Visualization:

- Creating charts, graphs, and visualizations to represent the information correctly.
- Designing dashboards or reports for imparting the insights in a visually appealing manner.

```
# Data Visualization
import matplotlib.pyplot as plt

# Histogram of a numerical variable
plt.hist(data['Age'], bins=20, color='skyblue')
plt.xlabel('Age')
plt.ylabel('Frequency')
plt.title('Distribution of Age')
plt.show()
```

## Predictive Analytics:

- Utilizing superior analytics techniques to make predictions or forecasts primarily based at the records.
- Building predictive fashions, the use of gadget learning algorithms to discover styles or developments inside the information.

```
# Step 5: Predictive Analytics
# Building a simple linear regression model
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

# Splitting data into training and testing sets
X = data[['Feature1', 'Feature2']]
y = data['Target']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,

# Building the model
model = LinearRegression()
model.fit(X_train, y_train)

# Making predictions
y_pred = model.predict(X_test)

# Evaluating the model
mse = mean_squared_error(y_test, y_pred)
print("Mean Squared Error:", mse)
```

Each step on this workflow is critical for carrying out a radical and effective information evaluation method, from facts import to predictive analytics. It includes a combination of records manipulation, exploration, and visualization techniques to derive actionable insights from the records.



| Step                    | Description  | Code Snippet   |
|-------------------------|--|--|
| 1. Data Import          | Importing data from a file into a DataFrame                                    | <code>python data = pd.read_csv('dataset.csv')</code>  |
| 2. Data Understanding   | Exploring the dataset, checking dimensions, data types, and summary statistics | <code>python print("Number of rows and columns:", data.shape) &lt;br&gt; python print("Data types of each column:", data.dtypes) &lt;br&gt; python print("Summary statistics:", data.describe())</code>  |
| 3. Data Cleaning        | Handling missing values, duplicates, outliers, and other inconsistencies       | <code>python data.dropna(inplace=True)</code>  |
| 4. Data Summarization   | Calculating summary statistics, means, medians, etc.                           | <code>python mean_age = data['Age'].mean() &lt;br&gt; python median_age = data['Age'].median()</code>  |
| 5. Predictive Analytics | Building predictive models, evaluating model performance                       | <code>python model = LinearRegression() &lt;br&gt; python model.fit(X_train, y_train) &lt;br&gt; python y_pred = model.predict(X_test) &lt;br&gt; python mse = mean_squared_error(y_test, y_pred)</code> |

## DataFrame Operations in Pandas:

- **Adding a Column:** Use the project operator (=) to add a brand new column to a DataFrame. You can calculate values based on existing columns or provide a list of values at once.

```
import pandas as pd

# Creating a DataFrame
data = {
    'country_name': ['United Kingdom', 'France', 'Italy'],
    'population': [281, 119, 206],
    'median_age': [40, 42, 46]
}

data_df = pd.DataFrame(data)

# Adding the 'province' column
data_df['province'] = ['RM', 'MI', 'NA']

# Printing the DataFrame
print(data_df)
```

|   | country_name   | population | median_age | province |
|---|----------------|------------|------------|----------|
| 0 | United Kingdom | 281        | 40         | RM       |
| 1 | France         | 119        | 42         | MI       |
| 2 | Italy          | 206        | 46         | NA       |

- **Mathematical Operations:** Pandas helps mathematical operations on columns. You can apply capabilities from libraries like NumPy without delay to columns.

```
import numpy as np
np.log(data_df["population"])
```

|   |          |
|---|----------|
| 0 | 5.638355 |
| 1 | 4.779123 |
| 2 | 5.327876 |

Name: population, dtype: float64



- **Iterating over Rows:** Use the `.iterrows()` approach to iterate over the rows of a **DataFrame**. It returns the index and row information as tuples.

```
for index, row in data_df.iterrows():  
    print(index)  
    print(row)
```

```
0  
country_name    United Kingdom  
population      281  
median_age      40  
province        RM  
Name: 0, dtype: object  
1  
country_name    France  
population      119  
median_age      42  
province        MI  
Name: 1, dtype: object  
2  
country_name    Italy  
population      206  
median_age      46  
province        NA  
Name: 2, dtype: object
```

- **Transposing and Iterating over Columns:** Transpose a DataFrame the use of `.T` after which use `.iteritems()` to iterate over its columns. It returns column names and statistics series.

```
for country_name, population in data_df.T.items():  
    print(country_name)  
    print(population)
```

```
0  
country_name    United Kingdom  
population      281  
median_age      40  
province        RM  
Name: 0, dtype: object  
1  
country_name    France  
population      119  
median_age      42  
province        MI  
Name: 1, dtype: object  
2  
country_name    Italy  
population      206  
median_age      46  
province        NA  
Name: 2, dtype: object
```

