

## **MAKE UP SESSION**

**NAME: PUSHKAR HARSHKUMAR PARAKH**

**PRN: 202401040313**

**CLASS:CS4**

**ROLL NO.:54**

```
import pandas as pd
import numpy as np

# Load the dataset
df = pd.read_csv('book_reviews.csv')

# Display basic info
df.info()
df.head()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32881 entries, 0 to 32880
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype  
---  --
 0   Unnamed: 0   32881 non-null  int64   
 1   Book         32881 non-null  object  
 2   Review       31772 non-null  object  
 3   Review Date  32881 non-null  object  
dtypes: int64(1), object(3)
memory usage: 1062.7+ KB
```

Unnamed: 0	Book	Review	Review Date
0	To Kill a Mockingbird	// gentle reminder that this is not the time ...	March 24, 2022
1	To Kill a Mockingbird	Very good story. I know I am riding a serious...	May 24, 2011
2	To Kill a Mockingbird	Very good looking for a new book but don't want to...	December 10, 2020
3	To Kill a Mockingbird	To Kill a Mockingbird Harper Lee (To Kill a Mo...	July 1, 2022
4	To Kill a Mockingbird	Why is it when I pick up To Kill a Mockingb...	October 25, 2009

1. Display the total number of reviews per book.

```
df['Book'].value_counts()
```

```
Book
To Kill a Mockingbird      30
Homeside Days              30
Aunt and Other Stories     30
Rabbit Redux (Harriet Angstrom, #2)  30
Sentimental Education      30
...
The Triple Mirror of the Self      5
The Trusting and the Misled       4
Paranoid - Volume 1              3
[책]신학 1권                     2
Hallowell                        2
Name: count, Length: 3896, dtype: int64
```

```
File Edit Selection View Go Run Terminal Help ← → Search
```

2. Find the number of missing reviews.

```
df['Review'].isna().sum()
```

```
389
```

3. Drop rows with missing reviews.

```
df_cleaned = df.dropna(subset=['Review'])
df_cleaned.shape
```

```
(31772, 4)
```

4. Convert 'Review Date' to datetime format.

```
df_cleaned['Review Date'] = pd.to_datetime(df_cleaned['Review Date'], errors='coerce')
df_cleaned['Review Date'].head()
```

```
C:\Users\shiva\AppData\Local\Temp\ipykernel_5844\5812x49379.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df_cleaned['Review Date'] = pd.to_datetime(df_cleaned['Review Date'], errors='coerce')
```

	Review Date
0	2022-02-24
1	2011-05-24
2	2020-12-10
3	2022-07-01
4	2009-10-25

Name: Review Date, dtype: datetime64[ns]

5. Find the earliest and latest review date.

```
df_cleaned['Review Date'].min(), df_cleaned['Review Date'].max()
```

```
(Timestamp('2007-02-10 00:00:00'), Timestamp('2022-07-04 00:00:00'))
```

6. Count the number of reviews per year.

```
df_cleaned['Year'] = df_cleaned['Review Date'].dt.year
df_cleaned['Year'].value_counts().sort_index()
```

```
C:\Users\shiva\AppData\Local\Temp\ipykernel_5844\5812x49379.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df_cleaned['Year'] = df_cleaned['Review Date'].dt.year
```

```
File Edit Selection View Go Run Terminal Help Search

6. Count the number of reviews per year.

df_cleaned['year'] = df_cleaned['review date'].dt.year
df_cleaned['year'].value_counts().sort_index()

C:\Users\shiva\AnacondaLocal\Tome\ipykernel_5846\118329961.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
df_cleaned['year'] = df_cleaned['review date'].dt.year

Year
2007.0    446
2008.0    739
2009.0    568
2010.0    525
2011.0    876
2012.0   1177
2013.0   1868
2014.0   1886
2015.0   1613
2016.0   2124
2017.0   2685
2018.0   2949
2019.0   2643
2020.0   3348
2021.0   3388
2022.0   3962
2023.0   2151
Name: count, dtype: int64

7. Get top 5 most reviewed books.

df_cleaned['book'].value_counts().head()

Book
The Pilgrim's Progress    38
The Perils of a Soldier (c. Auguste Dupin, #3)    38
Kidnapped (David Balfour, #1)    38
Of Love and Shadows    38
The Garden Party and Other Stories    38
Name: count, dtype: int64

8. Count how many reviews contain the word 'love'.
```

```
File Edit Selection View Go Run Terminal Help Search

9. Add a column with the length of each review.

df_cleaned['review length'] = df_cleaned['review'].str.len()
df_cleaned[['review', 'review length']].head()

C:\Users\shiva\AnacondaLocal\Tome\ipykernel_5846\164416654.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
df_cleaned['review length'] = df_cleaned['review'].str.len()

review review length
0 // gentle reminder that this is not the time... 8167
1 'yqjnh0 stars. I know I am raising a serious... 4289
2 'yqjnh0 Looking for a new book but don't want to... 2779
3 To Kill a Mockingbird, Harper Lee's To Kill a Mo... 3365
4 Why is it when I pick up 'To Kill A Mockingbi... 2542

10. Find average review length per book.

df_cleaned.groupby('book')['review length'].mean().sort_values(ascending=False).head()

Book
Cloud Atlas    6511.488889
Infinite Jest    6299.633333
Gravity's Rainbow    5181.888889
The Glass Bead Game    5155.433333
The Recognition    4764.500000
Name: review length, dtype: float64

11. Show standard deviation of review lengths using NumPy.

np.std(df_cleaned['review length'].dropna())

2289.849794523991

12. Create a column for word count in each review.

df_cleaned['word count'] = df_cleaned['review'].apply(lambda x: len(str(x).split()))
df_cleaned[['review', 'word count']].head()

C:\Users\shiva\AnacondaLocal\Tome\ipykernel_5846\164374763.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user\_guide/indexing.html#returning-a-view-versus-a-copy
df_cleaned['word count'] = df_cleaned['review'].apply(lambda x: len(str(x).split()))
```