Shahd Mohamed Amer 22010468 Ai



import pandas as pd

df-pd.read_csv(r'/content/books.csv')

1354 rows × 23 columns

	book_ld g	poodreads_book_1d be	est_book_id	work_1d	books_count	1stm	isbn13	authors	original_publication_year	original_title		ratings_count	work_ratings_count	work_text_reviews_count	ratings_1	ratings_2	ratings_3	rat
0	1	2757052	2767052	2792775	272	439023483	9 780439e+12	Suzanne Collins	2008 0	The Hunger Games		4780653	4942365	155254	66715	127936	560092	1
1	2	3	3	4640799	491	439554934	9 780440e+12	J.K. Rowling, Mary GrandPre	1997.0	Harry Potter and the Philosopher's Stone	÷	4602479	4800065	75867	75504	101676	455024	,
2	3	41865	41865	3212258	225	315015849	9.780316e+12	Stephenia Meyer	2005.0	Twilight	+	3866839	3916824	95009	456191	436802	793319	
3	6	11870085	11870085	16827462	225	525478817	9.780525e+12	John Green	2012.0	The Fault in Our Stars	#7	2346404	2478609	140739	47994	92723	327550	
4	12	13335037	13335037	13155899	210	62024035	9 780062e+12	Veronica Roth	2011.0	Divergent		1903563	2216814	101023	36315	82870	310297	1
-		***			-	100	-		-		++		· ·			++	-	l
1349	9925	86737	86737	3877968	52	1582349177	9.781582e+12	Mary Hoffman	2002.0	City of Masks	+	12048	13305	555	314	758	3154	
1350	9937	13010211	13010211	18171867	22	1595435712	9 781596e+12	Caragh M O'Brien	2012.0	Promised		11766	12854	1260	256	1098	3565	
1351	9942	16074758	16074758	21869436	18	1442486597	9 781442e+12	Abigali Haas, Abby McDonald	2013.0	Dangerous Girls		10439	12970	2631	203	553	2029	
1352	9947	21393626	21393526	40690062	19	62320521	9 780062e+12	Maria Dahvana Headley	2015 0	Magonia		12510	13652	2910	577	1440	3881	
1353	9955	13055327	13065327	18230950	25	E02734375	9.760803e+12	Simone Elkeles	2013.0	Wild Cards		13954	15400	Acti 1521	: Wi 495	DWS 965	3331	1

CS CamScanner

Go to Settings to activate Windows.

Date df.describe() 1.1 ≆ book id goodreads book id best book id work id books count isball original publication year average rating ratings count work ratings count work text reviews count ratings_1 ratings_2 ratings_3 ratings 1.354000e+03 1.354000e+03 1.354000e+03 1354.000000 1.310000e+03 1354000e+03 1354.000000 1.354000e+0 count 1354 000000 1351 000000 1354 000000 1.354000e+03 1354.000000 1354 000000 1354 000000 17528.918021 3.060591e+0 4453.584195 5.951852e+06 6.120589e+06 8.707028e+06 50.330871 9.766700e+12 2003.422650 3.999357 9.160429e+04 9.915569a+04 5151.093058 2297.409158 5005.615953 2894.277455 6.664595e+06 6.935008e+06 9.813696e+06 61.338867 3.572069e+11 16.779301 0.224263 2.871266e+05 3.023637e+05 10730.335273 13703 507239 16259.838433 43549.306920 8.427851e+0 1.000000 1.000000e+00 1.000000e+00 1.150000e+02 1.000000 7.678361e+10 1868 000000 3.230000 6.221000e+03 8.833000e+03 49.000000 33.000000 133 000000 826.000000 1.660000e+0 1850.250000 1.537868e+05 1.537962e+05 1.375035e+06 22.000000 9.780152e+12 2003 000000 3.850000 1.759325e+04 1.918150e+04 1162 500000 305.000000 978.000000 4140.500000 6.360500e+0 4177.500000 3.305318a+06 3.422646e+06 4.005716e+06 37.000000 9.780440a+12 2008.000000 4.0000000 2.943000e+04 3.255150a+04 2208.000000 619.000000 1732.500000 6557.000000 1.079550e+0 6814 500000 58 000000 9 780805e+12 6 073800e+04 6 681275e+04 1355 000000 3644 500000 13312 250000 2 227500e+0 9.917380e+06 1.019388e+07 1.435717e+07 2011 000000 4 160000 4650 750000 max 9955.000000 2017.000000 4.740000 4.780653e+06 155254.000000 455191.000000 436802.000000 793319.000000 1.481305e+0 3.207567e+07 3.360215e+07 4.963819e+07 1314.000000 9.789424e+12 4.942365e+06 [] df.shape T (1354, 23) [] columns list = list(ef.columns) columns_list

