

```

1 ## Loading the dataset
2 import pandas as pd
3 data = pd.read_csv('/content/all_kindle_review.csv')

```

```
1 data
```

	Unnamed: 0.1	Unnamed: 0	asin	helpful	rating	reviewText	reviewTime	reviewerID	reviewerName	summary
0	0	11539	B0033UV8HI	[8, 10]	3	Jace Rankin may be short, but he's nothing to ...	09 2, 2010	A3HHXRELK8BHQG	Ridley	Entertaining But Average
1	1	5957	B002HJV4DE	[1, 1]	5	Great short read. I didn't want to put it dow...	10 8, 2013	A2RGNZ0TRF578I	Holly Butler	Terrific menage scenes!
2	2	9146	B002ZG96I4	[0, 0]	3	I'll start by saying this is the first of four...	04 11, 2014	A3S0H2HV6U1I7F	Merissa	Snapdragon Alley
3	3	7038	B002QHWOEU	[1, 3]	3	Aggie is Angela Lansbury who carries pocketboo...	07 5, 2014	AC4OQW3GZ919J	Cleargrace	very light murder cozy
4	4	1776	B001A06VJ8	[0, 1]	4	I did not expect this type of book to be in li...	12 31, 2012	A3C9V987IQHOQD	Rjostler	Book
...
11995	11995	2183	B001DUGORO	[0, 0]	4	Valentine cupid is a vampire-Jena and lan ano...	02 28, 2014	A1OKS5Q1HD8WQC	lisa jon jung	jena
11996	11996	6272	B002JCSFSQ	[2, 2]	5	I have read all seven books in this series. Ap...	05 16, 2011	AQRSPXLNEQAMA	TerryLP	Peacekeepers Series
11997	11997	12483	B0035N1V7K	[0, 1]	3	This book really just wasn't my cuppa. The si...	07 26, 2013	A2T5QLT5VXOJAK	hwilson	a little creepy
11998	11998	3640	B001W1XT40	[1, 2]	1	tried to use it to charge my kindle, it didn't...	09 17, 2013	A28MHD2DDY6DXB	Allison A. Slater "Gryphon50"	didn't work
11999	11999	11398	B003370JUS	[5, 6]	3	Taking Instruction is a look into the often hi...	07 5, 2012	A3JUXLB4K9ZXCC	Dafna Yee	If you like BDSM with a touch of romance, this...

12000 rows × 11 columns

Next steps:

[Generate code with data](#)
[View recommended plots](#)
[New interactive sheet](#)

```
1 df = data[['reviewText', 'rating']]
```


```
1 df.head()
```

	reviewText	rating
0	Jace Rankin may be short, but he's nothing to ...	3
1	Great short read. I didn't want to put it dow...	5
2	I'll start by saying this is the first of four...	3
3	Aggie is Angela Lansbury who carries pocketboo...	3
4	I did not expect this type of book to be in li...	4


Next steps:

[Generate code with df](#)[View recommended plots](#)[New interactive sheet](#)

1 df.shape


 (12000, 2)

```
1 ## Missing values
2 df.isnull().sum()
```




	0
reviewText	0
rating	0

```
1 ## Unique values
2 df['rating'].unique()
```

 array([3, 5, 4, 2, 1])


1 df['rating'].value_counts()



rating	count
5	3000
4	3000
3	2000
2	2000
1	2000

```
1 ## Preprocesssing and Cleaning
2 ## Positive review is 1 and negative review is 0
3 df.loc[:, 'rating'] = df['rating'].apply(lambda x : 0 if x<3 else 1 )
```


1 df['rating'].value_counts()



rating	count
1	8000
0	4000

```
1 ## Preprocessing Steps
2 ## 1. Lower all text
```

```
1 df['reviewText'] = df['reviewText'].str.lower()
2 df.head()
```



<ipython-input-14-cc7f31db30b9>: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-vs-returning-a-copy

```
df['reviewText'] = df['reviewText'].str.lower()
```

	reviewText	rating
0	jace rankin may be short, but he's nothing to ...	1
1	great short read. i didn't want to put it dow...	1
2	i'll start by saying this is the first of four...	1
3	aggie is angela lansbury who carries pocketboo...	1
4	i did not expect this type of book to be in li...	1

Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

```
1 ## Removing Special Characters
2 import re
3 import nltk
4 from nltk.corpus import stopwords
5 nltk.download('stopwords')
```

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
True

```
1 from bs4 import BeautifulSoup
```

```
1 ## Removing special characters
2 df['reviewText']=df['reviewText'].apply(lambda x:re.sub('[^a-z A-z 0-9-]+', '',x))
3 ## Remove the stopwords
4 df['reviewText']=df['reviewText'].apply(lambda x: " ".join([y for y in x.split() if y not in stopwords.words('english')]))
5 ## Remove url
6 df['reviewText']=df['reviewText'].apply(lambda x: re.sub(r'(http|https|ftp|ssh):\/\/([\\w_-]+(?:\\.[\\w_-]+)+))([\\w.,@?^=%&:/~+#-]*\\S+)', '',x))
7 ## Remove html tags
8 df['reviewText']=df['reviewText'].apply(lambda x: BeautifulSoup(x, 'lxml').get_text())
9 ## Remove any additional spaces
10 df['reviewText']=df['reviewText'].apply(lambda x: " ".join(x.split()))
```

<ipython-input-17-86b4aeb70b16>:2: 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-vs-copy
df['reviewText']=df['reviewText'].apply(lambda x:re.sub('[^a-z A-z 0-9-]+', '',x))
<ipython-input-17-86b4aeb70b16>:4: 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-vs-copy
df['reviewText']=df['reviewText'].apply(lambda x: " ".join([y for y in x.split() if y not in stopwords.words('english')]))
<ipython-input-17-86b4aeb70b16>:6: 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-vs-copy
df['reviewText']=df['reviewText'].apply(lambda x: re.sub(r'(http|https|ftp|ssh):\/\/([\\w_-]+(?:\\.[\\w_-]+)+))([\\w.,@?^=%&:/~+#-]*\\S+)', '',x))
<ipython-input-17-86b4aeb70b16>:8: 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-vs-copy
df['reviewText']=df['reviewText'].apply(lambda x: BeautifulSoup(x, 'lxml').get_text())
<ipython-input-17-86b4aeb70b16>:10: 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-vs-copy
df['reviewText']=df['reviewText'].apply(lambda x: " ".join(x.split()))

```
1 df.head()
```

	reviewText	rating
0	jace rankin may short hes nothing mess man hau...	1
1	great short read didnt want put read one sitti...	1
2	ill start saying first four books wasnt expect...	1
3	aggie angela lansbury carries pocketbooks inst...	1
4	expect tvee book library pleased find price right	1

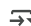
Next steps: [Generate code with df](#) [View recommended plots](#) [New interactive sheet](#)

```
1 ## Lematizer
2 from nltk.stem import WordNetLemmatizer
3 lemmatizer = WordNetLemmatizer()
4 nltk.download('wordnet')
```

[nltk_data] Downloading package wordnet to /root/nltk_data...
True

```
1 def lemmatize_words(text):
2     return " ".join([lemmatizer.lemmatize(word) for word in text.split()])
```

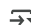

```
1 df['reviewText']=df['reviewText'].apply(lambda x:lemmatize_words(x))
```


 <ipython-input-25-7fa46fdadeba>: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-vs-copy

```
df['reviewText']=df['reviewText'].apply(lambda x:lemmatize_words(x))
```

```
1 df.head()
```

	reviewText	rating	
0	jace rankin may short he nothing mess man haul...	1	
1	great short read didnt want put read one sitti...	1	
2	ill start saying first four book wasnt expecti...	1	
3	aggie angela lansbury carry pocketbook instead...	1	
4	expect tvoe book librarv pleased find price right	1	

Next steps:

[Generate code with df](#)

[View recommended plots](#)

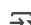
[New interactive sheet](#)

```
1 ## Train test split
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(df['reviewText'],df['rating'], test_size=0.20, random_state=42)
```

```
1 ## text to vector
2 from sklearn.feature_extraction.text import CountVectorizer
3 bow = CountVectorizer()
4 X_train_bow = bow.fit_transform(X_train).toarray()
5 X_test_bow = bow.transform(X_test).toarray()
```

