```
In [1]: #Importing Required libraries
        import pandas as pd
        import numpy as np
        import math
        import sys
        import findspark
        findspark.init()
In [2]: from pyspark import StorageLevel
        from pyspark.sql import SparkSession
        from pyspark.sql.functions import isnan,lit, col, explode, initcap, regexp_replace
        from pyspark.sql.types import StructType, StructField, StringType, IntegerType, F
        from pyspark.sql.dataframe import DataFrame
In [3]: spark = SparkSession \
                .builder \
                .appName("PlaysoreAppsPySpark") \
                .getOrCreate()
```

```
In [4]: # Reading the raw dataset from Cloud Storage bucketin to the python dataframe psd
        psapps_df = spark.read.csv("gs://ca675-assignment2-bucket/Raw_Data/Google-Playsto
        # Checking the first 10 rows of the dataframe
        psapps df.show(10)
```

```
-----
 -----+
                                        Category | Rating | Rating Cou
          App Name
                            App Id
nt|Installs|Minimum Installs|Maximum Installs|Free|Price|Currency|Size|Minimu
                                            Developer Email
m Android
              Developer Id | Developer Website |
Released Last Updated Content Rating
                                Privacy Policy Ad Supported In App
                        Scraped Time
Purchases | Editors Choice |
+-----
-----
          Gakondo| com.ishakwe.gakondo|
                                       Adventure|
0|
                                 15|true| 0.0|
                                                USD| 10M|
     10+|
                   10
1 and up|Jean Confident Ir...|https://beniyizib...|jean21101999@gmai...|Feb 2
6, 2020 | Feb 26, 2020 | Everyone | https://beniyizib... |
                                                  false
           false 2021-06-15 20:19:35
false
| Ampere Battery Info|com.webserveis.ba...|
                                          Tools
                                                 4.4
64 5,000+
                  5000
                                7662|true| 0.0|
                                                USD 2.9M
5.0 and up
                Webserveis|https://webservei...|webserveis@gmail.com|May
21, 2020 | May 06, 2021 |
                     Everyone|https://dev4phone...|
false|
           false 2021-06-15 20:19:35
           Vibook|com.doantiepvien.crm|
                                     Productivity|
                                                 0.0
01
     50+|
                   50 l
                                 58|true| 0.0|
                                                USD 3.7M
              Cabin Crew
                                   null | vnacrewit@gmail.com | Aug
3 and up
                     Everyone|https://www.vietn...|
9, 2019 Aug 19, 2019
                                                  false|
           false 2021-06-15 20:19:35
Smart City Trichy...|cst.stJoseph.ug17...| Communication|
                                                 5.0
5
     10+|
                   10
                                 19|true| 0.0|
                                                USD | 1.8M |
3 and up | Climate Smart Tech2|http://www.climat...|climatesmarttech2...|Sep 1
0, 2018|Oct 13, 2018|
                     Everyone
                                        null
                                                   truel
           false 2021-06-15 20:19:35
false
          GROW.me|com.horodyski.grower|
                                          Tools
                                                 0.0
01
    100+|
                  100
                                478|true| 0.0|
                                                USD 6.2M
1 and up|Rafal Milek-Horod...|http://www.horody...|rmilekhorodyski@g...|Feb 2
1, 2020 | Nov 12, 2018 |
                    Everyone|http://www.horody...|
                                                  false|
false
           false 2021-06-15 20:19:35
           IMOCCI|
                         com.imocci
                                         Social
                                                 0.0
0|
     50+|
                   50
                                 89|true| 0.0|
                                                USD| 46M|
0 and up
              Imocci GmbH|http://www.imocci...|
                                           info@imocci.com|Dec 2
4, 2018 | Dec 20, 2019 |
                        Teen|https://www.imocc...|
                                                  false
          false 2021-06-15 20:19:35
|unlimited 4G data...|getfreedata.super...|Libraries & Demo|
                                                 4.5
                                2567|true| 0.0|
                                                USD | 2.5M |
12 1,000+
                  1000
```

```
4.1 and up|android developer779|
                                        null|aitomgharfatimezz...|Sep
23, 2019|Sep 27, 2019|
                       Everyone|https://sites.goo...|
                                                       truel
            false 2021-06-15 20:19:35
The Everyday Cale... com.mozaix.simone...
                                         Lifestvle
                                                    2.0
                                   702|true| 0.0|
                                                   USD| 16M|
39
     500+
                    500
                                        null|elementuser03@gma...|Jun
5.0 and up
                 Mozaix LLC
21, 2019 Jun 21, 2019
                       Everyone|https://www.freep...|
                                                      false
            false 2021-06-15 20:19:35
false
          WhatsOpen
                    com.whatsopen.app
                                      Communication|
                                                    0.01
                                   18|true| 0.0|
                                                  USD | 1.3M |
01
     10+|
                    10|
                                                              4.
4 and up|Yilver Molina Hur...|http://yilvermoli...|yilver.mh1996@gma...|
null|Dec 07, 2018|
                       Teen|http://elcafedela...|
                                                  false
            false 2021-06-15 20:19:35
|Neon 3d Iron Tech...|com.ikeyboard.the...| Personalization|
                                                    4.7
                                                               8
20 | 50,000+|
                   50000
                                 62433|true| 0.0|
                                                   USD | 3.5M |
4.1 and up|Free 2021 Themes ...|https://trendytem...|trendyteme.888@gm...|Sep
22, 2019 Oct 07, 2020
                   Everyone|http://bit.ly/Emo...|
false
            false 2021-06-15 20:19:35
-----+
only showing top 10 rows
```