A Company of the Printers & County Committee Francisco Incolored

- ['book id'.

'goodreads_book_id',
'best_book_id',
'wark_id',
'books_count',
'isbn',
'isbn',
'authors',

'ratings_count',
'work ratings_count',

'retings 1'.

'original_sublication_year',
'original_title',
'title',
'language_code',
'average rating',

work text reviews count'.

Activate Windows

So to Self no. to activate Windows

CS CamScanner

```
[ ] df.duplicated().sum()
∓ 0
[ ] # check missing values in categorical variables
    df[categorical].isnull().sum()
→ isbn
                         52
    authors
                          0
    original_title
                         52
    title
                          0
    language_code
                        109
    image_url
                          0
    small_image_url
                          0
    dtype: int64
[ ] # view frequency counts of values in categorical variables
    for var in categorical:
         print(df[var].value_counts())
₹.
    isbn
    439023483
                   1
    689868235
                   1
    1416914234
    1423116186
                   1
    670012092
                   1
    451222385
                   1
    60575808
                   1
    1477823832
                   1
    316127256
                   1
                   1
    802734375
    Name: count, Length: 1302, dtype: int64
    authors
    Meg Cabot
                              27
                              25
    Tamora Pierce
    L.J. Smith
                              15
    John Flanagan
                              14
    Rick Riordan
                              14
                              ٠.
    Nicola Yoon
                               1
                               1
    Barry Lyga
    Annette Curtis Klause
                               1
    Patricia Reilly Giff
                               1
    Maria Dahvana Headley
                               1
    Name: count, Length: 555, dtype: int64
    original_title
    Once
                                   2
                                   2
    Destined
    Twirtad
```

```
1
50575808
              1
.477823832
              1
116127256
              1
102734375
lame: count, Length: 1302, dtype: int64
uthors
leg Cabot
                          27
amora Pierce
                          25
..J. Smith
                          15
John Flanagan
                          14
tick Riordan
                          14
                          4.4
licola Yoon
                           1
Jarry Lyga
                           1
Innette Curtis Klause
                           1
atricia Reilly Giff
                           1
laria Dahvana Headley
                           1
lame: count, Length: 555, dtype: int64
original_title
                               2
Ince:
Destined
                               2
                               2
wisted
leaven
                               2
)ark Reunion
                               2
he Wrath & the Dawn
                               1
Mild Magic (Immortals, #1)
                               1
hysik
                               1
Hat Full of Sky
                               1
                               1
hild Cards
lame: count, Length: 1289, dtype: int64
he Hunger Games (The Hunger Games, #1)
                                                     1
Dumplin' (Dumplin', #1)
                                                     1
                                                     1
Inspoken (The Lynburn Legacy, #1)
he Oueen (The Selection, #0.4)
                                                     1
                                                     1
hasing Vermeer (Chasing Vermeer, #1)
                                                     . .
incarceron (Incarceron, #1)
                                                     1
he Awakening (The Vampire Diaries, #1)
                                                     1
irelight (Firelight, #1)
                                                     1
he Face on the Milk Carton (Janie Johnson, #1)
                                                     1
                                                     1
lild Cards (Wild Cards, #1)
lame: count, Length: 1354, dtype: int64
.anguage code
:ng
        852
en-US
        365
en-GB
          14
         11
en-CA
```

```
for var in categorical:
        frequency_distribution = df[var].value_counts() / float(len(df))
        print(f"Frequency distribution for {var}:")
        print(frequency_distribution, "\n")
Frequency distribution for isbn:
    isbn
    439023483
                  0.000739
                  0.000739
    689868235
    1416914234
                  0.000739
    1423116186
                  0.000739
    670012092
                  0.000739
                    ...
    451222385
                  0.000739
    60575808
                  0.000739
    1477823832
                  0.000739
    316127256
                  0.000739
    802734375
                  0.000739
    Name: count, Length: 1302, dtype: float64
    Frequency distribution for authors:
    authors
    Meg Cabot
                             0.019941
    Tamora Pierce
                             0.018464
    L.J. Smith
                             0.011078
    John Flanagan
                             0.010340
    Rick Riordan
                             0.010340
                                ...
    Nicola Yoon
                             0.000739
    Barry Lyga
                            0.000739
    Annette Curtis Klause
                             0.000739
    Patricia Reilly Giff 0.000739
Maria Dahvana Headley 0.000739
    Name: count, Length: 555, dtype: float64
    Frequency distribution for original_title:
    original_title
    Once
                                   0.001477
    Destined
                                   0.001477
    Twisted
                                   0.001477
    Heaven
                                   0.001477
    Dark Reunion
                                   8.001477
    The Wrath & the Dawn
                                   0.000739
    Wild Magic (Immortals, #1)
                                  0.000739
    Physik
                                   0.000739
    A Hat Full of Sky
                                  8.000739
    Wild Cards
                                   0.000739
    Name: count, Length: 1289, dtype: float64
```