- **Adding a Row:** Use the `append()` technique to feature a brand-new row to a DataFrame. Specify the row data as a dictionary and use `ignore_index=True` to reset the index.

```
# New row to append
new_row = {"country_name": "Spain", "population": 467, "median_age": 43, "province": "MD"}

# Append the new row
data_df = pd.concat([data_df, pd.DataFrame([new_row])], ignore_index=True)

print(data_df)
```

|   | country_name   | population | median_age | province |
|---|----------------|------------|------------|----------|
| 0 | United Kingdom | 281        | 40         | RM       |
| 1 | France         | 119        | 42         | MI       |
| 2 | Italy          | 206        | 46         | NA       |
| 3 | Spain          | 467        | 43         | MD       |

- **Handling Index:** When adding rows, consider using `ignore_index=True` to reset the index of the DataFrame.

```
# Append the new row using concat
new_data_df = pd.concat([data_df, pd.DataFrame([new_row])], ignore_index=True)

print(new_data_df)
```

|   | country_name   | population | median_age | province |
|---|----------------|------------|------------|----------|
| 0 | United Kingdom | 281        | 40         | RM       |
| 1 | France         | 119        | 42         | MI       |
| 2 | Italy          | 206        | 46         | NA       |
| 3 | Spain          | 467        | 43         | MD       |
| 4 | Spain          | 467        | 43         | MD       |

These operations allow for flexible statistics manipulation and evaluation inside Pandas DataFrames.

## Concatenating DataFrame:

Concatenation is the system of combining or greater DataFrame along both axes. Let's consider some examples:

### Concatenating alongside rows:

```
import pandas as pd

# Creating two DataFrames
df1 = pd.DataFrame({"A": [1, 2, 3], "B": [4, 5, 6]})
df2 = pd.DataFrame({"A": [7, 8, 9], "B": [10, 11, 12]})

# Concatenating along rows
result = pd.concat([df1, df2])
print(result)
```

|   | A | B  |
|---|---|----|
| 0 | 1 | 4  |
| 1 | 2 | 5  |
| 2 | 3 | 6  |
| 0 | 7 | 10 |
| 1 | 8 | 11 |
| 2 | 9 | 12 |

## Concatenating along columns:

```
import pandas as pd

# Creating two DataFrames
df1 = pd.DataFrame({"A": [1, 2, 3], "B": [4, 5, 6]})
df2 = pd.DataFrame({"C": [7, 8, 9], "D": [10, 11, 12]})

# Concatenating along columns
result = pd.concat([df1, df2], axis=1)
print(result)
```

|   | A | B | C | D  |
|---|---|---|---|----|
| 0 | 1 | 4 | 7 | 10 |
| 1 | 2 | 5 | 8 | 11 |
| 2 | 3 | 6 | 9 | 12 |

## Joining DataFrames:

Joining is used to combine columns from two doubtlessly different-listed DataFrames right into a single result DataFrame. Here are a few examples:

### Inner Join:

```
import pandas as pd

# Creating two DataFrames
left = pd.DataFrame({"A": [1, 2], "B": [3, 4]}, index=['a', 'b'])
right = pd.DataFrame({"C": [5, 6], "D": [7, 8]}, index=['a', 'c'])

# Inner join based on index
result = left.join(right, how='inner')
print(result)
```

|   | A | B | C | D |
|---|---|---|---|---|
| a | 1 | 3 | 5 | 7 |

### Left Join:

```
import pandas as pd

# Creating two DataFrames
left = pd.DataFrame({"A": [1, 2], "B": [3, 4]}, index=['a', 'b'])
right = pd.DataFrame({"C": [5, 6], "D": [7, 8]}, index=['a', 'c'])

# Left join based on index
result = left.join(right, how='left')
print(result)
```

|   | A | B | C   | D   |
|---|---|---|-----|-----|
| a | 1 | 3 | 5.0 | 7.0 |
| b | 2 | 4 | NaN | NaN |

## Merging DataFrames:

Merging is just like becoming a member of; however, it is extra versatile and lets in becoming a member of on columns in addition to indexes. Here are some examples:

## Inner Merge

```
import pandas as pd

# Creating two DataFrames
left = pd.DataFrame({"key": ['a', 'b'], "value": [1, 2]})
right = pd.DataFrame({"key": ['a', 'c'], "value": [3, 4]})

# Inner merge based on the 'key' column
result = pd.merge(left, right, on='key', how='inner')
print(result)
```

|   | key | value_x | value_y |
|---|-----|---------|---------|
| 0 | a   | 1       | 3       |

## Left Merge:

```
import pandas as pd

# Creating two DataFrames
left = pd.DataFrame({"key": ['a', 'b'], "value": [1, 2]})
right = pd.DataFrame({"key": ['a', 'c'], "value": [3, 4]})

# Left merge based on the 'key' column
result = pd.merge(left, right, on='key', how='left')
print(result)
```

|   | key | value_x | value_y |
|---|-----|---------|---------|
| 0 | a   | 1       | 3.0     |
| 1 | b   | 2       | NaN     |

## Identifying Missing Data:

We can use various techniques to identify missing values (NaN) in a DataFrame.

