

Missing Values Approaches

Dealing with the missing values is most important step of Feature Engineering. We can do it in following ways-1. Check for Missing data 2. Drop NaN values 3. Mean/Median Imputation 4. Random Sample Imputation

```
In [1]: # Importing basic libraries
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: # Loading the dataset
df = pd.read_csv('Google_Cleaned.csv')
```

```
In [3]: # Getting sample of the data
df.sample(10)
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	Genres	Last Updated	Current Ver	Android Ver	Day	Month
7723	QRZ Assistant	COMMUNICATION	4.2	1044	3.5	100000	Free	0.00	Everyone	Communication	2018-05-30	2.0	4.4 and up	30	5
10466	Tassa & Finland	LIFESTYLE	3.6	346	7.5	50000	Free	0.00	Everyone	Lifestyle	2018-05-22	5.5	4.0 and up	22	5
10277	Fan App for Reading FC	SPORTS	3.6	9	13.0	500	Free	0.00	Teen	Sports	2018-08-06	000816	4.2 and up	6	8
9941	Tasker	TOOLS	4.6	43045	NaN	1000000	Paid	2.99	Everyone	Tools	2018-06-25	Varies with device	Varies with device	25	6
3620	Microsoft Edge	COMMUNICATION	4.3	27187	9.2	1000000	Free	0.00	Everyone	Communication	2018-07-28	42.0.0.2059	4.4 and up	28	7
1867	Mahalaam Dinandika 2018	PRODUCTIVITY	4.2	25427	9.2	1000000	Free	0.00	Everyone	Productivity	2018-01-05	18.0	4.0 and up	5	1
8990	DU GIF Maker: GIF Maker, Video to GIF & GIF Ed...	PHOTOGRAPHY	4.6	7695	4.3	500000	Free	0.00	Everyone	Photography	2018-04-10	1.2.2.2	4.3 and up	10	4
10627	PHOTOONE CREB Unleash of P	FINANCE	4.7	178	12.0	1000	Free	0.00	Everyone	Finance	2018-06-12	5.9.1.0	5.0 and up	12	6
8972	Flag of European Union UMP	PERSONALISATION	4.3	88	4.2	1000	Free	0.00	Everyone	Personalisation	2018-01-24	7.0	4.0.3 and up	24	1
2432	Logo Ninjab Wallpaper	MEDICAL	NaN	2	15.0	100	Free	0.00	Everyone	Medical	2018-06-02	1.2	3.0 and up	2	8

Checking for the Missing Values

```
In [4]: df.isna().sum()
```

App	0
Category	0
Rating	1474
Reviews	0
Size	1695
Installs	0
Type	1
Price	0
Content Rating	0
Genres	0
Last Updated	0
Current Ver	8
Android Ver	2
Day	0
Month	0
dtype: int64	

```
In [5]: df.isna().sum().sum()
```

3180

We have total 3180 nan values in the dataset.

```
In [6]: df.isna().sum().sort_values(ascending=False)
```

Size	1695	15.63636
Rating	1474	0.135978
Android Ver	2	0.000000
Type	1	0.000000
App	0	0.000000
Category	0	0.000000
Reviews	0	0.000000
Installs	0	0.000000
Price	0	0.000000
Content Rating	0	0.000000
Genres	0	0.000000
Last Updated	0	0.000000
Current Ver	8	0.000000
Day	0	0.000000
Month	0	0.000000
dtype: float64		

```
In [7]: df.shape[0]
```

28840

```
In [8]: null_df = pd.DataFrame({'Null Values': df.isna().sum().sort_values(ascending=False), 'Percentage Null Values': (df.isna().sum().sort_values(ascending=False)) / (df.shape[0]) * 100})
null_df
```

Null Values	Percentage Null Values
Size	15.63636
Rating	1474
Current Ver	8
Android Ver	2
Type	1
App	0
Category	0
Reviews	0
Installs	0
Price	0
Content Rating	0
Genres	0
Last Updated	0
Day	0
Month	0
dtype: float64	

```
In [9]: null_counts = df.isna().sum().sort_values(ascending=False)/len(df)
```

```
In [10]: null_counts
```

Size	0.156366
Rating	0.135978
Current Ver	0.000738
Android Ver	0.000085
Type	0.000002
App	0.000000
Category	0.000000
Reviews	0.000000
Installs	0.000000
Price	0.000000
Content Rating	0.000000
Genres	0.000000
Last Updated	0.000000
Day	0.000000
Month	0.000000
dtype: float64	

```
In [11]: len(null_counts)
```

15

```
In [12]: np.arange(len(null_counts))
```

array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14])

```
In [13]: # Plotting Null Values plots
plt.figure(figsize=(15,8))
plt.xticks(np.arange(len(null_counts))+0.5, null_counts.index, rotation='vertical')
plt.ylabel('Fraction of rows with missing data')
plt.bar(np.arange(len(null_counts)), null_counts)
plt.show()
```


1. Dropping Missing Values

In these method we are simply dropping the missing value datapoints.