view frequency distribution of categorical variables

CS CamScanner

```
0
   Destined
                                 0.001477
   Twisted
                                 0.001477
   Heaven
                                 0.001477
   Dark Reunion
                                 0.001477
                                   ...
    The Wrath & the Dawn
                                 0.000739
   Wild Magic (Immortals, #1)
                                 0.000739
    Physik
                                 0.000739
    A Hat Full of Sky
                                 0.000739
   Wild Cards
                                 0.000739
    Name: count, Length: 1289, dtype: float64
    Frequency distribution for title:
   title
   The Hunger Games (The Hunger Games, #1)
                                                      0.000739
   Dumplin' (Dumplin', #1)
                                                      0.000739
   Unspoken (The Lynburn Legacy, #1)
                                                      0.000739
    The Queen (The Selection, #0.4)
                                                      0.000739
    Chasing Vermeer (Chasing Vermeer, #1)
                                                      0.000739
                                                        . . .
    Incarceron (Incarceron, #1)
                                                      0.000739
   The Awakening (The Vampire Diaries, #1)
                                                      0.000739
    Firelight (Firelight, #1)
                                                      0.000739
    The Face on the Milk Carton (Janie Johnson, #1)
                                                      0.000739
   Wild Cards (Wild Cards, #1)
                                                      0.000739
 # check for cardinality in categorical variables
   for var in categorical:
        print(var, ' contains ', len(df[var].unique()), ' labels')
isbn contains 1303 labels
   authors contains 555 labels
    original title contains 1290 labels
    title contains 1354 labels
    language_code contains 8 labels
    image_url contains 1014 labels
    small image url contains 1014 labels
 # check missing values in categorical variables
   df[categorical].isnull().sum()
isbn
                       52
    authors
                        8
   original_title
                       52
   title
                        0
    language_code
                      109
    image_url
                        0
    small_image_url
                        0
    dtype: int64
```

U.0014//

Unce

```
[ ] # check missing values in numerical variables
    df[numerical].isnull().sum()

→ book_id

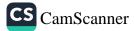
                                   0
    goodreads_book_id
                                   0
    best_book_id
                                   0
                                   θ
    work_id
                                   0
    books_count
                                  44
    isbn13
     original_publication_year
                                   3
    average_rating
                                   0
    ratings_count
                                   0
    work_ratings_count
                                   8
    work_text_reviews_count
    ratings_1
                                   8
    ratings_2
                                   0
                                   8
    ratings_3
                                   8
    ratings_4
                                   0
    ratings_5
    dtype: int64
[ ] # Calculate the percentage of missing values in the numerical variables
    missing_values = df[numerical].isnull().sum()
     percentage_missing = ((missing_values / len(df)) * 100).sort_values(ascending=False)[missing_values>0]
    percentage_missing

→ isbn13

                                  3.249631
    original_publication_year
                                  0.221566
    dtype: float64
[ ] df.isnull().sum()
                                    8
- book_id
                                    0
     goodreads_book_id
    best_book_id
                                    8
                                    0
    work_id
                                    8
    books_count
                                   52
    isbn
    isbn13
                                   44
                                    8
    authors
                                    3
    original_publication_year
                                   52
    original_title
    title
                                    8
                                  109
    language_code
                                    0
     average_rating
                                    0
    ratings count
    work_ratings_count
                                    8
    work_text_reviews_count
                                    0
    ratings 1
```