### Using isna() approach:

```
import pandas as pd

# Creating a DataFrame with missing values
data = {'A': [1, 2, None, 4, 5],
        'B': [None, 2, 3, 4, 5]}
df = pd.DataFrame(data)

# Check which elements are NaN
print(pd.isna(df))
```

|   | A     | B     |
|---|-------|-------|
| 0 | False | True  |
| 1 | False | False |
| 2 | True  | False |
| 3 | False | False |
| 4 | False | False |

Using notna() method:

```
import pandas as pd

# Creating a DataFrame with missing values
data = {'A': [1, 2, None, 4, 5],
        'B': [None, 2, 3, 4, 5]}
df = pd.DataFrame(data)

# Check which elements are not NaN
print(df['A'].notna())
```

|   |       |
|---|-------|
| 0 | True  |
| 1 | True  |
| 2 | False |
| 3 | True  |
| 4 | True  |

Name: A, dtype: bool

Using isnull() method (just like isna()):

```
import pandas as pd

# Creating a DataFrame with missing values
data = {'A': [1, 2, None, 4, 5],
        'B': [None, 2, 3, 4, 5]}
df = pd.DataFrame(data)

# Check which elements are NaN
print(df.isnull())
```

|   | A     | B     |
|---|-------|-------|
| 0 | False | True  |
| 1 | False | False |
| 2 | True  | False |
| 3 | False | False |
| 4 | False | False |

Dealing with Missing Data:

Once we have identified missing records, we will deal with it in numerous ways, such as dropping rows or columns containing NaN values.

Dropping NaN values in a column:

```
import pandas as pd

# Creating a DataFrame with missing values
data = {'A': [1, 2, None, 4, 5],
        'B': [None, 2, 3, 4, 5]}
df = pd.DataFrame(data)

# Drop NaN values in a specific column
print(df['A'].dropna())
```

|   |     |
|---|-----|
| 0 | 1.0 |
| 1 | 2.0 |
| 3 | 4.0 |
| 4 | 5.0 |

Name: A, dtype: float64



Dropping rows with any NaN values:

```
import pandas as pd

# Creating a DataFrame with missing values
data = {'A': [1, 2, None, 4, 5],
        'B': [None, 2, 3, 4, 5]}
df = pd.DataFrame(data)

# Drop rows with any NaN values
print(df.dropna())
```

|   | A   | B   |
|---|-----|-----|
| 1 | 2.0 | 2.0 |
| 3 | 4.0 | 4.0 |
| 4 | 5.0 | 5.0 |

Dropping rows where a specific column has NaN values:

```
import pandas as pd

# Creating a DataFrame with missing values
data = {'A': [1, 2, None, 4, 5],
        'B': [None, 2, 3, 4, 5]}
df = pd.DataFrame(data)

# Drop rows where column 'A' has NaN values
print(df[df['A'].notna()])
```

|   | A   | B   |
|---|-----|-----|
| 0 | 1.0 | NaN |
| 1 | 2.0 | 2.0 |
| 3 | 4.0 | 4.0 |
| 4 | 5.0 | 5.0 |

## DataFrame Methods

Now in the next instance, we can display some of the techniques we can practice in a DataFrame. Earlier we verified the sum technique, but pandas have lots more to offer, and here, we import the bundle seaborn and load the iris dataset that includes it giving us the facts in a DataFrame.

```
import seaborn as sns
iris = sns.load_dataset('iris')
iris.head()
```

|   | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1          | 3.5         | 1.4          | 0.2         | setosa  |
| 1 | 4.9          | 3.0         | 1.4          | 0.2         | setosa  |
| 2 | 4.7          | 3.2         | 1.3          | 0.2         | setosa  |
| 3 | 4.6          | 3.1         | 1.5          | 0.2         | setosa  |
| 4 | 5.0          | 3.6         | 1.4          | 0.2         | setosa  |



## DataFrame Exploration and Summary Statistics:

### head() and tail():

- **head():** Returns the first n rows of the DataFrame. Default is 5.
- **Tail():** Returns the ultimate n rows of the DataFrame. Default is 5.

```
iris.head()
iris.tail()
```

|     | sepal_length | sepal_width | petal_length | petal_width | species   |
|-----|--------------|-------------|--------------|-------------|-----------|
| 145 | 6.7          | 3.0         | 5.2          | 2.3         | virginica |
| 146 | 6.3          | 2.5         | 5.0          | 1.9         | virginica |
| 147 | 6.5          | 3.0         | 5.2          | 2.0         | virginica |
| 148 | 6.2          | 3.4         | 5.4          | 2.3         | virginica |
| 149 | 5.9          | 3.0         | 5.1          | 1.8         | virginica |

