

Data Analysis & Visualization

Data Science

Manipulation and analysis of data.

Grouping of data can be of two types: structured and unstructured.
structured forms a class unstructured forms a cluster

- structured data is ordered.
- the data is classified & uses label
- eg: different cars like XUV, sedan, etc come under same class 'car' because they all are cars.
- Unstructured data is unordered.
- there are no labels.
- objects under unstructured data are ~~categ~~ grouped together in a cluster.

example:

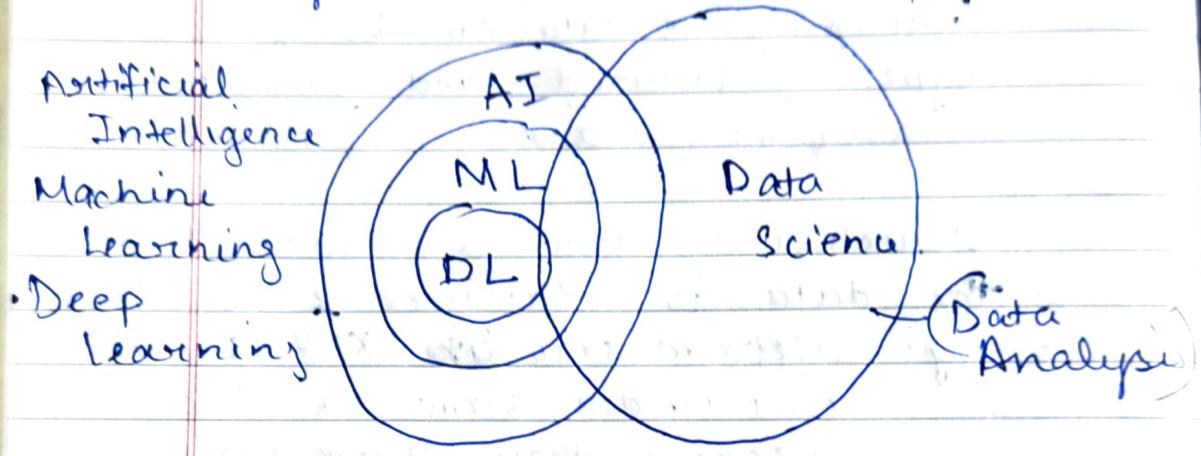
consider a bread.

In structured data, we will have class of breads. (similar objects) like wholewheat, brown bread, etc.

~~so~~ unstructured data will group together the things that go together with bread like eggs, butter, milk. so it will form a cluster of things which go together with bread.

Good Write

In structured data, we have labels like colour, shape, etc. but in unstructured data, we do not have labels, so we try to group together things based on similar properties.



Data Science contains AI, ML, and DL but is not limited to it. It is vast and can be done with other tools as well.

- EAD (exploratory data analysis).

- ↳ removal of unnecessary data
- ↳ find correlation b/w 2 data.

$$y = f(x)$$

$$f(n_1, n_2, n_3, \dots)$$

↳ feature vectors.

finding correlation b/w x and y

- data preprocessing: removal of noise (grainy) from photo comes under data preprocessing.

To reduce the calculation time,
we can consider everything other
than facial features to be noise
and crop out.

actual: $200 \times 200 \times 100 \times 60$

cropped: $60 \times 60 \times 100 \times 60$

less
calculation
time.

less no. of people
no. of img.

→ suppose we have n values and y values. If any of the y values is not there, what do we do? we drop it. we can't take avg. why take headache of prediction if we'll take avg. only & we can find correlation only when Good Write ~~avg~~ prediction is done. SO, we drop it.

That means, remove the data
(remove whole line including
 x & y values).

- If any x value is missing, we
can average or mod.

→ Exploratory Data Analysis

- what is data science?

Data science is a blend of various tools, algorithms and machine learning principles. Most simply, it involves obtaining meaningful information or insights from unstructured or unstructured data through a process of analysing programming and business skills. It is a field containing many elements like mathematics, statistics, computer science, etc.