In [5]: #Checking the structure of the numeric feature psapps\_df.describe()

DataFrame[summary: string, App Name: string, App Id: string, Category: string, Rating: string, Rating Count: string, Installs: string, Minimum Installs: strin g, Maximum Installs: string, Price: string, Currency: string, Size: string, Min imum Android: string, Developer Id: string, Developer Website: string, Develope r Email: string, Released: string, Last Updated: string, Content Rating: strin g, Privacy Policy: string, Scraped Time: string]

In [6]:	<pre>#Checking the total count null values in eah feature of the dataframe psapps_df psapps_null_df=psapps_df.select([count(when(col(c).isNull(),c)) for c in psapps_ psapps_null_df.show()</pre>
	[Stage 6:=====> (5 + 1) / 6]
	+
	+
	+
	+
	count(CASE WHEN (App Name IS NULL) THEN App Name END) count(CASE WHEN (App I d IS NULL) THEN App Id END) count(CASE WHEN (Category IS NULL) THEN Category END) count(CASE WHEN (Rating IS NULL) THEN Rating END) count(CASE WHEN (Rating Count IS NULL) THEN Rating Count END) count(CASE WHEN (Installs IS NULL) THEN Installs END) count(CASE WHEN (Minimum Installs IS NULL) THEN Minimum Installs END) count(CASE WHEN (Maximum Installs IS NULL) THEN Maximum Installs END) count(CASE WHEN (Free IS NULL) THEN Free END) count(CASE WHEN (Price IS NULL) THEN Price END) count(CASE WHEN (Currency IS NULL) THEN Currency END) count(CASE WHEN (Size IS NULL) THEN Size END) count(CASE WHEN (Minimum Android I S NULL) THEN Minimum Android END) count(CASE WHEN (Developer Id IS NULL) THEN Developer Website END) count(CASE WHEN (Developer Website IS NULL) THEN Developer Website END) count(CASE WHEN (Developer Email IS NULL) THEN Developer Email E ND) count(CASE WHEN (Released IS NULL) THEN Released END) count(CASE WHEN (La st Updated IS NULL) THEN Last Updated END) count(CASE WHEN (Content Rating IS NULL) THEN Content Rating END) count(CASE WHEN (Privacy Policy IS NULL) THEN Privacy Policy END) count(CASE WHEN (Ad Supported IS NULL) THEN Ad Supported END) count(CASE WHEN (In App Purchases IS NULL) THEN In App Purchases END) count(CASE WHEN (Editors Choice IS NULL) THEN Editors Choice END) count(CASE WHEN (Scraped Time IS NULL) THEN Scraped Time END)
	+
	++
	+

```
0|
0
                                                    0
22883
                                                            22883
107
                                                                    107
0
                                             0|
0
                                                   135
196
                                                                 6530
33|
                                                                  760835
                                                  71053
31|
                                                                0|
0
420953
                                                                 0
0
0
```

- In [7]: #Cleaning-Step1 #Removing the unwanted features abd storing the resultant in new dataframe 'psape psapps\_selected\_df=psapps\_df.drop("App Id","Minimum Installs","Currency","Minimum # Checking top 10 values after removal of unwanted features psapps selected df.head(10)
- Out[7]: [Row(App Name='Gakondo', Category='Adventure', Rating=0.0, Rating Count=0, Inst alls='10+', Maximum Installs=15, Free=True, Price=0.0, Size='10M', Released='Fe b 26, 2020', Content Rating='Everyone', Ad Supported=False, In App Purchases=Fa lse),

