

Basic Data Analysis

Learning Topics



- ✓ Display data
- ✓ Understanding data size
- ✓ Handling rows and indexes
- ✓ Handling columns and it's data types
- ✓ Observe null values
- ✓ Duplicate values.
- ✓ Access subset of data
- ✓ Boolean indexing

Understand the nature and characteristics of the data is the main step in Data Science or in Machine Learning. Some types of analysis techniques are available in pandas to understand the data. we will learn one by one in this chapter.

First load all the data as a pandas DataFrame to start the analysis. The below code loads all the csv file as pandas DataFrame. Now the analyzing data is in **df** (pandas DataFrame) variable.

```
import pandas as pd

df = pd.read_csv('File path')
df.head()
```

Display Data:

- There are three ways to display the data as follows.

<i>Keyword / Function</i>	<i>Use</i>
<code>print(df)</code>	<i>Display whole data</i>
<code>df.head(10)</code>	<i>Display starting rows</i>
<code>df.tail (10)</code>	<i>Display last rows</i>

- By default, head argument n = 5, so by default it displays first 5 rows. If we want to display specific number of rows then pass needful number.
- head() function can take negative values also. Suppose if n = -3 pass to head() then it will remove last 3 rows and display remaining rows.
- Vice versa tail. By default, tail also has n =5, so it displays last 5 rows.

Understanding data size:

- Size is the basic character of the data. We need to select the optimization coding techniques related to time and space based on the size information.

<i>Keyword / Function</i>	<i>Use</i>
<code>df.shape</code>	<i>Displays number of rows and columns.</i>
<code>df.size</code>	<i>Displays total number of labels / values</i>
<code>df.ndim</code>	<i>Displays data set dimension.</i>
<code>df.info()</code>	<i>Displays total information of the data</i>

Handling rows and indexes:

- Indexes of data is one of the important topics in the data handling. Here will learn about how to see the index numbers and how to set a particular column as index.

<i>Keyword / Function</i>	<i>Use</i>
<code>df.index</code>	<i>To know all dataset index numbers</i>
<code>df.set_index("clm_1" / 0)</code>	<i>To set a specific column value as index numbers.</i>
<code>df.drop(index=2, inplace=True)</code> <code>df.drop(index=[2,3,4], inplace=True)</code>	<i>Drop some specific rows based on index numbers.</i>

Handling columns and it's data types:

- Sometimes we need to change column names and data types of each column. In those cases, will use the below techniques.

<i>Keyword / Function</i>	<i>Use</i>
<code>df.columns</code>	<i>Provides all column names as series.</i>
<code>df.dtypes</code> / <code>df.info()</code>	<i>To know all columns data types</i>
<code>df['clm_1'].astype('dtype')</code>	<i>To change a particular column data type. (Column type casting)</i>
<code>df.astype({'clm_1': 'int', 'clm_2': 'str'})</code>	<i>Type caste the more than one column at a time.</i>
<code>df.drop(['clm_1', 'clm_2'], axis=1)</code> <code>df.drop(columns=['clm_1', 'clm_2'])</code>	<i>To drop the columns</i>
<code>df.info()</code>	<i>Displays all information of all columns.</i>

Observe null values:

- Understanding of null values is one of the important step in data preprocessing. The below functions will helpful to understand the null values structure and nature in data.

<i>Keyword / Function</i>	<i>Use</i>
<code>df.isna()</code> / <code>df.isnull()</code>	<i>Return the boolean values of data frame contains true where null values is present.</i>
<code>df.notna()</code> / <code>df.notnull()</code>	<i>Vise verse of isna() function</i>
<code>df.isna().sum()</code> / <code>df.notna().sum()</code>	<i>Return total number of null values in each column.</i>

Duplicate values:

- An important part of Data analysis is analyzing Duplicate Values and removing them. Pandas `duplicated()` method helps in analyzing duplicate values. It returns a Boolean Series denoting duplicate rows.

Syntax: `DataFrame.duplicated(subset=None, keep='first')`

subset: optional

- Only consider certain columns for identifying duplicates, by default use all of the columns.

Keep: {'first', 'last', False} : default 'first'

- first : Mark duplicates as True except for the first occurrence.
- last : Mark duplicates as True except for the last occurrence.

- False : Mark all duplicates as True.

Access subset of data

- Pandas provides 2 interesting functions ***iloc*** and ***loc*** for access subset of data.

Syntax: `DataFrame.iloc`

- `iloc[]` is purely integer index based selection by position. It takes row and column index numbers and returns subset of data based on those index numbers
- `.iloc[]` is primarily integer position based, but may also be used with a boolean array.
- **Note** : That contrary to usual python slices, the start included but the stop is excluded.

Keyword / Function	Use
<code>df.iloc[0]</code>	Returns 0 th row as pandas series.
<code>df.iloc[[0]]</code>	Returns 0 th row as pandas data frame.
<code>df.isna().sum() / df.notna().sum()</code>	Return total number of null values in each column.
<code>df.iloc[[rows_indexes], [cols_indexes]]</code> Ex: <code>df.iloc[[1,2], [0,5]]</code>	Return data frame with mentioned row and column indexes,
<code>df.iloc[1:3, 0:3]</code>	1 st row to 2 nd row and 0 th column to 2 nd column
<code>df.iloc[:3, 0:3]</code>	Starting row to 2 nd row and 0 th column to 2 nd column.
<code>df.iloc[10:, 0:3]</code>	10 th row to last row and 0 th column to 2 nd column
<code>df.iloc[:, 0:3]</code>	First row to last row and 0 th column to 2 nd column.

Syntax: `DataFrame.loc`

- `.loc[]` is purely integer index based selection by position. It takes row and column index numbers and returns subset of data based on those index numbers
- `.loc[]` is primarily label based, but may also be used with a Boolean array.
- `.loc[]` is very helpful when data set has string type of indexes.

- **Note** : that contrary to usual python slices, **both** the start and the stop are included.

<i>Keyword / Function</i>	<i>Use</i>
<code>df.loc[['str_1', 'str_2']]</code>	Returns data frame with selected string indexed rows only
<code>df.loc['str_1':'str_5', 'clm_2':'cls_5']</code>	Returns data frame with selected string indexed rows and columns only

Boolean indexing:

- Boolean indexing is the process of selecting the subset of data with boolean values based on some conditions.
- Comparison operators (<, >, <=, >=, ==, ==, !=, ~) are using to creating the boolean indexing.

Ex: `df[df['column_name'] > value]`

- `isin()` function Is one of the way to make *boolean Indexing*. This function is used to check multiple value in a single column.

Ex: `df[df['column_name'].isin([val_1, val_2, _ _ _, val_n])]`

- Applying boolean indexing on multiple columns is one of the scenarios for this requirement we need to separate each condition with parentheses “()”.

Ex: `df[(df['column_name'] > val) & () | () & ()]`

