```python
1  #%%
2  import re
3  import pickle
4  import operator
5  import numpy as np
6  import pandas as pd
7  import seaborn as sns
8  import matplotlib.pyplot as plt
9  from collections import Counter
10 from scipy.sparse import csr_matrix
11 from pandas.api.types import is_numeric_dtype
12 from sklearn.neighbors import NearestNeighbors
13 from sklearn.feature_extraction import
   DictVectorizer
14 from sklearn.metrics.pairwise import
   cosine_similarity
15 from sklearn.feature_extraction.text import
   TfidfVectorizer
16
17 import warnings
18 warnings.filterwarnings('ignore')
19 #%%
20 books = pd.read_csv(r'Books.csv', delimiter=';',
   encoding='ISO-8859-1', on_bad_lines='skip')
21 users = pd.read_csv(r'Users.csv', delimiter=';',
   encoding='ISO-8859-1', on_bad_lines='skip')
22 ratings = pd.read_csv(r'Book-Ratings.csv',
   delimiter=';', encoding='ISO-8859-1', on_bad_lines=
   'skip')
23
24
25 print('Books: ', books.shape)
26 print('Users: ', books.shape)
27 print('Book Rating: ', ratings.shape)
28
29 #%%
30 books.head()
31 #%%
32 # Pre-Processing
33
34 print("Coluns: ", books.columns)
```

```python
35  #%%
36  books.drop(['Image-URL-L', 'Image-URL-M', 'Image-
    URL-S'], axis=1, inplace=True)
37  books.head()
38  #%%
39  books.isnull().sum()
40  #%%
41  books.loc[books['Book-Author'].isnull()]
42  #%%
43  books.loc[books['Publisher'].isnull()]
44  #%%
45  books.at[187689, 'Book-Author'] = 'Other'
46  books.at[118033, 'Book-Author'] = 'Other'
47  books.at[128890, 'Publisher'] = 'Other'
48  books.at[129037, 'Publisher'] = 'Other'
49  #%%
50  books.loc[books['Book-Author'].isnull()] #Null
    authors are removed
51  #%%
52  books['Year-Of-Publication'].unique()
53  #%%
54  pd.set_option('display.max_colwidth', -1)
55  #%%
56  books.loc[books['Year-Of-Publication'] == 'DK
    Publishing Inc']
57  #%%
58  books.loc[books['Year-Of-Publication'] == '
    Gallimard']
59  #%%
60  # Change Year of Publication to actual Year and not
     Publisher Name
61  books.at[209538, 'Publisher'] = 'DK Publishing Inc'
62  books.at[209538, 'Year-Of-Publication'] = 2000
63  books.at[209538, 'Book-Title'] = 'DK Readers:
    Creating the X-Men, How It All Began (Level 4:
    Proficient Readers)'
64  books.at[209538, 'Book-Author'] = 'Michael
    Teitelbaum'
65
66
67  books.at[221678, 'Publisher'] = 'DK Publishing Inc'
```

```
68 books.at[221678, 'Year-Of-Publication'] = 2000
69 books.at[221678, 'Book-Title'] = 'DK Readers:
   Creating the X-Men, How Comic Books Come to Life (
   Level 4: Proficient Readers)'
70 books.at[221678, 'Book-Author'] = 'James Buckley'
71
72
73 books.at[220731, 'Publisher'] = 'Gallimard'
74 books.at[220731, 'Year-Of-Publication'] = 2003
75 books.at[220731, 'Book-Title'] = 'Peuple du ciel,
   suivi de les bergers'
76 books.at[220731, 'Book-Author'] = 'Jean-Marie,
   Gustave Le ClÃ?Â©zio'
77
78
79 #%%
80 #Converting year of publication in Numbers
81 books['Year-Of-Publication']=books['Year-Of-
   Publication'].astype(float)
82
83 print(sorted(list(books['Year-Of-Publication'].
   unique())))
84 #%%
85 #Replacing invalid years with max year
86 count = Counter(books['Year-Of-Publication'])
87 [k for k, v in count.items() if v ==max(count.
   values())]
88 #%%
89
```