Row(App Name='Ampere Battery Info', Category='Tools', Rating=4.4, Rating Count =64, Installs='5,000+', Maximum Installs=7662, Free=True, Price=0.0, Size='2.9 M', Released='May 21, 2020', Content Rating='Everyone', Ad Supported=True, In A pp Purchases=False),

Row(App Name='Vibook', Category='Productivity', Rating=0.0, Rating Count=0, In stalls='50+', Maximum Installs=58, Free=True, Price=0.0, Size='3.7M', Released ='Aug 9, 2019', Content Rating='Everyone', Ad Supported=False, In App Purchases =False),

Row(App Name='Smart City Trichy Public Service Vehicles 17UCS548', Category='C ommunication', Rating=5.0, Rating Count=5, Installs='10+', Maximum Installs=19, Free=True, Price=0.0, Size='1.8M', Released='Sep 10, 2018', Content Rating='Eve ryone', Ad Supported=True, In App Purchases=False),

Row(App Name='GROW.me', Category='Tools', Rating=0.0, Rating Count=0, Installs ='100+', Maximum Installs=478, Free=True, Price=0.0, Size='6.2M', Released='Feb 21, 2020', Content Rating='Everyone', Ad Supported=False, In App Purchases=Fals e),

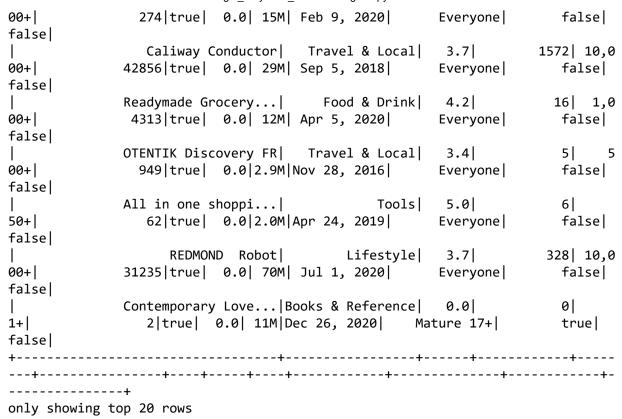
Row(App Name='IMOCCI', Category='Social', Rating=0.0, Rating Count=0, Installs ='50+', Maximum Installs=89, Free=True, Price=0.0, Size='46M', Released='Dec 2 4, 2018', Content Rating='Teen', Ad Supported=False, In App Purchases=True), Row(App Name='unlimited 4G data prank free app', Category='Libraries & Demo', Rating=4.5, Rating Count=12, Installs='1,000+', Maximum Installs=2567, Free=Tru e, Price=0.0, Size='2.5M', Released='Sep 23, 2019', Content Rating='Everyone',

Row(App Name='The Everyday Calendar', Category='Lifestyle', Rating=2.0, Rating Count=39, Installs='500+', Maximum Installs=702, Free=True, Price=0.0, Size='16 M', Released='Jun 21, 2019', Content Rating='Everyone', Ad Supported=False, In App Purchases=False),