**Columns:** The columns characteristic returns the column labels (names) of the DataFrame.

**count():** count() returns the wide variety of non-null observations for every column inside the DataFrame.

```
iris.count()
```

```
sepal_length    150
sepal_width     150
petal_length    150
petal_width     150
species         150
dtype: int64
```

**describe():** The describe() method generates descriptive data that summarize the central tendency, dispersion, and shape of a dataset's distribution.

```
iris.describe()
```

|       | sepal_length | sepal_width | petal_length | petal_width |
|-------|--------------|-------------|--------------|-------------|
| count | 150.000000   | 150.000000  | 150.000000   | 150.000000  |
| mean  | 5.843333     | 3.057333    | 3.758000     | 1.199333    |
| std   | 0.828066     | 0.435866    | 1.765298     | 0.762238    |
| min   | 4.300000     | 2.000000    | 1.000000     | 0.100000    |
| 25%   | 5.100000     | 2.800000    | 1.600000     | 0.300000    |
| 50%   | 5.800000     | 3.000000    | 4.350000     | 1.300000    |
| 75%   | 6.400000     | 3.300000    | 5.100000     | 1.800000    |
| max   | 7.900000     | 4.400000    | 6.900000     | 2.500000    |

It offers many, many, well-known deviation, minimal, twenty fifth percentile (Q1), median (fiftieth percentile), seventy fifth percentile (Q3), and most values.

**max()**, **min()**, **imply()**, **median()**, **mode()**, **std()**, **sum()**, **var()**:

These strategies compute various precise statistics for the DataFrame or a specific column. For example:

**max()**: Returns the most price.

**min()**: Returns the minimum value.

**mean()**: Returns the suggest fee.

**median()**: Returns the median price.

**mode()**: Returns the maximum frequent fee.

**std()**: Returns the same old deviation.

**sum()**: Returns the sum of values.

**var()**: Returns the variance.

```
# Assuming iris is your DataFrame containing the Iris dataset

# Select numeric columns only
numeric_columns = iris.select_dtypes(include=['float64', 'int64'])

# Maximum
iris_max = numeric_columns.max()

# Minimum
iris_min = numeric_columns.min()

# Mean
iris_mean = numeric_columns.mean()

# Median
iris_median = numeric_columns.median()

# Mode
iris_mode = numeric_columns.mode().iloc[0] # Mode may return multiple values, so we take the first

# Standard Deviation
iris_std = numeric_columns.std()

# Sum
iris_sum = numeric_columns.sum()

# Variance
iris_var = numeric_columns.var()

# Printing the results
print("Maximum:")
print(iris_max)
print("\nMinimum:")
print(iris_min)
print("\nMean:")
print(iris_mean)
print("\nMedian:")
print(iris_median)
print("\nMode:")
print(iris_mode)
print("\nStandard Deviation:")
print(iris_std)
print("\nSum:")
print(iris_sum)
```

```
Maximum:
sepal_length    7.9
sepal_width     4.4
petal_length    6.9
petal_width     2.5
dtype: float64

Minimum:
sepal_length    4.3
sepal_width     2.8
petal_length    1.0
petal_width     0.1
dtype: float64

Mean:
sepal_length    5.843333
sepal_width     3.057333
petal_length    3.758000
petal_width     1.199333
dtype: float64

Median:
sepal_length    5.80
sepal_width     3.00
petal_length    4.35
petal_width     1.30
dtype: float64

Mode:
sepal_length    5.0
sepal_width     3.0
petal_length    1.4
petal_width     0.2
Name: 0, dtype: float64

Standard Deviation:
sepal_length    0.828066
sepal_width     0.435866
petal_length    1.765298
petal_width     0.762238
dtype: float64

Sum:
sepal_length    876.5
sepal_width     458.6
petal_length    563.7
petal_width     179.9
dtype: float64

Variance:
sepal_length    0.685694
sepal_width     0.189979
petal_length    3.116278
petal_width     0.581006
dtype: float64
```