+ + = = 5 5 5

-								*******												
		- ;;-	-	-	-			416				-		-						
	- 4			4940/99	***	*****		Country.	2007	the Philosophers					*****	*****	455024	*******	wines	
	2	,	,	4543199	491	438334834	9./804400-12	Mary	1997.0	Stone		4002479	4833,065	73007	75504	101010	433024	1156318	3011543	esse
								GrandPre		******										***
								TK												
								Fillwirth.		Horry Porter and										
					***			Mag		Phone: Planter and				-		*****	****	******		
		2	2.0					termorrale.		Azkaban				75						ismo
								BOCK		Ambunia										
								3.10		Many Probes sunt										
	41			-	341	*393300110	-	Mary	2403.0	-		******	1044/244		****	31411	-	-	1144000	200
	**	,	•	2002202	347	*****	0.780470++47	- COURSE	3003.0	- Shoom		*****		THERE	06.30	74577	****	101177	*******	
								Rovino.		Hung France wie										
	777			=======			270000-17	- test	1222.0		-			31177	2000	27	717715	5:000	:::::::::::::::::::::::::::::::::::::::	tens
								(I mathat		· · · · · · · · · · · · · · · · · · ·										
	**			****	****			JK					***			40,000	****	11-4 1400-		
	24	6		3046572	332	#30130600	9.7801399+12	Ruwlery.	2000.0	Harry Potter and the Gobiet of		1753043	1868642	21004	6676	20210	151785	191926	1195045	
		•		3910312	222	43613600	B.1801388-12	Mary	1500,0	Fau			-	31001	00.0	20210	131103	101020	1193013	
								Granillya												
								1.10		Harry Potter and										
	25	136251	135251	2953218	263	545010225	9.7805459+12	Rowling	2007.0		1	1745574	1847395	51942	9353	22245	113545	383914	1318227	alb
								Grandine		PLINE					2000		+			95
		-	-					1.1												**
								Rowling.		Harry Patter and										
	27	,		4133542	2**	439705960	9.700440e+12	Mary	2005.0	the flatf-Clood		1670023	1705075	27520	7300	21515	120222	459026	1151421	**
								GrandPile		France										
	em .	262041	262041	2962402	76	\$450agnes	0.7805456-12	-	1002.0	Culpinning		100060	201125	6602	1106	1295	7020	30666	164040	
	_	*****	802041			940044201	0.7605450-12	Rowling	1552.5	Potter Bared Set		1,0000		****	1100	1200	7020	30000	104040	#5
								0.207.0		Livery Serline										
	1/31	70	10	2145/19/8	•	A HELF THEM	W /901440m+17	Rowling	SAME I			740 TH	AUIA	100	2014	1995	1145	48911	21114K	**
				******	***			-	*****	Potter #1-6			* ****	***		Acres	ate th	intritution		
								_		The Magical						en in	and the same	11:5:00		-
	/Una	465449	41144	4/1/92		DECTION OF	9.7804200+12	:=.:Dure!	2901.0	Winnish, orl Hanne		13820	10140	491	329	110	S. 111.	3793	9336	
	1310		-6,5-40	40 11 92		0-27 (305 (A	3.004236712	Collect	2001,0	Nomes: V		13040	13143	701	369	1123	3100	2,593	9332	asset
		-			7.			- ++	Connected to Python 1 Gr	Tracery .				***		1.75	11.24			



[] columns = ['books_count','average_rating','ratings_count','work_ratings_count','work_text_reviews_count']
columnsdrop = df.columns.difference(columns)