Ad Supported=True, In App Purchases=False),

Row(App Name='WhatsOpen', Category='Communication', Rating=0.0, Rating Count= 0, Installs='10+', Maximum Installs=18, Free=True, Price=0.0, Size='1.3M', Rele ased=None, Content Rating='Teen', Ad Supported=False, In App Purchases=False), Row(App Name='Neon 3d Iron Tech Keyboard Theme', Category='Personalization', R ating=4.7, Rating Count=820, Installs='50,000+', Maximum Installs=62433, Free=T rue, Price=0.0, Size='3.5M', Released='Sep 22, 2019', Content Rating='Everyon e', Ad Supported=True, In App Purchases=False)]

```
In [8]: #Cleaning-Step2
       # Replacing the space in column names with '_' and storing the new fieldnames in df = [col(column).alias(column.replace('', '_')) for column in psapps_selected_c
       psapps selected df.select(*df).show()
       df1=psapps selected df.select(*df)
       +-----
       App Name
                                                Category | Rating | Rating | Count | Insta
       1ls|Maximum Installs|Free|Price|Size|
                                            Released | Content Rating | Ad Supported | I
       n App Purchases
       +-----
       -----+
                                Gakondo|
                                                           0.0
                                               Adventure
       10+|
                       15|true|
                                0.0 | 10M | Feb 26, 2020 |
                                                          Everyone|
                                                                        falsel
       falsel
                      Ampere Battery Info
                                                   Tools
                                                                       64 5,0
                                                           4.4
       00+1
                      7662|true|
                                0.0|2.9M|May 21, 2020|
                                                          Everyone|
                                                                         true
       false|
                                            Productivity|
                                                           0.0
                                 Vibook|
                                                                        0|
       50+|
                       58|true|
                                0.0|3.7M| Aug 9, 2019|
                                                          Everyone|
                                                                        falsel
       false|
                     Smart City Trichy...
                                           Communication|
                                                           5.0
                                                                        5|
                       19|true| 0.0|1.8M|Sep 10, 2018|
       10+|
                                                          Everyone|
                                                                         true
       false|
                                GROW.me|
                                                           0.0
                                                  Tools
                                                                        0|
                                                                              1
                                0.0|6.2M|Feb 21, 2020|
       00+l
                       478|true|
                                                          Everyone|
                                                                        false
       falsel
                                 IMOCCI
                                                  Social
                                                           0.0
                                                                        0|
                                0.0 | 46M | Dec 24, 2018 |
       50+|
                       89|true|
                                                             Teen
                                                                        false
       true
                     unlimited 4G data... | Libraries & Demo |
                                                           4.5
                                                                       12 1,0
       00+1
                      2567 true | 0.0 | 2.5 M | Sep 23, 2019 |
                                                          Everyone|
                                                                         truel
       false
                     The Everyday Cale...
                                               Lifestyle
                                                           2.0
                                                                       39|
                       702|true| 0.0| 16M|Jun 21, 2019|
                                                          Everyone|
                                                                        falsel
       00+1
       false|
                               WhatsOpen|
                                           Communication|
                                                           0.01
                                                                        01
                       18|true| 0.0|1.3M|
       10+|
                                                null
                                                             Teen
                                                                        false
       falsel
                     Neon 3d Iron Tech... | Personalization |
                                                                      820 | 50,0
                                                           4.7
       00+|
                     62433 | true | 0.0 | 3.5 M | Sep 22, 2019 |
                                                          Everyone|
                                                                         true
       false|
                         Dodge The Cars!
                                                  Racing|
                                                           4.9
                                                                       55
                       329|true| 0.0| 51M|Jul 30, 2020|
       00+1
                                                          Everyone|
                                                                        false
       falsel
                                Parents | Maps & Navigation |
                                                           0.01
                                                                        0|
       00+|
                       330|true|
                                0.0|2.7M|Jan 10, 2018|
                                                          Everyone|
                                                                        false
       falsel
       |桃園機場捷運時刻表 - 捷運轉乘路...|
                                        Travel & Local
                                                        3.9
                                                                    118 | 10,000+
                  37763|true| 0.0|7.6M| Apr 3, 2018|
                                                       Everyone|
                                                                      true|
       false|
                            be.MOBILISED | Maps & Navigation |
                                                           0.0
                                                                        0 l
                                                                              1
```



## In [9]: # Checking data structure of the dataframe df1 df1.printSchema()

|-- App\_Name: string (nullable = true)

```
|-- Category: string (nullable = true)
|-- Rating: double (nullable = true)
|-- Rating_Count: integer (nullable = true)
|-- Installs: string (nullable = true)
|-- Maximum Installs: long (nullable = true)
```

root

- |-- Free: boolean (nullable = true) |-- Price: double (nullable = true)
- |-- Size: string (nullable = true)
- |-- Released: string (nullable = true)
- |-- Content\_Rating: string (nullable = true)
- |-- Ad Supported: boolean (nullable = true)
- |-- In App Purchases: boolean (nullable = true)

```
In [10]: #Cleaning-Step3
         # Changing the datatype of Rating column to Double, Installs to Integer, Price to
         df1 = df1 \setminus
              .withColumn("Rating", col("Rating").cast(DoubleType())) \
              .withColumn("Rating_Count", col("Rating_Count").cast(IntegerType())) \
              .withColumn("Installs", regexp_replace(col("Installs"), "[^0-9]", "")) \
              .withColumn("Installs", col("Installs").cast(IntegerType())) \
              .withColumn("Price", regexp_replace(col("Price"), "[$]", "")) \
              .withColumn("Price", col("Price").cast(DoubleType())) \
              .withColumn("Released", to_date('Released', 'MMM d, yyyy'))
```

