

Machine Learning associate

CLASSMATE

Date _____
Page _____

Numpy:

- Numpy stands for numerical python, which is used to perform numerical calculation
- Numpy is used to handle arrays and matrices.

What is numpy?

- Numpy in python is a library/package that is used to work with arrays
- It is a data structure used to store multiple values.

Why numpy?

- Faster and more compact than lists
- consumes less memory
- data type can be specified
- high level mathematical functions can be used
- used to solve problems related to linear algebra, fourier transform, vectors etc.

What is an array

A)

Numpy arrays commonly used data structures in python that store data as a grid or a matrix and easy to access.

`import numpy as np`

`arr1 = np.array([6, 7, 8, 9])`
`print(arr1)`

O/P: [6 7 8 9]

Built-in function to generate numpy array

- ① arange: It is used to generate evenly spaced values within a given interval

Note: end value not included.

Ex: `np.arange(4)`

O/P: [0, 1, 2, 3]

- ② ones: returns an array with all ones

Ex: `np.ones(5)`

O/P: [1, 1, 1, 1, 1]

- ③ zeros: returns zeros

Ex: `np.zeros(5)`

O/P: [0, 0, 0, 0, 0]

- ④ Eye -- (Identity matrix)

Ex: `np.eye(2)` O/P [[1 0]

20. [0 1]]

(5) full:- It is used to generate array with some elements.

np. full ((2, 2), 10)

0 IP $\begin{bmatrix} [10, 10] \\ [10, 10] \end{bmatrix}$

(6) linspace:- It is used to generate linearly spaced values.

np.linspace (1, 10, retstep=True)

$\begin{bmatrix} 1 & 2 & 3 & \dots & 10 \end{bmatrix}$ different b/w
b/w 1 & 10 each value

(7) Random:-

a) randint:- It is used to generate random integers from given interval.

np.random(), randint (1, 100, 10)
 \hookrightarrow 10 number of Element

b) rand:- It generates random values b/w the range [0, 1]

⑩ [uniform distribution]

np.random.rand(10)

\hookrightarrow 10 elements print
b/w 0, 1

c) randn:- It is used to generate random numbers from normal distribution within a range [-3, 3]

np.random.randn(10)

\hookrightarrow 10 element in [-3, 3]

Indexing with 2-D array:-

`arr = np.array([1, 2, 3], [4, 5, 6], [7, 8, 9])`

arr

o/p: $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}$

fetching first row

`i> arr[0]` o/p [1 2 3]

fetching first column

`i> arr[:, 0]`

o/p : $\begin{bmatrix} 1 \\ 4 \\ 7 \end{bmatrix}$

1-D Array:-

`arr = np.array([1, 2, 3, 4, 5, 6])`

arr.ndim: dimension of an array

arr.size: number of elements

arr.dtype: data type of an array element

arr.shape: order of matrix

Ex: $(1, 6)$, $(2, 3)$

type casting: arr.astype("int")

list will convert to integers.

It is used to convert one data type into another data type.

→ Matrix multiplication: A @ B

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} @ \begin{bmatrix} 5 & 2 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} (15+2) & (2+0) \\ (5+4) & (6+12) \end{bmatrix}$$

→ Element wise multiplication: A * B

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} * \begin{bmatrix} 5 & 2 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 5 & 4 \\ 3 & 12 \end{bmatrix}$$

(VVT) Difference b/w matrix multiplication & element wise multiplication.

Pandas:-

- Pandas is a Python library used for working with data sets.
- It has functions for analyzing, cleaning, exploring, and manipulating data.

Why is Pandas?

Pandas allows us to analyze big data and make conclusions based on statistical theories.

Pandas can clean messy data sets, and make them readable and relevant.

- * Import pandas as pd

pandas series :- Its one dimensional array holding data of any type.

1) `A = pd.Series()`

2) `B = pd.Series([30], index=[30])`

pandas data frame :-

A data frame is a data structure that organizes data into a 2-dimensional table of rows and columns, much like a spreadsheet.

```
dict = { "name": ["John", "Sam"], "age": [21, 22],  
        "salary": [3000, 4000]}
```

```
data = pd.DataFrame(dict)
```

data

id	name	age	salary
0	John	21	3000
1	Sam	22	4000

* How to create copy :- `data1 = data.copy()`

* sort data :- `data1.sort_values("Salary")` → default asc

* permanent change :- use `Inplace = True`.

* to sort in desc :- `data1.sort_values("age", ascending = False, inplace = True.)`

* rename column :- `data1.rename({ "name of column" : "Renamed by emp" }, onlix = 1)`

* fetching column :- `data[["age"]]`

Loc and iloc

loc :- label based indexing which helps in fetching rows and columns.

iloc :- index / integer based indexing

$[:, 1:4]$

$\rightarrow \text{data}.loc[:, \text{column}] \rightarrow [1, 2, 3], :]$

$\rightarrow \text{data}.iloc[:, \text{column}] \rightarrow [:, ["\text{age}", "\text{job}"]]$
 $\rightarrow [103, 105], :]$

Reading file in pandas

`data = pd.read_csv("train.csv")`

• shape :- It is used to find no row & columns
`data.shape`

• Head :- `data.head()`

• tail :- `data.tail()`

• data.columns :- To fetch all the names of columns.

