

PDS Assignment 3

1) What do you mean by missing values? Explain the different way to handle the missing value with example & explain how to deal with missing

→ Data can in pandas.

→ Data can have missing values for a number of reasons such as observation that were not recorded and data corruption. Handling missing data is important as many machine learning algorithms do not support data with missing values.

• In Python specifically pandas, numpy and sklearn we mark missing values as `nan` values with `nan` value are ignored from arithmetic operation.

• Pandas dataframe provide a function `isnull()` it returns a new dataframe of same size as calling dataframe it contains only True or False with True at place `nan` in original dataframe and False at other place.

→ Ex

→

```
import pandas as pd
import numpy as np
```

```
df = pd.DataFrame(np.random.randn(3,3),
                  index=['a', 'c', 'e'], columns=['one', 'two', 'three'])
```

```
df = df.reindex(['b', 'c', 'e'])
```

```
Print: ("now replace with 0.")
```

```
Print (df.fillna(0))
```

o/p

	one	two	three
a	0.523	0.869	0.725
b	nan	nan	nan
c	0.869	0.324	0.123

- 2] Explain Dataframe in Pandas
- A Pandas frame is a 2 dimensional data structure like a 2-D array or a table with row and column.
 - Pandas dataframe is 2D size-mutable potentially heterogeneous tabular data with table axes.
 - A Data frame is a 2-dimensional data structure i.e. ~~data~~ data aligned in a tabular fashion in rows and columns.
 - Pandas dataframe consists of three principal components: the data, row and column.

Ex

```
import pandas as pd
```

```
data = {'calories': [420, 380, 390]}
```

```
duration = [50, 45, 40]
```

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

```
print(df)
```

O/p	calories	duration
0	420	50
1	380	40
2	390	45

3) Compare numpy vs Pandas

Pandas	numpy
<ul style="list-style-type: none">• Pandas is mostly use for data analysis.	<ul style="list-style-type: none">• numpy is mostly work with numerical values.
<ul style="list-style-type: none">• It consumes more memory.	<ul style="list-style-type: none">• It consumes less memory.
<ul style="list-style-type: none">• Indexing of Pandas series is very slow.	<ul style="list-style-type: none">• Indexing of numpy array is fast.
<ul style="list-style-type: none">• Pandas have 2d label object called DataFrame.	<ul style="list-style-type: none">• numpy provides a multidimensional array.
<ul style="list-style-type: none">• DataFrame and series are most powerful tools for Pandas.	<ul style="list-style-type: none">• Array is a most powerful for numpy.
<ul style="list-style-type: none">• It is used in higher industry.	<ul style="list-style-type: none">• It is used in lower industry.

4) Explain TF-IDF technique.

→ TF-IDF stands for Term Frequency Inverse data frequency of method.

- It can be defined as calculation of how relevant a word in series of dataset.

- The meaning increases proportionally to the number of time in text a word frequency in the dataset.

• Time Frequency (TF) - It gives w. frequency of word in each document.

• Inverse Date Frequency (IDF) - IDF is u/e to calculate the weight of word across all documents.

→ ~~eg~~ Import numpy & pandas, num from sklearn. feature - extraction text import. String = ['geeks geeks'] * 5
 +fidf = TfidfVectorizer()
 result = +fidf.fit_transform(string)
 report =
 Print (+fidf.vocabulary)
 print (result)

O/p

word index	tf-idf value	
0	(0,0)	1.0
1	(1,0)	1.0
2	(1,0)	1.0
3	(3,0)	1.0
4	(4,0)	1.0

5] Differentiate rand() and randn() function in numpy

rand()	randn()
- generate uniform random number	- generate a gaussian random number
- rand() have one random number	- A randn() have n random number
- A rand() have a positive random	- A randn() have negative random