```
In [14]: df_copy = df.copy()
```

```
In [15]: cols = [var for var in df_copy.columns if df_copy[var].isnull().mean()*100 > 0.1]
```

['Rating', 'Size', 'Type', 'Current Ver', 'Android Ver']

```
In [16]: df_copy[cols].isnull().sum()
```

1695

```
In [17]: df_copy[cols].isnull().mean()*100
```

15.63631365313653

```
In [18]: df_copy[cols].isnull().mean()*100
```

0.0

```
In [19]: if df_copy[cols].isnull().sum():
    print('yes')
else:
    print('no')
```

yes

```
In [20]: if df_copy[cols].isnull().sum()*100 > 0.1:
    print('yes')
else:
    print('no')
```

no

```
In [21]: df_copy.shape
```

(16840, 15)

```
In [22]: # We are dropping missing values
```

```
In [23]: drop_df = df_copy[cols].dropna()
```

Rating	Size	Type	Current Ver	Android Ver
0	4.1	19.000	Free	1.0.0
1	3.9	14.000	Free	2.0.0
2	4.7	0.700	Free	1.2.4
3	4.5	26.000	Free	Varies with device
4	4.3	2.800	Free	1.1
10832	4.8	0.619	Free	0.8
10833	4.0	2.600	Free	1.0.0
10835	4.5	53.000	Free	1.48
10836	5.0	3.600	Free	1.0
10839	4.5	1.000	Free	Varies with device

7723 rows x 5 columns

```
In [24]: df_copy.shape, drop_df.shape
```

((16840, 15), (7723, 5))

```
In [25]: 16840-7723
```

9117

```
In [26]: # We are going to drop this much data
```

```
In [27]: plt.figure(figsize=(15,7))
plt.subplots(2,2)
drop_df[cols].plot.density(colors='red')
df_copy[cols].plot.density(colors='green')
```



```
In [28]: drop_df[cols].plot.density(colors='red')
df_copy[cols].plot.density(colors='green')
```


insights

- As we can observe from the above plots 1. Drop NA technique is changing our distribution pattern. 2. So we reject Drop NA Technique

2. Mean/Median Imputation

In these method we can fill the nan value with the Mean or Median of the particular feature.

```
In [29]: df_copy_m_me = df_copy()
```

```
In [30]: df_copy_m_me[df_copy_m_me['Size'].isnull()]
```

App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	Genres	Last Updated	Current Ver	Android Ver	Day	Month
0	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
1	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
...
10825	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
10826	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
10827	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
10828	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
10829	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

10840 rows x 15 columns

```
In [31]: df_copy_m_me[df_copy_m_me['Size'].isnull()]
```

App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	Genres	Last Updated	Current Ver	Android Ver	Day	Month	
37	Floor Plan Creator	ART_AND_DESIGN	4.1	39639	NaN	1000000	Free	0.0	Everyone	Art & Design	2018-07-14	Varies with device	2.3.3 and up	14	7
42	Testgram - write on photos	ART_AND_DESIGN	4.4	295221	NaN	10000000	Free	0.0	Everyone	Art & Design	2018-07-30	Varies with device	Varies with device	30	7
58	Used Cars and Trucks for Sale	AUTO_AND_VEHICLES	4.6	17057	NaN	1000000	Free	0.0	Everyone	Auto & Vehicles	2018-07-30	Varies with device	Varies with device	30	7
67	Ulysee Speedometer	AUTO_AND_VEHICLES	4.3	40211	NaN	5000000	Free	0.0	Everyone	Auto & Vehicles	2018-07-30	Varies with device	Varies with device	30	7
68	REPUVE	AUTO_AND_VEHICLES	3.9	356	NaN	100000	Free	0.0	Everyone	Auto & Vehicles	2018-05-25	Varies with device	Varies with device	25	5
...	
10712	My Earthquake Alerts - US & Worldwide Earthquakes	WEATHER	4.4	3471	NaN	100000	Free	0.0	Everyone	Weather	2018-07-24	Varies with device	Varies with device	24	7
10724	Posta App	MAPS_AND_NAVIGATION	3.6	8	NaN	1000	Free	0.0	Everyone	Maps & Navigation	2017-09-27	Varies with device	4.4 and up	27	9
10764	Chat For Strangers - Video Chat	SOCIAL	3.4	622	NaN	10000	Free	0.0	Mature 17+	Social	2018-05-23	Varies with device	Varies with device	23	5
10825	Prim: get new friends on local chat rooms	SOCIAL	4.0	89486	NaN	5000000	Free	0.0	Mature 17+	Social	2018-03-23	Varies with device	Varies with device	23	3
10838	The SCP Foundation DB fr mfin	BOOKS_AND_REFERENCE	4.5	114	NaN	1000	Free	0.0	Mature 17+	Books & Reference	2015-01-19	Varies with device	Varies with device	19	1

1695 rows x 15 columns