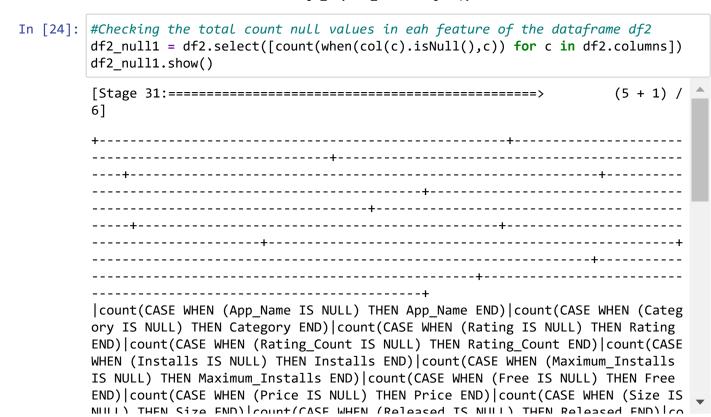
## In [11]: #Checking data structure of the dataframe df1 df1.printSchema()

```
root
```

```
|-- App_Name: string (nullable = true)
|-- Category: string (nullable = true)
|-- Rating: double (nullable = true)
|-- Rating Count: integer (nullable = true)
-- Installs: integer (nullable = true)
|-- Maximum Installs: long (nullable = true)
|-- Free: boolean (nullable = true)
|-- Price: double (nullable = true)
-- Size: string (nullable = true)
|-- Released: date (nullable = true)
|-- Content_Rating: string (nullable = true)
|-- Ad Supported: boolean (nullable = true)
|-- In App Purchases: boolean (nullable = true)
```

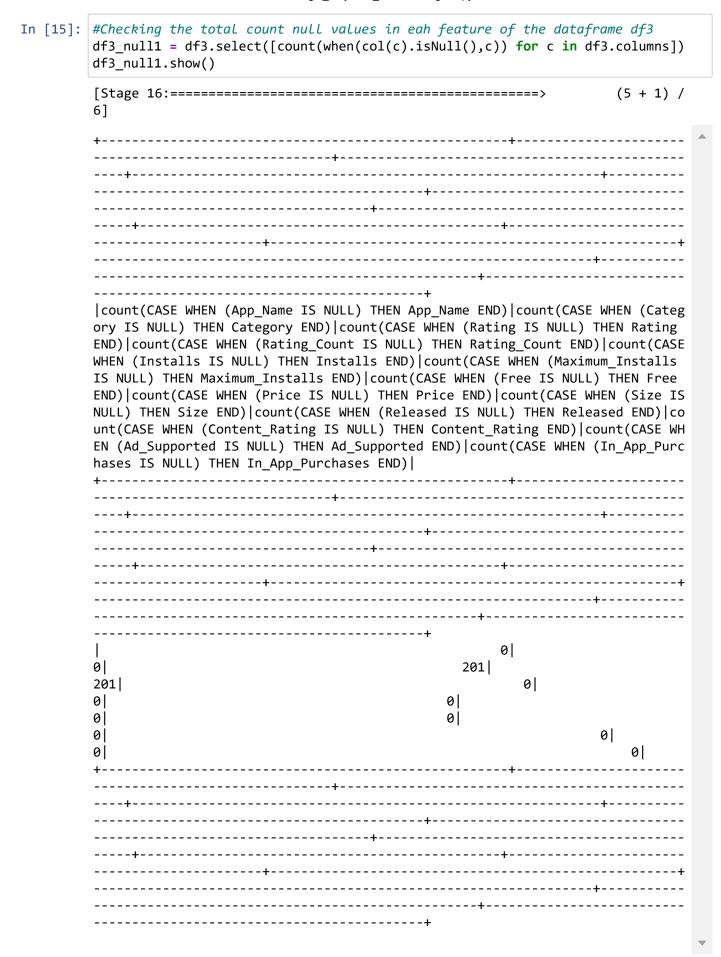
```
In [23]: #Cleaning-Step4
         # Dropping rows having null values in Release column and storing the result in do
         df2=df1.dropna(subset=["Released"])
         df2.show(10)
```

+			
++++	ng Rating_Count Installs Maximum_Ins ting Ad_Supported In_App_Purchases		
+++			
Gakondo  Adventure  0	.0  0  10		
15 true  0.0  10M 2020-02-26  Everyone	e  false  false		
Ampere Battery Info  Tools  4	.4  64  5000		
7662 true  0.0 2.9M 2020-05-21  Everyo			
Vibook  Productivity  0	• • • • • • • • • • • • • • • • • • • •		
58 true  0.0 3.7M 2019-08-09  Everyone			
Smart City Trichy  Communication  5			
19 true  0.0 1.8M 2018-09-10  Everyone			
GROW.me  Tools  0			
478 true  0.0 6.2M 2020-02-21  Everyor			
IMOCCI  Social 0.	.0  0  50		
89 true  0.0  46M 2018-12-24  Teer			
unlimited 4G data Libraries & Demo  4.	·		
2567 true  0.0 2.5M 2019-09-23  Everyo	·		
The Everyday Cale  Lifestyle  2.			
702 true  0.0  16M 2019-06-21  Everyor			
Neon 3d Iron Tech   Personalization   4.	·		
62433 true  0.0 3.5M 2019-09-22  Every Dodge The Cars!  Racing  4.			
329 true  0.0  51M 2020-07-30  Everyone  false  false			
+++			
only showing top 10 rows			



```
In [14]: #Cleaning-Step5
          # Dropping rows having null values in Installs column and storing the result in \mathfrak c
          df3=df2.na.drop(subset=["Installs"])
          df3.show(10)
```

