

```

# Import pandas and libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import sys

# Loading the dataset
df = pd.read_csv('googleplaystore.csv')

# Getting top 5 rows of the dataset
df.head()

# Getting last 5 rows of the dataset
df.tail()

# Getting datatypes of the features
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10841 entries, 0 to 10840
Data columns (total 13 columns):
 #   Column          Non-Null Count  Dtype
---  --
 0   App             10841 non-null   object
 1   Category        10841 non-null   object
 2   Rating          9987 non-null    float64
 3   Reviews         10841 non-null   object
 4   Size            9145 non-null    float64
 5   Installs        10841 non-null   object
 6   Type            10841 non-null   object
 7   Price           10841 non-null   object
 8   Content Rating  10841 non-null   object
 9   Genres          10841 non-null   object
10   Last Updated   10841 non-null   object
11   Current Ver    10832 non-null   object
12   Android Ver    10838 non-null   object
dtypes: float64(1), object(12)
memory usage: 1.1+ MB

# Getting number of rows and number of columns
df.describe()

(10841, 13)

# Checking if our dataset is having duplicates values or not
df.duplicated()

# Getting top 10 apps
df.sort_values('Reviews', ascending=False).head(10)

# Getting basic statistical terms
df.describe(include='all')

count      App      Category      Rating      Reviews      Size      Installs      Type      Price      Content Rating      Genres      Last Updated      Current Ver      Android Ver
unique    9650      34      NaN      6002      462      22      3      93      6      120      178      2832      33
top  ROLIOX      FAMILY      NaN      0      Values with device      100000+      Free      0      Everyone      Tools      August 13, 2018      Values with device      4.1 and up
freq      App      9      1972      NaN      595      1695      1579      10039      10040      8714      842      325      1459      2451
mean      App      NaN      4.133338      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
std      App      NaN      0.574311      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
min      App      NaN      1.000000      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
max      App      NaN      10.000000      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
50%      App      NaN      4.300000      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
75%      App      NaN      4.300000      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
max      App      NaN      10.000000      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN

# Getting reviews
df['Reviews'].head()

0      155
1      967
2      87519
3      215644
4      967
Name: Reviews, dtype: object

df['Reviews'].shape
(10841, 1)

df['Reviews'].dtype
dtype('O')

df['Reviews']
0      155
1      967
2      87519
3      215644
4      967
...
10836      38
10837      3
10838      3
10839      114
10840      963937
Name: Reviews, Length: 10841, dtype: object

df['Reviews'].str.isnumeric().sum()
10840

df['Reviews'].str.isnumeric()
0      False
1      False
2      False
3      False
4      False
...
10836      False
10837      False
10838      False
10839      False
10840      False
Name: Reviews, Length: 10841, dtype: bool

# (~) denotes the complement
df[~df['Reviews'].str.isnumeric()]

# Copying the original dataset
df_copy = df.copy()

# Dropping the 10472 index of the review columns
df_copy = df_copy.drop(df_copy.index[10472])

df_copy.shape
(10840, 13)

df_copy['Reviews'].dtype
dtype('O')

# Changing the datatype into int
df_copy['Reviews'] = df_copy['Reviews'].astype('int')

# It is good practice to copy our data
df_copy['Reviews'].dtype
dtype('int64')

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10840 entries, 0 to 10839
Data columns (total 13 columns):
 #   Column          Non-Null Count  Dtype
---  --
 0   App             10840 non-null   object
 1   Category        10840 non-null   object
 2   Rating          9986 non-null    float64
