

HANDLING MISSING VALUES :

MISSING VALUES :

- In python, specifically Pandas, Numpy and scikit learn, we mark missing values as NaN
- Values with a NaN value are ignored from operations like sum, count etc.
- We can mark values as NaN easily with the Pandas DataFrame by using the replace() function on a subset of the columns we are interested in.

HOW TO IDENTIFY MISSING VALUES ?

- df.isnull() - It returns True if the values are missing.

HOW TO DEAL WITH MISSING VALUES?

→ Drop the missing values.

↳ Either drop the entire row.

```
df.dropna('col-name', axis=0, inplace=True)
```

↳ Or drop the entire column.

```
df.dropna('col-name', axis=1, inplace=True)
```

→ Replace the missing value.

↳ Use Business understanding.

↳ Statistical methods - Imputation

(Mean, Median, Mode)

```
df['col-name'].fillna(df.col-name.mean(),
```

inplace=True)

(127)

I/P - import numpy as np.

Import pandas as pd.

I/P - df = pd.read_csv('E:\data.csv')

df

O/P -

	FirstName	lastName	Age	Sex	Post Test score	Test score	locatn
0	abc	mno	12.0	m	90.0	65.0	NAN
1	NAN	NAN	NAN	NAN	90.0	?	NAN
2	ghi	pqr	12.0	f	-	65.0	?
3	jkl	stu	12.0	f	90.0	62.0	NAN
4	mno	rwx	12.0	fm	89.0	63.0	NAN

I/P - df.shape

O/P - (5, 7)

I/P - df.info()

128

1/p - missing_val = ['n/a', '-', '?']

```
df = pd.read_csv('E:\data.csv', na_values  
                  ('NaN', 'n/a'))
```

df.head()

Op-

	First Name	Last Name	Age	Sex	Pre Test Score	Post Test Score	Location
0	abc	mno	12.0	m	90.0	65.0	NaN
1	Nan	Nan	Nan	Nan	90.0	Nan	NaN
2	ghi	pqr	12.0	f	Nan	65.0	NaN
3	jkl	stu	12.0	f	90.0	62.0	NaN
4	mno	vwz	12.0	m	89.0	63.0	NaN

HOW TO IDENTIFY MISSING VALUES?

1/p - df. esnull()

O/P - First Last Age Sex Pre Test Post Test location
Name Name Score Score

(129)

Identifying missing values in columns

1/p - df.isnull.sum()

O/p - First Name 1

Last Name 1

Age 1

Sex 1

PreTest Score 1

PostTest Score 1

location 5

dtype: int64

Columns with atleast one missing value

1/p - df.isnull.any(axis=0)

O/p - First Name True

Last Name True

Age True

Sex True

PreTest Score True

PostTest Score True

Location True

dtype: bool True

Columns with all missing values.

I/p - df.isnull().all(axis=0)

O/p - First Name False
Last Name False
Age False
Sex False
Pre Test Score False
Post Test Score False
Location True
dtype: bool

Number of columns with all missing values

I/p - df.isnull.all(axis=0).sum()

O/p - 1

Rows with atleast one missing value

I/p - df.isnull.any(axis=1)

O/p - 0 True
1 True
2 True
3 True
4 True
dtype: bool

Rows with all missing values

1/p - df.isnull().all(axis=1)

O/p - 0 False

1 False

2 False

3 False

4 False

dtype: bool

Number of rows with all missing values

1/p - df.isnull.all(axis=1).sum()

O/p - 0

MISSING VALUES TREATMENT IN COLUMNS :

1/p - round(100*(df.isnull().sum()/len(df.index)), 2)

O/p - First Name 20.0

Last Name 20.0

Age 20.0

Sex 20.0

PreTest Score 20.0

Post Test score 20.0

Location 100.0

dtype : float64

removing the location column

I/p - df.dropna(axis=1, how='all', inplace=True)

round(100 * (df.isnull().sum() / len(df.index)),
2)

O/p - First Name 20.0

last Name 20.0

Age 20.0

Sex 20.0

Pre Test score 20.0

Post Test score 20.0

dtype: float64

MISSING VALUES TREATMENT IN ROWS

I/p - df[df.isnull().sum(axis=1) >= 4]

O/p -

	First Name	last Name	Age	Sex	Pre Test Score	Post Test Score
1	Nan	Nan	Nan	Nan	90.0	Nan

retaining the rows having ≤ 4 NaN's

I/P - df = df[df.isnull().sum(axis=1) $\leq 4]$

look at the summary again

round(100*(df.isnull.sum()/len(df.index)),2)

O/P - FirstName

0.0

lastName

0.0

Age

0.0

Sex

0.0

PreTest Score

25.0

PostTest score

0.0

dtype: float64

I/P - df.head()

	First	last	Age	Sex	Pre Test Score	Post Test Score
	Name	Name				
0	abc	mno	12.0	m	90.0	65.0
1	ghi	pqr	12.0	f	Nan	65.0
2	jkl	stu	12.0	f	90.0	62.0
3	mno	vwx	12.0	m	89.0	63.0
4						

I/P - df['preTestScore'].describe()

O/P - count 3.0000

mean 89.6667

std 0.5773

min 89.0000

25% 89.5000

50% 90.0000

75% 90.0000

max 90.0000

Name : preTestScore, dtype : float64

Imputing pretestscore by mean values

I/P - df['preTestScore'].fillna(df['PreTestScore'].mean(),
 inplace=True)

round(100*(df.isnull().sum()/len(df.index)),2)

O/P - FirstName 0.0

LastName 0.0

Age 0.0

Sex 0.0

PreTestScore 0.0

PostTestScore 0.0

dtype : float64

I/p - df.head()

(135)

1309-10-2022

	First Name	Last Name	Age	Sex	Pre Test score	Post Test score
0	abc	mno	12.0	m	90.00	65.0
2	ghi	pqr	12.0	f	89.67	65.0
3	fgk	stc	12.0	f	90.00	62.0
4	mno	rwx	12.0	m	89.00	63.0

fraction of rows lost

I/p → 1 - len(df.index)/5

O/p - 0.19996