## DataFrame Operations:

**corr():** The corr() method computes the pairwise correlation of columns, indicating the strength and course of the linear dating among variables.

```
# Assuming iris is your DataFrame containing the Iris dataset

# Select numeric columns only
numeric_columns = iris.select_dtypes(include=['float64', 'int64'])

# Calculate the correlation matrix
correlation_matrix = numeric_columns.corr()

# Printing the correlation matrix
print(correlation_matrix)
```

|              | sepal_length | sepal_width | petal_length | petal_width |
|--------------|--------------|-------------|--------------|-------------|
| sepal_length | 1.000000     | -0.117570   | 0.871754     | 0.817941    |
| sepal_width  | -0.117570    | 1.000000    | -0.428440    | -0.366126   |
| petal_length | 0.871754     | -0.428440   | 1.000000     | 0.962865    |
| petal_width  | 0.817941     | -0.366126   | 0.962865     | 1.000000    |

**cov():** The cov() method computes the covariance matrix for the DataFrame, which measures how plenty random variables change together.

```
# Assuming iris is your DataFrame containing the Iris dataset

# Select numeric columns only
numeric_columns = iris.select_dtypes(include=['float64', 'int64'])

# Calculate the covariance matrix
covariance_matrix = numeric_columns.cov()

# Printing the covariance matrix
print(covariance_matrix)
```

|              | sepal_length | sepal_width | petal_length | petal_width |
|--------------|--------------|-------------|--------------|-------------|
| sepal_length | 0.685694     | -0.042434   | 1.274315     | 0.516271    |
| sepal_width  | -0.042434    | 0.189979    | -0.329656    | -0.121639   |
| petal_length | 1.274315     | -0.329656   | 3.116278     | 1.295609    |
| petal_width  | 0.516271     | -0.121639   | 1.295609     | 0.581006    |

**cumsum():** The cumsum() method returns the cumulative sum of the elements along a given axis. Specific Column Operations: count number() on column:

```
iris.cumsum()
```

|     | sepal_length | sepal_width | petal_length | petal_width | species   |
|-----|--------------|-------------|--------------|-------------|---|
| 0   | 5.1          | 3.5         | 1.4          | 0.2         | setosa  |
| 1   | 10.0         | 6.5         | 2.8          | 0.4         | setosasetosa  |
| 2   | 14.7         | 9.7         | 4.1          | 0.6         | setosasetosasetosa                                      |
| 3   | 19.3         | 12.8        | 5.6          | 0.8         | setosasetosasetosasetosa                                |
| 4   | 24.3         | 16.4        | 7.0          | 1.0         | setosasetosasetosasetosasetosa                          |
| ... | ...          | ...         | ...          | ...         | ...   |
| 145 | 851.6        | 446.7       | 543.0        | 171.9       | setosasetosasetosasetosasetosasetosasetosasetosaseto... |
| 146 | 857.9        | 449.2       | 548.0        | 173.8       | setosasetosasetosasetosasetosasetosasetosasetosaseto... |
| 147 | 864.4        | 452.2       | 553.2        | 175.8       | setosasetosasetosasetosasetosasetosasetosasetosaseto... |
| 148 | 870.6        | 455.6       | 558.6        | 178.1       | setosasetosasetosasetosasetosasetosasetosasetosaseto... |
| 149 | 876.5        | 458.6       | 563.7        | 179.9       | setosasetosasetosasetosasetosasetosasetosasetosaseto... |

150 rows x 5 columns

**count()** can also be carried out to unique columns, providing the rely of non-null observations in that column.

```
iris['sepal_length'].count()  
150
```

## Accessing precise column for other operations:

Specific columns may be accessed the use of dot notation (e.G., iris.sepal\_length.mean()) or bracket notation (e.G., iris['sepal\_length'].mean()).

```
iris.sepal_length.mean()  
5.8433333333333334
```

These DataFrame techniques are important for facts exploration, summary data computation, and gaining insights into the dataset's traits and relationships between variables.

## Handling Missing Data in Pandas DataFrames:

### dropna():

This approach drops rows or columns containing missing values (NaN) from the DataFrame.

axis=0 drops' rows with NaN values.

axis=1 drops columns with NaN values.

### fillna():

**fillna()** replaces missing values with certain values, along with a consistent or the imply of the data.

data.fillna(value) fills all lacking values in the DataFrame with the specified value.

`data.fillna(data.mean())` fills lacking values with the imply of each column.

`data.fillna(approach='pad')` fills missing values with the preceding non-null cost.

`data.fillna(method='bfill')` fills missing values with the next non-null value.

## **interpolate():**

`interpolate()` is used to interpolate lacking values within the DataFrame.

**data.interpolate()** plays linear interpolation by means of default, filling NaN values with values linearly interpolated among non-NaN values.

Additional interpolation techniques may be specific the usage of the technique parameter, together with 'barycentric', 'pchip', 'akima', 'spline', or 'polynomial'.

## **Handling lacking records in Series:**

The identical strategies (`fillna()`, `interpolate()`) may be applied to Series items.

Interpolation can be in addition custom designed the use of non-obligatory arguments like 'limit', 'limit\_direction', and 'limit\_area' to control the interpolation behavior.