- How data science works?

Data science is not a one step process, it passes from many stages and every element is important.
Following are the steps for the development of a data science model

I. Problem Statement.

- It is really important to declare or formulate your problem statement very clearly and it should be directed towards the pivotal goal to be achieved.

II. Data Collection.

- After defining the problem statement, the next obvious step is to go in search of data that your model might require. Data can be in any form, i.e. unstructured or structured.
It can be in various forms like Videos, spreadsheets, coded form, etc.

III. Data Cleaning

- Data cleaning is all about the removal of missing, redundant, unnecessary and duplicate data from your collection various tools such as numpy and pandas from python are used for data cleaning.

IV. Data analysis and exploration

- It's all about analysing the structure of data, finding hidden patterns in them, studying behaviour, visualising the effects of one variable over the others and then concluding we can use matlib and seaborn libraries of python for visualisation

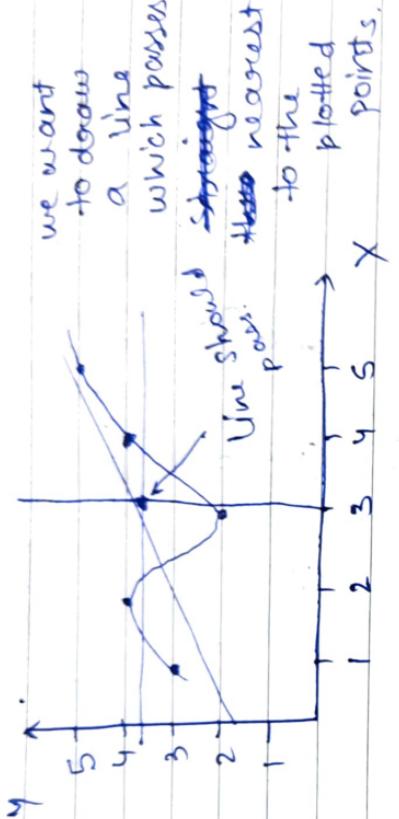
V Data Model

Once the study have been performed and detection have been made for data visualization , a hypothesis model must be build that may yield a good prediction for the target value .
There are different types of algorithms from regression to classification -
SPN ,

decision tree , or
you can also form clusters based on unsupervised learning .
After creation of the model , the model must be trained using training data and its predictive analysis can be done using test data .

Numerical:

$$\begin{aligned} X &= [1, 2, 3, 4, 5] \\ Y &= [3, 4, 2, 4, 3] \end{aligned}$$



$$y = mx + c . \quad m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$\bar{x} = 3, \quad \bar{y} = 3.6$$

$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$
-2	3	-0.6
-1	4	0.4
0	2	-1.6
1	4	0.4
2	5	1.4

$$\sum (x - \bar{x})^2 = 4 + 1 + 1 + 4 = 10.$$

$$m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{4}{10} = \underline{\underline{0.4}}$$

Good Write

$$\hat{y} = my + c.$$

$$\hat{y} = m\bar{x} + c.$$

$$3.6 = 0.4 \times 3 + c.$$

$$3.6 = 1.2 + c$$

$$c = 2.4$$

\bar{x}	\bar{y}	D_{xy}
1	2.8	2.8
2	3.2	3.2
3	3.6	2
4	4	4
5	4.4	5

f predicted

$$R^2 \text{ score : } \frac{\sum (y_p - \bar{y})^2}{\sum (y_i - \bar{y})^2}$$

If predicted value of y_p is equal to y_i , then it is best. So, it should be near 1.