Drop all columns except the specified ones harrypotter = harrypotter.drop(columns=columnsdrop) harrypotter

•	books_count	average_rating	ratings_count	work_ratings_count	work_text_reviews_count
1	491	4.44	4002479	4800065	75857
6	376	4.53	1832823	1969375	35099
	307	4.46	1735358	1840548	28685
9	398	4.37	1779331	1906199	34172
10	332	4.53	1753043	1868642	31084
11	263	4.61	1746574	1847395	51942
12	275	4.54	1678823	1785676	27520
96	76	4.74	190050	204125	8503
613	0	4.73	24018	25274	882
103	6 42	3.96	13820	15145	267
126	5 5	4.40	10738	11732	185

the_most_selling_books = Harry Potter.sort_values(by='work_ratings_count', ascending=False)
the_most_selling_books

•					
	books_count	average_rating	ratings_count	work_ratings_count	work_text_reviews_count
1	491	4.44	4502479	4800065	75857
6	370	4.53	1832823	1989375	35099
9	305	4.37	1779331	1906199	34172
10	332	4.53	1753043	1868642	31084
11	263	4.61	1740574	1847395	51942
8	307	4 45	1735368	1840548	28685
12	275	4.54	1678823	1785070	27520
96	76	4.74	190050	204125	6508
613	8	4.73	24618	26274	882
1036	42	3.98	13820	15145	267
1266	5	4.40	10735	11732	185

^[] avg_rate=the_most_selling_books['average_rating'].mean()
avg_rate

[]

^{4.482727272727273}

```
| | # Calculate the percentage of missing values in the categorical variables
   missing_values - df[categorical].isnull().sum()
   percentage missing * ((missing values / len(df)) * 100).sort_values(escending-felse)[missing values+0]
   percentage_missing
== language_code
                    8.050222
   1shn
                     3.848473
   original_title
                    3.846473
   dtype: float64
[ ] # find numerical variables
   numerical = [ver for ver in df.columns if df[var].dtype(='0']
   print('There are () numerical variables\n'.format(len(numerical)))
   print('The numerical variables are :', numerical)
There are 16 numerical variables
```

The numerical variables are : ['book_id', 'goodreads_book_id', 'best_book_id', 'books_count', 'istral3', 'original_publication_year', 'average_ratings_count', 'work_text_reviews_count', 'ratings_1', 'ratings_2', 'ratings_3',

2345404

1903563

2478509

2215814

140739

101023

47994

36315

92723

82370

327550

310297

598471

673028

1311871

1114304

∓•	book_ld	goodreads_book_1d	best_book_id	work_1d	books_count	istm13	original_publication_year	average_rating	ratings_count	work_ratings_count	work_text_reviews_count	ratings_1	ratings_2	ratings_3	ratings_4	ratines_5
	1	2767052	2767052	2792775	272	9.780439++12	2008.0	434	4780653	4942365	155254	56715	127936	560092	1481305	2706317
	2	3	3	4640799	491	9.7804408-12	1997.0	4.44	4602479	4800065	75857	75504	101676	455024	1156318	3011543
	2 3	41955	41065	3212258	226	9 780316e+12	2005 0	3.57	3066839	3916824	95009	456191	436802	793319	875073	1355439

4.25

4.24

2012.0

2011.0

Code + lext

[] # view the numerical variables
 df[numerical].head()