## **replace():**

`replace()` is used to replace values within the DataFrame.

`data.replace(to_replace, value)` replaces values laid out in `to_replace` with `value`.

It may be used to update precise NaN values or other values within the DataFrame.

These techniques provide flexibility in coping with missing statistics in Pandas DataFrames, permitting users to drop, fill, or interpolate lacking values based on their necessities.

Additionally, `replace()` presents a way to update precise values, which include NaN, with desired values.

```
#Handling Missing Data in Pandas DataFrames:
import pandas as pd
import numpy as np
import seaborn as sns

# Load the Iris dataset
iris = sns.load_dataset('iris')

# Displaying the first few rows of the dataset
print("Original DataFrame:")
print(iris.head())

# Creating some NaN entries in the DataFrame
iris.iloc[2:4, 0] = np.nan # Setting NaN values in the 'sepal_length' column
iris.iloc[1:3, 1] = np.nan # Setting NaN values in the 'sepal_width' column

# Displaying the DataFrame with NaN values
print("\nDataFrame with NaN values:")
print(iris.head())

# Dropping rows with NaN values
dropped_rows = iris.dropna(axis=0)
print("\nDataFrame after dropping rows with NaN values:")
print(dropped_rows.head())

# Dropping columns with NaN values
dropped_columns = iris.dropna(axis=1)
print("\nDataFrame after dropping columns with NaN values:")
print(dropped_columns.head())
```



```
# Filling NaN values with the mean of each column
filled_with_mean = iris.fillna(iris.mean())
print("\nDataFrame after filling NaN values with the mean of each column")
print(filled_with_mean.head())

# Interpolating NaN values
interpolated = iris.interpolate()
print("\nDataFrame after interpolating NaN values:")
print(interpolated.head())

# Replacing specific values
replaced = iris.replace(2.3, 2)
print("\nDataFrame after replacing specific values (2.3 with 2):")
print(replaced.head())
```

Original DataFrame:

|   | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1          | 3.5         | 1.4          | 0.2         | setosa  |
| 1 | 4.9          | 3.0         | 1.4          | 0.2         | setosa  |
| 2 | 4.7          | 3.2         | 1.3          | 0.2         | setosa  |
| 3 | 4.6          | 3.1         | 1.5          | 0.2         | setosa  |
| 4 | 5.0          | 3.6         | 1.4          | 0.2         | setosa  |

DataFrame with NaN values:

|   | sepal_length | sepal_width | petal_length | petal_width | species |
|---|--------------|-------------|--------------|-------------|---------|
| 0 | 5.1          | 3.5         | 1.4          | 0.2         | setosa  |
| 1 | 4.9          | NaN         | 1.4          | 0.2         | setosa  |
| 2 | NaN          | NaN         | 1.3          | 0.2         | setosa  |
| 3 | NaN          | 3.1         | 1.5          | 0.2         | setosa  |
| 4 | 5.0          | 3.6         | 1.4          | 0.2         | setosa  |

DataFrame after dropping columns with NaN values:

|   | petal_length | petal_width | species |
|---|--------------|-------------|---------|
| 0 | 1.4          | 0.2         | setosa  |
| 1 | 1.4          | 0.2         | setosa  |
| 2 | 1.3          | 0.2         | setosa  |
| 3 | 1.5          | 0.2         | setosa  |
| 4 | 1.4          | 0.2         | setosa  |

## Grouping:

Grouping entails splitting the information into organizations based on a few standards, which includes the values of one or more columns. The **groupby()** feature in pandas is used to perform grouping. When you apply **groupby()** to a DataFrame, it returns a GroupBy item, that is an intermediate step that allows you to carry out diverse operations at the agencies.

### Key points approximately grouping:

- **Splitting:** The information is break up into agencies primarily based on the specified criteria.
- **Applying:** A function is carried out to each group independently.
- **Combining:** The outcomes of the characteristic applied to every organization are mixed into a resulting facts shape.

## Aggregation:

Aggregation entails computing a precis statistic (e.g., sum, mean, depend) for every institution. It condenses the records right into a smaller, extra doable shape, considering less difficult evaluation and interpretation. Pandas affords several techniques for aggregation, which includes sum(), mean(), be count(), min(), max(), etc.



## Key factors about aggregation:

- **Summary facts:** Aggregation computes summary statistics for every group.
- **Reduces dimensionality:** Aggregation reduces the dimensionality of the facts by using summarizing multiple values into a single value.
- **Operates on businesses:** Aggregation capabilities operate on character agencies created at some point of the grouping step.

## Pivot Tables:

Pivot tables are a powerful device for information summarization and evaluation, commonly used in spreadsheet software program like Excel. In pandas, **pivot\_table()** is used to create pivot tables from Data-Frames. Pivot tables assist you to reorganize and summarize statistics, making it less difficult to investigate styles and relationships.