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 tfidf = TfidfVectorizer()
3 X_train_tfidf = tfidf.fit_transform(X_train).toarray()
4 X_test_tfidf = tfidf.transform(X_test).toarray()
5
```

```
1 X_train_bow
```


 array([[0, 0, 0, ..., 0, 0, 0],
[0, 0, 0, ..., 0, 0, 0],
[0, 0, 0, ..., 0, 0, 0],
...,
[0, 0, 0, ..., 0, 0, 0],
[0, 0, 0, ..., 0, 0, 0],
[0, 0, 0, ..., 0, 0, 0]])

```
1 from sklearn.naive_bayes import GaussianNB
2 nv_model_bow = GaussianNB().fit(X_train_bow, y_train)
3 nv_model_tfidf = GaussianNB().fit(X_train_tfidf, y_train)
4
```

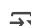
```
1 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
2
```

```
1 y_pred_bow = nv_model_bow.predict(X_test_bow)
2 y_pred_tfidf = nv_model_tfidf.predict(X_test_tfidf)
```

```
1 confusion_matrix(y_test, y_pred_bow)
```

 array([[499, 304],
[717, 880]])

```
1 print("BOW Accuracy", accuracy_score(y_test, y_pred_bow))
```

 BOW Accuracy 0.5745833333333333

```
1 confusion_matrix(y_test, y_pred_tfidf)
```

```
array([[488, 315],
       [696, 901]])
```

```
1 print("TFIDF Accuracy", accuracy_score(y_test, y_pred_tfidf))
```

```
TFIDF Accuracy 0.57875
```

```
1 ## Word2Vec
2 import gensim
```

Double-click (or enter) to edit

```
1 from nltk.tokenize import word_tokenize
2
3 # Tokenizing train and test data
4 X_train_tokens = [word_tokenize(sent.lower()) for sent in X_train]
5 X_test_tokens = [word_tokenize(sent.lower()) for sent in X_test]
```

```
1 from gensim.models import Word2Vec, KeyedVectors
```

```
1 import gensim.downloader as api
2 wv = api.load('word2vec-google-news-300')
```

```
[=====] 100.0% 1662.8/1662.8MB downloaded
```

```
1 from gensim.models import Word2Vec
2
3 # Train Word2Vec on training data only
4 w2v_model = Word2Vec(sentences=X_train_tokens, vector_size=100, window=5, min_count=1, workers=4)
5
```

```
1 import numpy as np
2 def sentence_vector(sentence, model, vector_size=100):
3     vectors = [model.wv[word] for word in sentence if word in model.wv]
4     return np.mean(vectors, axis=0) if vectors else np.zeros(vector_size)
5
6 # Convert train and test sets separately
7 X_train_vectors = np.array([sentence_vector(sent, w2v_model) for sent in X_train_tokens])
8 X_test_vectors = np.array([sentence_vector(sent, w2v_model) for sent in X_test_tokens])
9
```

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
3
```

```
1 # Train Random Forest Classifier
2 rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
3 rf_model.fit(X_train_vectors, y_train)
```

```
RandomForestClassifier
RandomForestClassifier(random_state=42)
```

```
1 # Predict on test data
2 y_pred = rf_model.predict(X_test_vectors)
3
```

```
1 # Evaluate model
2 accuracy = accuracy_score(y_test, y_pred)
3 print(f"Accuracy: {accuracy:.2f}")
```

```
Accuracy: 0.76
```

Generated code may be subject to a licence | djangotraining/Project

```
1 confusion_matrix(y_test, y_pred)
```

```
array([[ 411,  392],
       [ 192, 1405]])
```

```
1 classification_report(y_test, y_pred)
```

```
↵ '          precision    recall  f1-score   support\n\n  0.78         0.88         0.83        1597\n accuracy\n  0.71         2400\nweighted avg          0.75         0.76         0.75        2400\n'
```

1 Start coding or [generate](#) with AI.