12

11870085

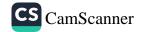
13335037

11870085 16827462

13335037 13155899

225 9.7805250+12

210 9.7800528+12



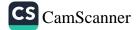
```
print('There are {} categorical variables\n'.format(len(categorical)))
    print('The categorical variables are :\n\n', categorical)
There are 7 categorical variables
    The categorical variables are :
     ['isbn', 'authors', 'original title', 'title', 'language code', 'image url', 'small image url']
[ ] # view the categorical variables
    df[categorical].head()
₹
              isbn
                                      authors
                                                                    original title
                                                                                                                         title language_code
                                                                                                                                                                                     image_url
                                                                                                                                                                                                                              small_image_url
     0 439023483
                               Suzanne Collins
                                                                 The Hunger Games
                                                                                       The Hunger Games (The Hunger Games, #1)
                                                                                                                                            eng https://images.gr-assets.com/books/1447303603m...
                                                                                                                                                                                                 https://images.gr-assets.com/books/1447303603s.
     1 439554934 J.K. Rowling, Mary GrandPré Harry Potter and the Philosopher's Stone Harry Potter and the Sorcerer's Stone (Harry P.
                                                                                                                                            eng https://images.gr-assets.com/books/1474154022m__
                                                                                                                                                                                                 https://images.gr-assets.com/books/1474154022s.
     2 316015849
                               Stephenie Meyer
                                                                            Twillight
                                                                                                             Twilight (Twilight, #1)
                                                                                                                                          en-US https://images.gr-assets.com/books/1361039443m
                                                                                                                                                                                                 https://images.gr-assets.com/books/1361039443s
     3 525478817
                                   John Green
                                                               The Fault in Our Stars
                                                                                                           The Fault in Our Stars
                                                                                                                                            eng https://images.gr-assets.com/books/1360206420m.
                                                                                                                                                                                                 https://images.gr-assets.com/books/1360206420s.
         62024035
                                 Veronica Roth
                                                                          Divergent
                                                                                                         Divergent (Divergent, #1)
                                                                                                                                            ang https://images.gr-assets.com/books/1328559506m__
                                                                                                                                                                                                 https://images.gr-assets.com/books/1328559506a
```

] # find categorical variables

categorical = [var for var in df.columns if df[var].dtype=='0']

] # check missing values in categorical variables

df[categorical].isnull().sum()



[]	df.head	()													-			
₹	book	k_ld go	oodreads_book_1d	best_book_1d	work_1d	books_count	İsbn	1sbn13	authors	original_publication_year	original_title	ratings_cour	t work_ratings_count	work_text_reviews_coum	t_ratings_1	ratings_2	ratings_3	rating
	0	1	2767052	2767052	2792775	272	439023483	9.780439++12	Suzanne	2008 0	The Hunger Games	478069	3 4942365	15625	4 66715	127936	560092	14811
	1	2	3	3	4640799	491	439554934	9.780440e+12	J K Rowling, Mary GrandPre	1997.0	Harry Potter and the Philosopher's Stone	460247	9 4800065	7586	7 75504	101676	455024	11563
	2	3	41865	41865	3212258	226	316015849	9.780316e+12	Stephenie Mayer	2005 0	Twilight	386683	9 3915824	9600	9 456191	436802	793319	8751
	3	6	11870085	11870085	16827462	226	525478817	9.780525e+12	Jahn Green	2012.0	The Fault in Our Stars	234640	4 2478609	14073	9 47994	92723	327550	698-
	4	12	13335037	13335037	13155899	210	62024035	9.780062+12	Veronica Roth	2011.0	Divergent	190356	3 2216814	10102	3 36315	82870	310297	6731
	5 rows =	23 colum	nns															
	•																	
[]	df.teil	L()																
₹	1	book_id	goodreads_book_id	d best_book_	id work_	id books_com	int :	ishn ish	er13 author	ors original_publication_	year original_ti	itle ratings	count work_ratings_c	ount work_text_reviews_	count ratio	ngs_1 ratin	gs_2 rating	p_3 re
	1349	9925	8673	7 867	37 387796	58	52 1582345	9177 9.7815826	s+12 Hoffe	lary 20	02.0 City of Ma	asks	12048 1	3385	555	314	758 3	1154
	1350	9937	1301021	1 130102	1817186	57	22 159643	5712 9 781596e	0-12 0-12	FM 20	12.0 Prom	ited	11766 1	2884	1260	256	1098 3	665

18 1442486597 9.7814420+12

25 802734375 9.780803e+12

62320521 9.780062e+12

Abby McDonald Maria

Dahvana

Headey

Simone

1351

1352

1353

9942

9947

9955

16074758

21393526

13065327

16074758 21869436

21393526 40690062

13065327 18230950

2013.0 Dangerous Girls ...

Magonia

Wild Cards

2015.0

2013.0

553

965

2029

3331

2631

Actizante Wisolows 1440

1514 406

Go to Settings to activate Windows:

12970

13652

15400

10439

12510

13954