## Key factors about pivot tables:

- **Reshaping facts:** Pivot tables reshape information by means of aggregating and summarizing it consistent with person-certain criteria.
- **Multidimensional analysis:** Pivot tables allow multidimensional evaluation with the aid of permitting users to specify rows, columns, and values to combination.
- **Customizable:** Pivot tables are customizable, permitting users to specify distinct aggregation features, rows, columns, and values based on their analysis necessities.

## Applications:

Grouping and aggregation are generally used for exploratory records evaluation, summarizing information for reporting purposes, and generating insights from datasets.

Pivot tables are useful for studying relationships between variables, identifying trends, and summarizing big datasets in a compact and interpretable format.

In precis, grouping, aggregation, and pivot tables are vital techniques in pandas for organizing, summarizing, and reading facts, allowing users to gain valuable insights and make informed choices primarily based on their facts.

```
import seaborn as sns
import pandas as pd
import numpy as np

# Load the Iris dataset
iris = sns.load_dataset('iris')

# Grouping by species and calculating the sum
groupby = iris.groupby('species')
sum_result = groupby.sum()
print("Sum:")
print(sum_result)

# Grouping by species and calculating the mean
mean_result = groupby.mean()
print("\nMean:")
print(mean_result)

# Looping over the group and printing the name and contents of each group
for name, group in groupby:
    print("\n" + name)
    print(group.head())

# Applying aggregate function to the groupby object
aggregate_result = groupby.agg(np.sum)
print("\nAggregate:")
print(aggregate_result)

# Grouping by species with as_index set to False
sum_as_index_false = iris.groupby('species', as_index=False).sum()
print("\nSum with as_index=False:")
print(sum_as_index_false)

# Applying size() method to the groupby object
size_result = groupby.size()
print("\nSize:")
print(size_result)

# Applying describe() method to the groupby object on a specific column
describe_result = groupby['sepal_length'].describe()
print("\nDescribe:")
print(describe_result)
```

```
# Applying multiple numpy methods to the groupby object
multi_agg_result = groupby['sepal_length'].agg([np.sum, np.mean, np.std])
print("\nMultiple Aggregations:")
print(multi_agg_result)

# Applying a lambda function to the groupby object
lambda_agg_result = groupby.agg({'lambda x: np.std(x, ddof=1)})
print("\nLambda Aggregation:")
print(lambda_agg_result)

# Getting the largest and smallest values in each group
largest_result = groupby['sepal_length'].nlargest(3)
smallest_result = groupby['petal_length'].nsmallest(4)
print("\nLargest Values:")
print(largest_result)
print("\nSmallest Values:")
print(smallest_result)

# Applying a custom function to the groupby object using apply()
def custom_function(group):
    return pd.DataFrame({'original': group, "demeaned": group - group.mean()})

apply_result = groupby['petal_length'].apply(custom_function).head()
print("\nApply:")
print(apply_result)

# Using qcut to create equal-sized buckets and then grouping by it
factor = pd.qcut(iris['sepal_length'], [0, 0.25, 0.5, 0.75, 1])
qcut_grouped = iris.groupby(factor)
qcut_mean_result = qcut_grouped.mean()
print("\nQcut Grouped:")
print(qcut_mean_result)

# Loading the tips dataset
tips = sns.load_dataset('tips')

# Grouping by multiple columns
multi_groupby = tips.groupby(['sex', 'smoker'])
multi_groupby_sum = multi_groupby.sum()
print("\nMultiple Groupby:")
print(multi_groupby_sum)
```

```
# Using pivot_table to pivot the data
pivot_table_result1 = pd.pivot_table(tips, index=["sex"])
print("\nPivot Table 1:")
print(pivot_table_result1)

pivot_table_result2 = pd.pivot_table(tips, index=["sex", "smoker", "day"])
print("\nPivot Table 2:")
print(pivot_table_result2)

# Using pivot_table with aggfunc and values arguments
pivot_table_result3 = pd.pivot_table(tips, index=["sex", "smoker"], values=["tip"], columns=["day"], aggfunc=[np.mean])
print("\nPivot Table 3:")
print(pivot_table_result3)

# Using pivot_table with margins argument
pivot_table_result4 = pd.pivot_table(tips, index=["sex", "smoker"], values=["tip"], columns=["day"], aggfunc=[np.mean], margins=True)
print("\nPivot Table 4:")
print(pivot_table_result4)
```

|            |              |             |              |             |
|------------|--------------|-------------|--------------|-------------|
| Sum:       | sepal_length | sepal_width | petal_length | petal_width |
| species    |              |             |              |             |
| setosa     | 250.3        | 171.4       | 73.1         | 12.3        |
| versicolor | 296.8        | 138.5       | 213.0        | 66.3        |
| virginica  | 329.4        | 148.7       | 277.6        | 101.3       |
| Mean:      | sepal_length | sepal_width | petal_length | petal_width |
| species    |              |             |              |             |
| setosa     | 5.006        | 3.428       | 1.462        | 0.246       |
| versicolor | 5.936        | 2.770       | 4.260        | 1.326       |
| virginica  | 6.588        | 2.974       | 5.552        | 2.026       |