Because

$$(y_p - \bar{y})^2 : \begin{aligned} 2.8 - 3.6 &= -0.8 \\ 3.2 - 3.6 &= -0.4 \\ 3.6 - 3.6 &= 0 \\ 4 - 3.6 &= 0.4 \\ 4.4 - 3.6 &= 0.8 \end{aligned}$$

$$\text{Sum} = 0.8$$

$$1.6$$

$$(y - \hat{y})^2 : \begin{aligned} 3 - 3.6 &= -0.6 = 0.36 \\ 4 - 3.6 &= 0.4 = 0.16 \\ 2 - 3.6 &= -1.6 = 0.25 \\ 4 - 3.6 &= 0.4 = 0.16 \\ 5 - 3.6 &= 1.4 = 1.96 \end{aligned}$$

$$0.36 + 0.16 + 0.25 + 0.16 + 1.96 = 3.20$$

$$\text{R-Score} = \frac{1.60}{5.20} = \frac{+6}{13} = 0.3$$

→ play.tensorflow

→ plt.scatter (x, y) is used to plot the true values (or expenditure) to be

→ plt.plot : used to plot predicted value of y w.r.t x

→ (agai, search cat and dog .CSV data)
 it will have series of x's and y.
 figure and which one are important
 find correlation b/w x & y.
 methods :
 1. CSV
 2. plot.

Pandas

Series vs Dataframe.

* Series

- a series is a one-D array like an object, a sequence of values and an associated array of data labels called as its index.

eg: `serarray = pd.Series([1, 2, 3, 4, 5])`

- The string representation of a series shows the index on the left and the values on the right. since the index is not specified, a default one consisting of the integers 0 through $n-1$ is created.

— manual indexing eg:-

```
obj = pd.Series([2, -3, 1, 4],  
               index = ['a', 'b',  
                      'd', 'c'])
```

we can manually specify the index by writing the keyword `index` and specifying the indices.
Also, `names = [1, 2, 3, 4, 5]`
`index = names`

* Dataframe.

- A Dataframe represents a rectangular table of data and contains an ordered collections of columns, each of which can be a different value type. The Dataframe has both a row & column index. Under the hood, the data is stored as one or more 2-D blocks rather than a list, dictionary or some other collection of one-dimensional arrays.

- eg: obj = pd.DataFrame([1, 2, 3], [4, 5, 6])

```
data = {'state': ['Ohio', 'Delhi', 'Texas'],
        'year': [1996, 2000, 2004],
        'pop': [15, 17, 32]}
```

```
obj2 = pd.DataFrame(data)
```

- Columns have a name, rows are generally indexed.

- frame = pd.DataFrame(data, columns=['State', 'Year', 'Pop'], index=[1, 2, 3, 4])

* Index objects

 ↑ arrange

```
obj = pd.Series([range(3), index=[  
                  'A', 'B', 'C']])
```

```
index = obj.index // reading index  
// reading particular index.  
index[1]          → a  
index[1:]        → b c
```

// changing index.

```
index[1] = 'd' // not possible.
```

→ In pandas, indexes are immutable. We can not change them.

* Index objects are immutable and thus cannot be modified by the user.

```
labels = pd.Index([np.arange(3)])
```

```
obj2 = pd.Series([1.5, 2.5, 3.0],  
                  index=labels)
```

* Essential functionality

1. Re-indexing: Re-index creates a new object with the data conformed to a new index

```
obj = pd.Series(['a:2', 1, 3, 4, 6],  
                index=[pd.Timestamp('2013-1-1'),  
                       pd.Timestamp('2013-1-2'),  
                       pd.Timestamp('2013-1-3'),  
                       pd.Timestamp('2013-1-4'),  
                       pd.Timestamp('2013-1-5')])  
  
obj2 = obj.reIndex(['d','b','c','e'],  
                    method='ffill')  
  
# index of obj originally is not  
# changed. rather, a new object  
# with new indexes is created.
```

```
obj2:  
d    0.2  
b    1.0  
a    3.4  
c    6.0  
e    NaN  
  
# fill  
1. ffill  
2. bfill  
3. bfill.
```

Operations

- ffill - forward fill , fills upward value.
- bfill - backward fill , fills downward value.