 3   Reviews         10840 non-null   int64
 4   Size            9145 non-null    float64
 5   Installs        10840 non-null   object
 6   Type            10840 non-null   object
 7   Price           10840 non-null   object
 8   Content Rating  10840 non-null   object
 9   Genres          10840 non-null   object
10   Last Updated   10840 non-null   object
11   Current Ver    10832 non-null   object
12   Android Ver    10838 non-null   object
dtypes: float64(1), int64(1), object(11)
memory usage: 1.1+ MB

# Cloning the dataset
df_copy['Size'].unique()

array(['15M', '14M', '8.7M', '25M', '12.8M', '5.6M', '20M', '13M', '1.3M',
       '28M', '12M', '20M', '12M', '37M', '2.7M', '5.5M', '17M', '39M',
       '3.2M', '4.2M', '7.4M', '13M', '6.8M', '1.9M', '10M', '1.5M', '1.8M',
       '5.3M', '13M', '24M', 'Values with device', '9.4M', '15M', '10M',
       '1.2M', '20M', '8.8M', '9.9M', '20M', '10M', '3M', '5M', '20K',
       '1.8M', '5.7M', '8.6M', '2.4M', '23M', '2.5M', '15M', '3.4M',
       '1.9M', '3.5M', '2.9M', '3.9M', '32M', '1.4M', '18M', '1.1M',
       '1.2M', '3.7M', '22M', '7.4M', '6.4M', '3.2M', '8.2M', '9.9M',
       '4.9M', '10.3M', '5.9M', '1.5M', '13M', '73M', '6.8M', '3.9M',
       '1.4M', '12.3M', '7.2M', '12.1M', '42M', '7.3M', '3.4M', '55M',
       '1.6M', '1.4M', '1.5M', '1.5M', '53M', '12M', '48M', '6.5M', '44M',
       '6.3M', '4.3M', '4.7M', '1.3M', '49M', '7.8M', '8.8M', '1.6M',
       '5.3M', '61M', '68M', '73M', '8.4M', '18M', '44M', '65M',
       '6.2M', '18K', '63M', '1.4M', '5.8M', '5.8M', '3.8M', '1.9M',
       '49M', '63M', '49M', '77M', '1.4M', '8.8M', '70M', '6.9M', '1.3M',
       '18.0M', '8.2M', '36M', '64M', '37M', '2.0M', '1.9M', '1.8M',
       '5.3M', '47M', '155K', '52K', '76M', '7.0M', '15M', '9.7M', '73M',
       '72M', '43M', '7.7M', '6.3M', '33K', '14M', '92M', '60M', '70M',
       '188M', '55M', '59M', '68M', '64M', '67M', '66M', '84M', '22K',
       '99M', '65M', '94M', '8.5M', '14M', '70K', '11M', '69M', '12.7M',
       '74M', '62M', '60M', '75M', '98M', '85M', '82M', '56M', '87M',
       '73M', '69M', '92M', '12M', '10M', '10M', '10M', '10M', '10M',
       '89M', '37K', '26K', '37K', '1.3M', '975K', '898K', '4.1M',
       '89M', '69K', '54K', '52K', '52K', '77K', '85K', '72K',
       '73K', '72K', '31K', '58', '24K', '15K', '87K', '93',
       '93K', '85K', '25K', '93K', '54K', '73K', '74K', '10K',
       '28', '14K', '33K', '37K', '22K', '73M', '75K', '93K',
       '28K', '17K', '74K', '14K', '137K', '78K', '92K', '60K', '81K',
       '51', '93K', '10K', '37K', '10K', '10K', '10K', '10K',
       '28K', '65K', '74K', '93K', '87K', '12K', '32K', '1.0M',
       '91K', '17K', '23K', '14M', '14M', '20K', '95K', '44K', '77K',
       '21K', '69K', '38K', '78K', '30K', '90K', '47K', '17K',
       '30K', '38K', '45K', '42K', '78K', '82K', '42K', '84K',
       '67K', '42K', '45K', '47K', '83K', '83K', '75K', '72K', '43K',
       '19K', '19K', '26K', '46K', '72K', '49K', '83K', '83K', '44K',
       '56K', '87K', '61K', '24K', '56M', '77K', '68K', '59K',
       '31K', '18K', '84K', '64K', '15K', '37K', '43K', '59K',
       '73K', '10K', '12K', '12K', '10K', '10K', '10K',
       '69K', '15K', '93K', '93K', '2K', '54K', '35K', '7K',
       '52', '29K', '92K', '14K', '53K', '2K', '82K', '89K',
       '74K', '11K', '13K', '29K', '83K', '49K', '17K', '72K',
       '89M', '12K', '41K', '60K', '78M', '78M', '60M', '60M',
       '64K', '96K', '97K', '51K', '87K', '78K', '78K', '78K',
       '34K', '92K', '25K', '18K', '45K', '24K', '62K', '78K',
       '28K', '77K', '78K', '63K', '18K', '94M', '39K', '75K',
       '91K', '93K', '68K', '50K', '54K', '56K', '84K', '90K',
       '68K', '91K', '27K', '19K', '10K', '52K', '92K', '10K',
       '98K', '74K', '28K', '74K', '51K', '75K', '82K', '45K',
       '68M', '84K', '87K', '82K', '18K', '10K', '10K',
       '17K', '14K', '88K', '44K', '14K', '10K', '37K',
       '44K', '67K', '67K', '55K', '85K', '182K', '56K', '61K',
       dtype=object)

# Replacing 'M' with '000'
df_copy['Size'] = df_copy['Size'].str.replace('M', '000')

df_copy['Size']
0      15000
1     
```