• data.dtypes :- used to find data type of each column.

- * `data.drop("column_name", axis=1, inplace=True)`
- * `data.select_dtypes(include=["float64", "int64"])`
↳ fetching numerical columns
- * `data.select_dtypes(include="object")`
fetching categorical data
- * for above TWO codes add `.columns` → To see only column name
- * `data.info()`
- * `data.describe()` → by default fetch numerical data.
- * `data.describe(include='object')`
→ To fetch categorical data
- * `data.columnname.unique()` → To fetch unique element [distinct element]
- * ~~data.~~ `columnname.value_counts()`
It is used to count of unique categories
- * `data.isnull().sum()` → To check missing value

How to filter data :-

- ① `data.loc[data['sex'] == 'male']`
fetch records related to male passenger
- ② `data.loc[data['Age'] > 60]`
fetching records where age is greater than 60
- ③ `shape` → add to above see no. of records
④ `count`
- ④ `data.loc[data['Age'] < 10].shape`
o/p : (62, 12)

`data1 = data.copy()`. → Creating a copy

How to create new column

- `data["column name"] = "Data instead column"`
- Eg:- `data["course"] = "Data science"`
- To combine Two column to one different column

`data["Family"] = data["SibSp"] + data["Parch"]`

drop the column:

data.drop ("course", axis=1, inplace=True)
 ↳ column name.

data.drop (["Col1", "Col2"], axis=1, inplace=True)
 ↳ ↳ dropping multiple columns

sort data

data.sort_values ("Age", axis=0, inplace=True)
 ↳ In ascending order.

data.sort_values ("Age", axis=0, ascending=False,
 inplace=True, ignore_index=True)
 ↳ In descending order.

Type casting in pandas:

- ① int can be converted into float and vice versa
- ② Object cannot be converted into float if data inside column is text.
- ③ Object can be converted into int or float if data inside column is numbers
- ④ int, float can be converted into object.
- ⑤ Type casting works only when data doesn't have missing values

① `data1.dtypes` → To check data types

② `data1["Fair"] = data1["Fair"].as_type("float64")`

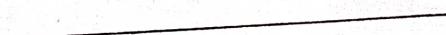
③ `data1["Family"] = data1["Family"].as_type("float64")`

Column name → To check data type to convert

Missing values

How to replace missing values

- ① numerical → mean →  → normal distribution

numerical → median →  (skewed):

- ② Categorical \rightarrow mode.

VVS If you find more than 80% data missing drop columns.

- * How to replace missing value.

- ① check for the distribution

```
import seaborn as sns
```

```
sns.histplot(x=data["Age"], kde=True)
```

0/0-

② Since distribution of age is not normal we replace missing values in age column with median.

`data["Age"].median()` \rightarrow 38.0

③ replace missing value in age column with median

`data.loc[data["Age"].isnull() == True, "Age"] = 38`

replaced successfully

④ replacing categorical column for missing value

⑤ Embarked data as missing value

`data[④ "Embarked"].mode()`

⑥ `data.loc[data["Embarked"].isnull() == True, "Embarked"] = 'S'`

If you get multiple mode means replace with any one of them.

combine two data frame.

`pd.concat([df1, df2], axis=0, ignore_index=True)`

throws manipulation error

Seaborn :- [Exploratory data Analysis (EDA)]

Page

① Univariate :- Analysis single variable

e.g: To check distribution :- histplot, distplot, countplot, boxplot.

② Bivariate :- Analysis two variables :-

e.g: line plot, scatterplot, pie, barplot, Box plot, violin, strip, swarm

③ Multivariate :- Analysing multiple variables

e.g: pairplot
heat map

Seaborn has a inbuilt dataset

→ sns.get_dataset_names()

→ Loading data:-

data = sns.load_dataset("tips")

① → Relation plot.

relation plots are used to understand relationship b/w two variables

① line plot (2 numerical)

② scatter - - -

③ relplot - - -

② Categorical plots:-

It is used to analyze one numerical and one categorical column.

-) count plot (1 cat & 1 cat)
-) box plot (1 num ⑥ 1 num & 1 cat)
-) violin plot (1 cat and 1 num) ⑥ 1 num
-) strip plot (1 num ⑥ 1 num & 1 cat)
-) swarm plot
-) cat plot.

③ Distribution plots:-

-) histplot
-) dist plot

These two graphs are used to check distribution of data.

④ Multivariant:-

-) pairplot
-) heat map

What is seaborn?

Seaborn is a python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics.

- Seaborn is built on top of matplotlib.