- Dropping entries from axis

→ Dropping one or more entries from an axis is easy if you already have an index array or list without those entries. For doing this, the 'Drop' method will be used.

→ The Drop method will return a new object with the indicated value deleted from an axis.

```
obj = pd.Series([1, 2, 3, 4, 5],  
               index=['a', 'b',  
                      'c',  
                      'd', 'e'])
```

```
obj2 = obj.drop('a') // To drop  
OR  
      obj.drop(['a', 'd']) // To drop two items
```

```
dataframe = np.arange(16).reshape(4, 4)  
obj3 = pd.DataFrame(data=  
                     {'a': [1, 2, 3, 4],  
                      'b': [5, 6, 7, 8],  
                      'c': [9, 10, 11, 12],  
                      'd': [13, 14, 15, 16]},  
                     columns=['one', 'two', 'three', 'four'])  
  
// dropping column.  
obj3.drop('three')
```

→ You can drop values from the column by passing axis = 1 or axis = columns.

- iloc : index location
takes numerical value.
- loc : takes name values.
 - data loading → storage and file formats:

- Handling dates and other custom types require extra offset.
for eg: data saved in any format
for example a .csv, a .txt,
can be easily read using
python's pandas library.
The only thing that is kept in
mind is how the values are
separated. They can be
separated by a tab, ','
, space or even a ' ; '.
All formats are readable using
pandas library.

- Reading a .csv file :

```
cars = pd.read_csv("cars.csv")
```

- only ',' separated files are
read as .csv.
excel is just a way to access
a csv file.

Good Write

- we can also use 'gredd-table' for reading such values but we have to specify the delimiter.

→ if label is missing, then
specify header = None.
otherwise first row will be
considered as header names.

→ changing in file :

```
cars = pd.read_csv("cars.csv",  
                   names = [" ", " ", " ", " ", " "],  
                   index_col = 'message')
```

if the separator between values is a variable space, we can use the regular expression "S+" written

as 'ls+'

Suppose you want to skip some rows while reading a csv file use

skiprows = [2, 0, 3]

↑
These rows will be skipped.

→ Reading txt files into pieces.

df = pd.read_csv("cars.csv",
nrows=10)

↑ no. of rows from start.

- Suppose we want to divide the given csv into several chunks :

chunksize = 5

↓
last chunk will contain all remaining rows.

→ Storing data

- one of the easiest way to store data efficiently in a serialized manner, we can use python's built-in command, i.e. 'pickle'

data.to_pickle("data.pkl")

↳ or any path.

if you want to read retrieve it.

`pd.read_pickle ("cars.csv")`

→ Reading excel file.

for reading an excel file, we need the syntax: `pd.read_excel()`

`data = pd.read_excel ("")`

→ Discretization and Binning.

- Binning is the process of creating bins or containers of specific ranges.

`ages = [20, 21, 32, 29, 41, 18, ...]`

`bins = [18, 25, 35, 60, 100]`

`cats = pd.cut(ages, bins)`

`Cats` :

Output : $(18, 25]$, $(18, 25]$, $(25, 35]$...
 (20 belongs) (21 belongs) (32 belongs)

each value of age is categorized according to the bin range.

pd.value_counts (cats)

- calculates values contained in each bin.
- // labeling the bins

group-names = ["Young", "Young-adult",
"adult", ""]

pd.cut (ages, bins, labels= group-names)

// to make 4 random bins.

~~dates~~ pd.cut (data, 4, precision:
↑
4 random bins.)

// To make bins with equal no. of
// values in them.

data1 = np.random.randint(1000)
cats = pd.qcut (data1, 4).

→ qcut will fail if we have
the same number in abundance.
suppose we have the no.'2' 8
times, and the rest of the no.'s
a single time. In that case,
qcut won't divide into equal
bins & its purpose will fail!