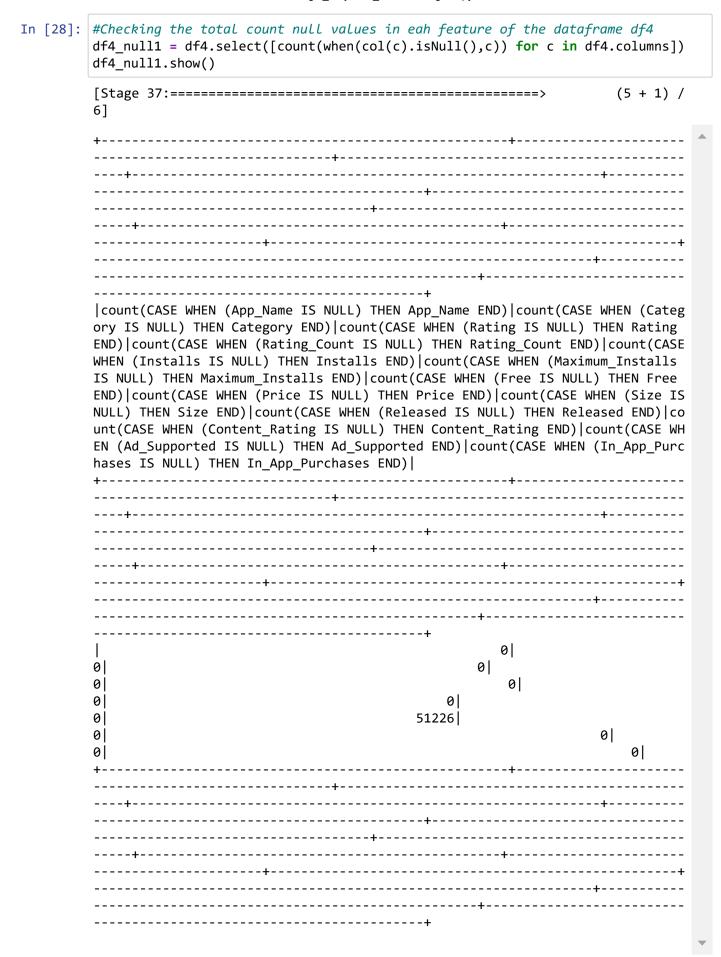
```
+-----
 App Name|
                    Category | Rating | Rating | Count | Installs | Maximum | Ins
talls|Free|Price|Size| Released|Content Rating|Ad Supported|In App Purchases|
+-----
Gakondo
                   Adventure
                             0.0
15|true| 0.0| 10M|2020-02-26|
                         Everyone|
                                    falsel
                                                false
| Ampere Battery Info|
                      Tools 4.4
                                      64
                                           5000
7662|true| 0.0|2.9M|2020-05-21|
                                      true|
                         Everyone
                                                 false|
          Vibook
                 Productivity|
                             0.0
                                       0
                                             50
58|true| 0.0|3.7M|2019-08-09|
                         Everyone|
                                    false
                                                false|
|Smart City Trichy...| Communication|
                             5.0
                                             10
                                       5|
19|true| 0.0|1.8M|2018-09-10|
                         Everyone
                                                false
                                    truel
                             0.0
         GROW.me
                      Tools
                                       01
                                            100
478|true| 0.0|6.2M|2020-02-21|
                         Everyone|
                                                false|
                                    false
                             0.0
          IMOCCI
                      Social|
                                             50
                                       01
89|true| 0.0| 46M|2018-12-24|
                                    false
                                                true|
                           Teen
|unlimited 4G data...|Libraries & Demo|
                             4.5
                                      12
                                           1000
2567|true| 0.0|2.5M|2019-09-23|
                         Evervonel
                                      truel
                                                 falsel
|The Everyday Cale...|
                                      39|
                   Lifestyle|
                             2.0
                                            500
702|true| 0.0| 16M|2019-06-21|
                                                false|
                         Everyone|
                                    false
|Neon 3d Iron Tech...| Personalization|
                            4.7
                                      820
                                           50000
62433|true| 0.0|3.5M|2019-09-22|
                          Everyone
                                      true
                                                  false
    Dodge The Cars!
                      Racing|
                             4.9
                                      55
                                            100
329|true| 0.0| 51M|2020-07-30|
                         Everyone
                                    false|
                                                false|
+-----
only showing top 10 rows
```



```
#Cleaning-Step6
In [16]:
         # Modifying the values of Size Column in dataframe df3
         df3 = df3 \setminus
             .withColumn("Size", regexp replace(col("Size"), 'k', 'e+3')) \
             .withColumn("Size", regexp_replace(col("Size"), 'M', 'e+6')) \
             .withColumn("Size", regexp_replace(col("Size"), 'G', 'e+9')) \
             .withColumn("Size", regexp_replace(col("Size"), "Varies with Device", 'nan'))
             .withColumn("Size", regexp_replace(col("Size"), "[+]", "")) \
             .withColumn("Size", regexp replace(col("Size"), "[^0-9]", "")) \
             .withColumn("Size", col("Size").cast(IntegerType()))
In [17]: #Checking the total count nan values in Size feature of the dataframe df3
         df3.select([count(when(isnan('Size'),True))]).show()
         Stage 19:=========>>
                                                                             (5 + 1) /
         6]
         |count(CASE WHEN isnan(Size) THEN true END)|
In [18]: # Checking the data structure of dataframe df3
         df3.printSchema()
         root
          |-- App_Name: string (nullable = true)
          |-- Category: string (nullable = true)
          |-- Rating: double (nullable = true)
          |-- Rating Count: integer (nullable = true)
          |-- Installs: integer (nullable = true)
           |-- Maximum_Installs: long (nullable = true)
          |-- Free: boolean (nullable = true)
          |-- Price: double (nullable = true)
          |-- Size: integer (nullable = true)
          |-- Released: date (nullable = true)
           |-- Content Rating: string (nullable = true)
          |-- Ad_Supported: boolean (nullable = true)
           |-- In App Purchases: boolean (nullable = true)
```

