

## BOOK RECOMMENDATION SYSTEM USING K-NEAREST NEIGHBOUR

AUTHOR : SYED SHAHID 



In [ ]:

```
dataset link ="https://www.kaggle.com/datasets/jealousleopard/goodreadsbooks"  
Githublink = "https://github.com/syedshaahid02/syed"
```

In [ ]:

```
# Importing Libraries
```

In [1]:

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
import seaborn as sns  
## for interactive plots  
import ipywidgets  
from ipywidgets import interact  
from ipywidgets import interact_manual  
## For Ignoring Warning ErrorMessage  
from warnings import filterwarnings  
filterwarnings('ignore')
```

In [2]:

```
# Importing Dataset
```

In [3]:

```
df =pd.read_csv("books.csv",error_bad_lines = False)
```

```
b'Skipping line 3350: expected 12 fields, saw 13\nSkipping line 4704: expected 12 fields, saw 13\nSkipping line 5879: expected 12 fields, saw 13\nSkipping line 8981: expected 12 fields, saw 13\n'
```





In [16]:

df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11123 entries, 0 to 11122
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   bookID          11123 non-null   int64  
 1   title            11123 non-null   object  
 2   authors          11123 non-null   object  
 3   average_rating   11123 non-null   float64 
 4   isbn             11123 non-null   object  
 5   isbn13           11123 non-null   int64  
 6   language_code    11123 non-null   object  
 7   num_pages        11123 non-null   int64  
 8   ratings_count   11123 non-null   int64  
 9   text_reviews_count 11123 non-null   int64  
 10  publication_date 11123 non-null   object  
 11  publisher        11123 non-null   object  
dtypes: float64(1), int64(5), object(6)
memory usage: 1.0+ MB
```

## Feature engineering

#extract important features #reducing the size of features #creating new features from existing onesw

In [17]:

#dropping unnecessary column present in database

In [18]:

df.drop(['bookID', 'isbn', 'isbn13'], axis=1, inplace=True)

In [19]:

df.publication\_date

Out[19]:

```
0      9/16/2006
1      9/1/2004
2      11/1/2003
3      5/1/2004
4      9/13/2004
       ...
11118    12/21/2004
11119    12/1/1988
11120    8/1/1993
11121    2/27/2007
11122    5/28/2006
Name: publication_date, Length: 11123, dtype: object
```

In [20]:

## Creating new year columns

In [21]:

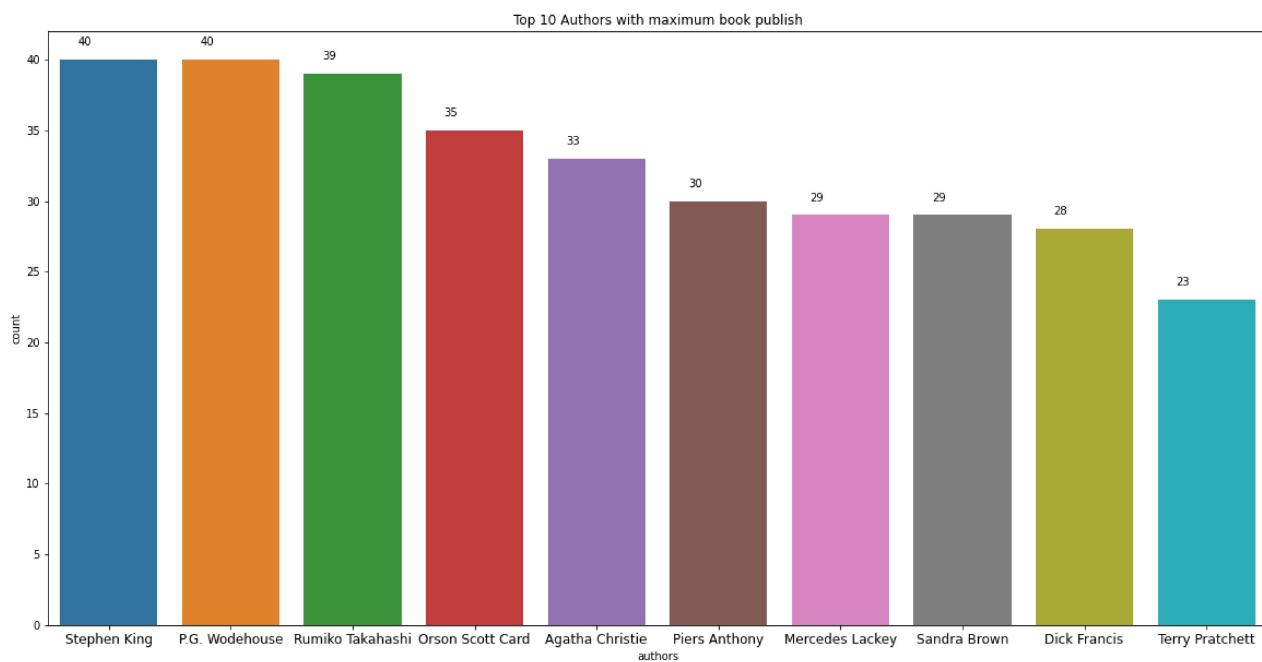
```
df['year'] = df['publication_date'].str.split('/')
df['year'] = df['year'].apply(lambda x: x[2])
```