→ Permutation & random sampler

`df = pd.DataFrame(np.arange(5*4).
reshape((5,4)))`

`sampler = np.random.permutation(5)`

`df.`

// sampler will give random order
// permutation is used for taking
5 numbers randomly
// random.permutation is used
for taking numbers randomly
which are ordered randomly
as well.

eg: instead of 0,1,2,3,4
we have 0,4,3,1,2

→ 1-HOT encoding.

Suppose we have data categorical
data with strings and values.
We cannot use numpy & pandas
library on it. So, we convert
the categorical data to numerical
data.

Let us consider following eg.

Good Write

India
China
Japan

now, if
data
Ja - 3
library
value n
But it

Now,
data
from
Jute

In 1 1
ch 2 0
Ja 3 0

this
0 is
no

synt

df =

pd
Good

India	Jute
China	silk
Japan	cotton

now, if we convert it to numeric data and assign In-1, Ch-2, Ja-3, Numpy or pandas library will consider the value numerically like 1 < 2 < 3. But it is not that kind of value.

Now, we will encode the data into a matrix. In the form of true and false values
 Jute silk cotton.

In 1.	1	0	0
Ch 2	0	1	0
Ja 3	0	0	1.

this way, we can say that 0 is false & 1 is true, now we can use numeric libraries.

Syntax:

```
df = pd.DataFrame(  
    'key': ['b', 'b', 'a', 'c', 'a', 'b',  
    'data2': orange(6) ] ).
```

pd.get_dummies(df["key"])

→ Regular Expression.

```
import re  
txt = 'on the horizon'  
x = re.findall('on', txt)  
print(x)  
// finding all occurrences of "on" in tri  
['on', 'on']
```

// checking if a string is present or not

```
txt1 = "- - "  
x1 = re.search("python", txt1)  
print(x1.start())  
↳ // tells what is the "post" of the first occurrence
```

// Separating words of a string

```
txt = " "
```

```
x2 = re.split("\s", txt)
```

// x2 is a list which contains each word of the string.

// to check if a string starts with a particular word, use '^'.

$x_3 = \text{re.findall} ("^python", txt)$
↳ checking this word.

```
if  $x_3$  :  
    print ("Yes")  
else :  
    print ("No match").
```

— \$ // to check if a string ends with a word, use '\$'

$x_4 = \text{re.findall} ("tutorial$", txt)$
if x_4 :
 print ("Yes")
else :
 print ("No match").

// to check if any of the words exist in string, use OR '|'.
↳

$x_5 = \text{re.findall} ("python|tutorial", txt)$
if x_5 :
 print ("Yes")
else :
 print ("No match").

// to check if there is one or more occurrence of the word, use '+'.

$x_7 = \text{re.findall}("thi+", \text{txt})$.

if x_7 :

 print ("Yes")

else:

 print ("No match").

→ group by

To find the mean for a particular column,

$\text{df}["\text{column name}"].mean()$

Now to find mean of a particular column which has a particular value in another column.

→ $\text{df}[\text{df}["\text{symbols}"] = "RIVN"] ["\text{close}"].m$

∴ now, if we have different categories (RIVN, LCID ...), then we have to write the statement multiple times to get mean of each category separately.

→ $\text{df}[\text{df}["\text{symbols}"] = "LCID"] ["\text{close}"].mean$

Good Write

Simplifying this expression using
groupby

→ df.groupby(["symbol"])["close"].mean

→ read titanic.csv file.

- to fill null age values.

titanic['age'] = titanic["age"].

replace(['?', [None]]).

astype['float']

∴ (float type for months).

reading csv and select 4 columns.

ite. → df = pd.read_csv("titanic").

~~titanic~~ → titanic = df[["Survived", "name", "Sex", "Age"]].

gbo = titanic.groupby('sex')

↳ or (by='sex')

- gbo.ngroup (to mention no. of groups)

- gbo.get_group("female")

- to do any type of work with a particular category.