```
In [32]: df_copy_m_me['Size'].mean()
```

21.50654062329852

```
In [33]: df_copy_m_me['mean_size'] = df_copy_m_me['Size'].fillna(df_copy_m_me['Size'].mean())
df_copy_m_me['median_size'] = df_copy_m_me['Size'].fillna(df_copy_m_me['Size'].median())
df_copy_m_me['mean_rating'] = df_copy_m_me['Rating'].fillna(df_copy_m_me['Rating'].mean())
df_copy_m_me['median_rating'] = df_copy_m_me['Rating'].fillna(df_copy_m_me['Rating'].median())
```

```
In [34]: print('Original Size Variance', df_copy_m_me['Size'].var())
```

Original Size Variance 518.5801557864865

```
In [35]: print('Size Variance after mean imputation', df_copy_m_me['mean_size'].var())
```

Size Variance after mean imputation 430.7357638630519

```
In [36]: print('Size Variance after median imputation', df_copy_m_me['median_size'].var())
```

Size Variance after median imputation 440.28217654805237

Here we can observe that the data has less variance in mean imputation in compare to median imputation. Therefore, mean will be good option because of less variance.

```
In [37]: print('Original Rating Variance', df_copy_m_me['Rating'].var())
print('Rating Variance After mean imputation', df_copy_m_me['mean_rating'].var())
print('Rating Variance After median imputation', df_copy_m_me['median_rating'].var())
```

Original Rating Variance 0.28545847227541496
Rating Variance After mean imputation 0.22935175503821595322
Rating Variance After median imputation 0.22872842263353322

```
In [38]: df_copy_m_me['Size'].plot.density()
```



```
In [39]: df_copy_m_me['mean_size'].plot.density()
```



```
In [40]: df_copy_m_me['median_size'].plot.density()
```


insights

-As we can observe from the above plots 1. Mean and Median imputation technique is changing our distribution pattern which is going to create problems on Assumptions of the original data. 2. So we will reject Mean and Median Imputation Technique also.

3. Random Sample Imputation

Here we are going to fill the nan values with that much number of random sample of datapoints.

```
In [41]: df_random = df_copy()
```

```
In [42]: df_random['Size'].dropna().sample(20)
```

6527 0.2
18660 29.0
3565 15.0
1884 18.0
6108 4.0
18259 33.0
3279 1.9
38322 21.0
9551 60.0
5838 29.0
2385 37.0
18282 4.6
6518 6.2
6215 6.6
282 15.0
1657 46.0
38391 61.0
4785 38.0
16781 11.0
4288 27.0
Name: Size, dtype: float64

```
In [43]: # Getting index of the nan value row
df_random[df_random['Size'].isnull()].index
```

Int64Index([37, 42, 52, 67, 68, 73, 85, 88, 89, 92, 10646, 10678, 10688, 10706, 10711, 10712, 10724, 10764, 10825, 10838], dtype='int64', length=1695)

Size feature has 1695 nan values thus we need to take that much random sample datapoints

```
In [44]: df_random['Size'].dropna().sample(1695)
```

18411 29.000
7119 15.000
9252 4.000
10835 53.000
6708 9.778
3886 4.100
6130 3.400
2178 29.000
6508 4.200
2297 18.000
Name: Size, length: 1695, dtype: float64

```
In [45]: df_random.isnull().sum()
```

App	0
Category	0
Rating	1474
Reviews	0
Size	1695
Installs	0
Type	1
Price	0
Content Rating	0
Genres	0
Last Updated	0
Current Ver	8
Android Ver	2
Day	0
Month	0
dtype: int64	

```
In [46]: df_random.sample Imputation(feature):
random_sample = df_random[feature].dropna().sample(df_random[feature].isnull().sum())
random_sample.index = df_random[feature].isnull().index
df_random.loc[df_random[feature].isnull(), feature] = random_sample
```

```
In [47]: for col in df_random:
    Random_Sample Imputation(col)
```

```
In [48]: print('Original Size Variance', df['Size'].var())
print('Size Variance After Random Imputation', df_random['Size'].var())
```

Original Size Variance 518.5801557864865
Size Variance After Random Imputation 515.721818038196

```
In [49]: print('Original Size Variance', df['Rating'].var())
print('Size Variance After Rating Imputation', df_random['Rating'].var())
```