

Assignment - 2

Q-1

We have Student dataset with attributes Name, Enrollment No, Semester, CGPA, Result status.

→

First we have to import Pandas module & the dataset.

```
import Pandas as pd.
```

```
df = pd.read_csv("sample.csv")
```

→

There are mainly 4 techniques to handle missing data

1) Delete the rows with missing values.

```
# deleting rows if all vals are missing
```

```
df.dropna(how = "all")
```

2) Replace the missing vals with

mean or median

```
df[['CGPA']].fillna(value = df['CGPA'].mean(), inplace = True)
```

```
df[['Semester']].fillna(value = df
```

```
['Semester'].median(),
```

```
inplace = True)
```

3) Imputation method for categorical columns

```
df['name'].fillna('Unknown',  
inplace = True).
```

Q4) If some data has longitudinal behaviour, then it will make sense to use the last valid observation

df['semester'], fillna(method='ffill')

Q-2

Random sampling refers to randomly selecting rows from the dataframes.

→ sample() function is used to get sample data from the provided data frame.

→ Example (i) :-

```
import pandas as pd
df = pd.read_csv("Data.csv")
# To get 1 row randomly from df
```

```
df.sample(1)
```

To get 3 rows of sample data
df.sample(n=3)

To get 70% random data.

```
df.sample(frac=0.7)
```

when we want no. of rows of

Sample that may exceed available rows.

df.sample(n=100, replace = True).

Q-3

rand & randn in numpy :

→ randn() : It generates an array of specified shape & fills it with random float value as per standard normal distribution.

→ Mean of all values is 0 & Variance is 1.

→ rand() : It generates an array of specified shape & fills it with random val 1.

It doesn't fill value as per std normal distribution.

(2) Join and Merge in Pandas :-

Join () : It combines 2 data frames on the bases of their indices.

Merge() :- It is more versatile & allows us to specify columns beside the index to join on in both dataframes.

Q-4 Attributes of the database :- Name, En. no, Sem, CGPA, Result Status.

```
import Pandas as pd
idt = pd.read_csv("sample.csv")
```

Grouping the data by semester & counting

Freq. for each sem

```
grp = idt.groupby('sem').count()
print(grp)
```

Q-5 n-grams :- They are continuous sequences of words or symbols or tokens in a document.

- Bag of words simply refers to a matrix where rows are docs & columns are words.

- whereas, TF gives the freq. of the word in each doc. & IDF calculate the weight of rare words.

Q-6) Regular expressions are specially encoded text strings used as patterns for matching sets of strings.

→ Conditions for valid MGTU enrollment no. :-

- length Should be 12

- contains only no.

→ Code for validating enrollment no :-

```
import re
test = input()
# for second condition,
def check1(test):
    if re.match("^[0-9]", test)
        return False.
    else return True.
```

For first condition

```
def check2(test):
    if (test... len(c)) == 12
        return True
    else return False.
```

if (check 1 (test) and check 2 (test)):

 Paint ("Valid enrollment no.")

else

 Paint ("Invalid enrollment no.")