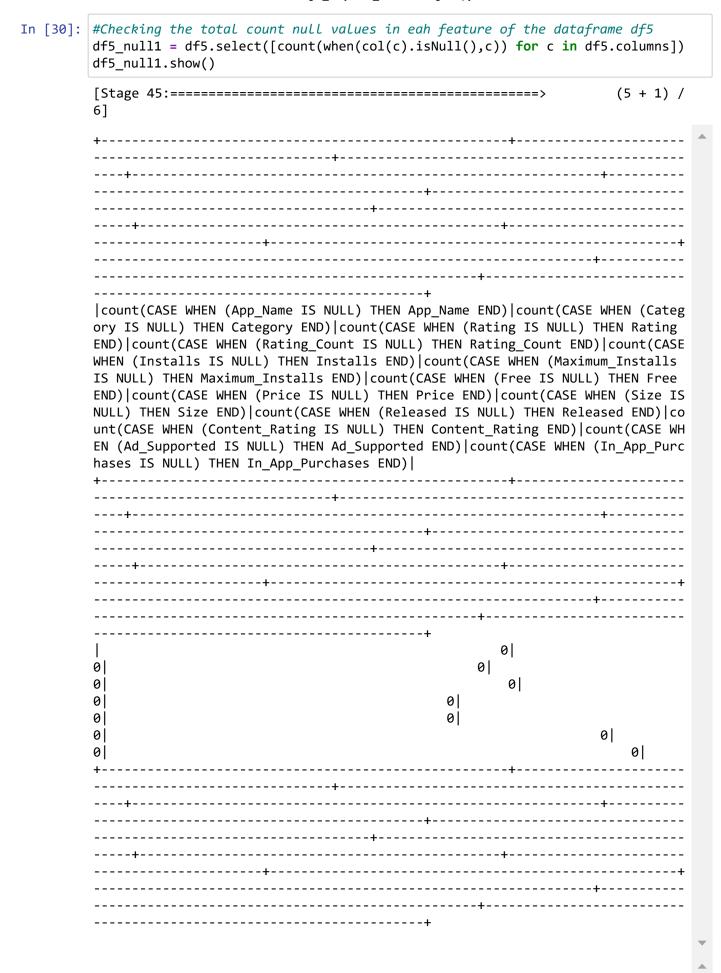
```
In [25]: #Cleaning-Step7
         # Dropping rows having null values in columns Rating and Rating Count
         df4=df3.na.drop(subset=["Rating", "Rating Count"])
         df4.show(10)
```

```
+-----
 App Name
                    Category | Rating | Rating | Count | Installs | Maximum | Ins
talls|Free|Price|Size| Released|Content Rating|Ad Supported|In App Purchases|
+-----
Gakondo | Adventure | 0.0 |
15|true| 0.0| 106|2020-02-26|
                        Everyone
                                   falsel
                                               false
| Ampere Battery Info|
                      Tools 4.4
                                      64
                                           5000
7662|true| 0.0| 296|2020-05-21|
                                     true|
                        Everyone
                                                false|
          Vibook
                Productivity|
                            0.0
                                       0
                                            50
58|true| 0.0| 376|2019-08-09|
                        Everyone|
                                   false
                                               false|
|Smart City Trichy...| Communication|
                            5.0
                                            10
                                       5|
19|true| 0.0| 186|2018-09-10|
                        Evervone
                                               false
                                    truel
                      Tools
                            0.0
                                            100|
         GROW.me
                                       0
478|true| 0.0| 626|2020-02-21|
                         Everyone|
                                                false|
                                    false
                            0.0
          IMOCCI
                     Social|
                                            50
                                       01
89|true| 0.0| 466|2018-12-24|
                                   false
                                                true|
                           Teen
|unlimited 4G data...|Libraries & Demo|
                            4.5
                                      12
                                           1000
2567|true| 0.0| 256|2019-09-23|
                        Everyone
                                     truel
                                                falsel
|The Everyday Cale...|
                                      39|
                 Lifestyle|
                            2.0
                                            500
702|true| 0.0| 166|2019-06-21|
                                                false|
                         Everyone|
                                    false
|Neon 3d Iron Tech...| Personalization| 4.7|
                                     820
                                          50000
62433|true| 0.0| 356|2019-09-22|
                         Everyone
                                      true
                                                 false
    Dodge The Cars!
                     Racing|
                            4.9
                                      55
                                            100
329|true| 0.0| 516|2020-07-30|
                         Everyone
                                    false|
                                                false|
+-----
only showing top 10 rows
```



```
In [29]:
         #Cleaning-Step8
         # Dropping rows having null values in column Size