## Read Files

| Functionality              | Method                              | Description                                     |
|----------------------------|-------------------------------------|---|
| Reading CSV Files          | <code>'pd.read_csv()'</code>        | Reads data from CSV files into a DataFrame.     |
|                            | <code>'DataFrame.to_csv()'</code>   | Writes DataFrame data to a CSV file.            |
| Reading Other File Formats | <code>'pd.read_json()'</code>       | Reads data from JSON files into a DataFrame.    |
|                            | <code>'pd.read_excel()'</code>      | Reads data from Excel files into a DataFrame.   |
| Additional Writing Methods | <code>'DataFrame.to_json()'</code>  | Exports DataFrame data to a JSON file.          |
|                            | <code>'DataFrame.to_dict()'</code>  | Converts DataFrame data to a Python dictionary. |
|                            | <code>'DataFrame.to_html()'</code>  | Exports DataFrame data to an HTML file.         |
|                            | <code>'DataFrame.to_latex()'</code> | Converts DataFrame data to LaTeX format.        |
| Direct Writing to File     | <code>'DataFrame.to_csv()'</code>   | Writes DataFrame data directly to a CSV file.   |
|                            | <code>'DataFrame.to_json()'</code>  | Writes DataFrame data directly to a JSON file.  |
|                            | <code>'DataFrame.to_html()'</code>  | Writes DataFrame data directly to an HTML file. |
|                            | <code>'DataFrame.to_latex()'</code> | Writes DataFrame data directly to a LaTeX file. |

The furnished textual content offers a comprehensive assessment of studying and writing documents with pandas, covering diverse report codecs together with CSV, Excel, JSON, and extra. Here's a breakdown of the key factors protected inside the textual content:

### Reading and Writing CSV Files:

Pandas affords to\_csv() approach to jot down DataFrame to a CSV record, specifying index=False prevents the DataFrame index from being written.

To read records from a CSV record, read\_csv() technique is used, creating a DataFrame with the contents of the record.

### Reading Other File Formats:

Pandas gives strategies like read\_excel() and read\_json() to read Excel and JSON documents into DataFrames, respectively.

Example code demonstrates loading information from JSON and Excel files into DataFrames.

### General Delimited Files:

For widespread delimited documents, the read\_table() approach can be used.

## Additional Methods for Writing Data:

Besides `to_csv()`, pandas give methods like `to_json()`, `to_html()`, `to_latex()`, and greater for exporting facts to numerous formats.

Examples reveal converting DataFrame statistics to JSON, dictionary, HTML, LaTeX, and string codecs.

## Direct Writing to File:

Data can be directly written to documents the usage of methods like `to_json()`, `to_html()`, etc.

## Advantages of pandas:

Pandas simplifies information manipulation and analysis, presenting functionalities like becoming a member of, merging, grouping, and pivoting records.

It handles lacking statistics successfully and is capable of coping with big datasets.

## Compatibility and Integration:

Pandas integrates well with other Python applications, making it essential for Python programmers.

Overall, the text emphasizes the flexibility and comfort of pandas for information manipulation, analysis, and record I/O duties, highlighting its significance in the Python surroundings.

```
import pandas as pd
import seaborn as sns

# Reading and Writing CSV Files
tips = sns.load_dataset('tips')
tips.to_csv('myfile.csv', index=False)
data = pd.read_csv('myfile.csv')

# Reading Other File Formats
file_name_json = '/path/to/boston.json'
data_json = pd.read_json(file_name_json)

file_name_excel = '/path/to/boston.xlsx'
data_excel = pd.read_excel(file_name_excel)

# Additional Methods for Writing Data
json_data = tips.head().to_json()
dict_data = tips.head().to_dict()
html_data = tips.head().to_html()
latex_data = tips.head().to_latex()

# Direct Writing to File
tips.head().to_json('tips.json')
tips.head().to_html('tips.html')
tips.head().to_latex('tips.latex')
tips.head().to_latex('tips.tex')
```

## Conclusion:

DataFrame in Pandas are effective equipment for dealing with tabular facts, presenting a wide range of functionalities for records manipulation, evaluation, and visualization. Understanding the way to create, get admission to, and manipulate DataFrame is crucial for everybody running with information in Python.