In [33]:

```
plt.figure(figsize = (20, 10))
ax = sns.countplot(x = 'authors', data = df,
order = df['authors'].value_counts().iloc[:10].index)
plt.title("Top 10 Authors with maximum book publish")
plt.xticks(fontsize = 12)
for p in ax.patches:
    ax.annotate(format(p.get_height()), (p.get_x() + 0.15, p.get_height() + 1))
plt.show()
```



In [ ]:

```
# Sort All Value Count of Language_code.
```

In [34]:

```
df.language_code.value_counts()
```

Out[34]:

```
eng      8908
en-US    1408
spa      218
en-GB    214
fre      144
ger      99
jpn      46
mul      19
zho      14
grc      11
por      10
en-CA     7
ita      5
lat      3
enm      3
rus      2
swe      2
nor      1
tur      1
ale      1
srp      1
nl       1
msa      1
glg      1
ara      1
gla      1
wel      1
Name: language_code, dtype: int64
```

In [ ]:

```
##Creating Groupby Function base on Language_code Column and getting Required Output.
```

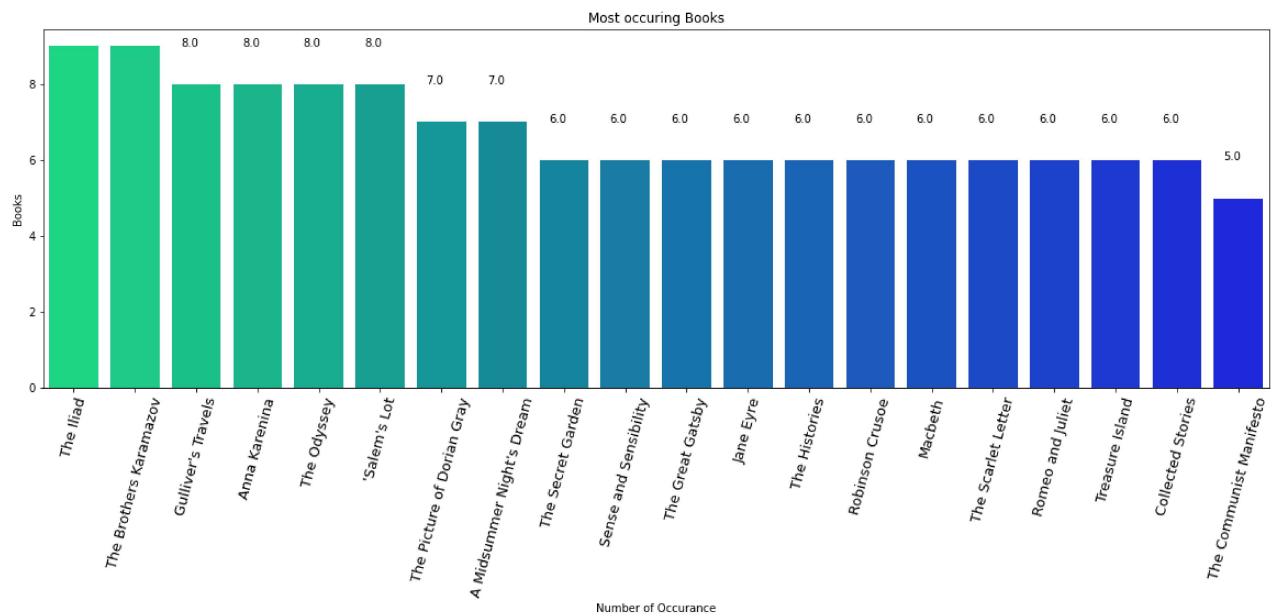


In [ ]:

#Plotting Barplot to find most occuring book in our data.

In [38]:

```
### Plotting BarPlot to find most occuring book in our data
plt.figure(figsize = (20, 6))
book = df['title'].value_counts()[:20]
ax = sns.barplot(x = book.index, y = book,
                  palette = 'winter_r')
plt.title("Most occuring Books")
plt.xlabel("Number of Occurance")
plt.ylabel("Books")
plt.xticks(rotation = 75, fontsize = 13)
for p in ax.patches:
    ax.annotate(format(p.get_height()), (p.get_x()+0.15, p.get_height()+1))
plt.show()
```

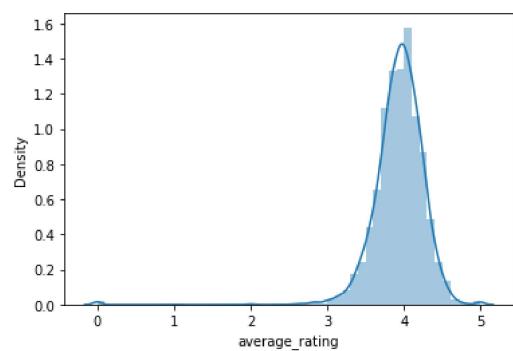


In [ ]:

#Plotting Distribution Graph on Average\_Rating.m

In [39]:

```
### Ploting Distribution Graph on Average_Rating.
sns.distplot(df['average_rating'])
plt.show()
```



In [40]:

#Sorting Dataset related with maximum Average\_Rating Column