         df5=df4.na.drop(subset=["Size"])
         df5.show(10)
```

```
+-----
  Category | Rating | Rating | Count | Installs | Maximum | Ins
         App Name
talls|Free|Price|Size| Released|Content_Rating|Ad_Supported|In_App_Purchases|
 -----
Gakondo
                   Adventure
                             0.0
                                       0|
                                             10
15|true| 0.0| 106|2020-02-26|
                                    false
                                                false|
                        Everyone
| Ampere Battery Info|
                       Tools
                             4.4
                                       64
                                            5000
7662|true| 0.0| 296|2020-05-21|
                         Everyone
                                      truel
                                                 falsel
          Vibook | Productivity |
                             0.0
                                       0|
                                             50|
58|true| 0.0| 376|2019-08-09|
                         Everyone
                                    false
                                                false
|Smart City Trichy...| Communication|
                                             10
                             5.0
                                        5|
19|true| 0.0| 186|2018-09-10|
                         Everyone
                                                false|
                                     true
         GROW.me
                       Tools
                             0.0
                                            100
                                       01
478|true| 0.0| 626|2020-02-21|
                         Everyone
                                                 falsel
                                     falsel
          IMOCCI|
                      Social|
                             0.0
                                        0|
                                             50|
89|true| 0.0| 466|2018-12-24|
                            Teen
                                    false
                                                 true
|unlimited 4G data...|Libraries & Demo|
                             4.5
                                            1000
                                       12|
2567|true| 0.0| 256|2019-09-23|
                         Everyone
                                      true
                                                 false|
|The Everyday Cale...|
                 Lifestyle|
                             2.0
                                       39|
                                            500
702|true| 0.0| 166|2019-06-21|
                         Everyone
                                     falsel
                                                 false
|Neon 3d Iron Tech...| Personalization|
                             4.7
                                      820
                                           50000
62433 | true | 0.0 | 356 | 2019 - 09 - 22 |
                                                  false
                           Everyone
                                       true
    Dodge The Cars!
                      Racing|
                             4.9
                                       55
                                            100
329|true| 0.0| 516|2020-07-30|
                         Everyone|
                                     false
                                                 false
+-----
only showing top 10 rows
```



In [31]: df5.count() Out[31]: 2190445