

Analysis of goodreads.com

-Ankita Sawant

-Shradha Nayak

-Vandna Yadav

MOTIVATION

Encouraging authors to write books which will be liked
by goodread readers



GOAL OF OUR PROJECT

- Predicting if reviewers will like a book or not.



ASSUMPTION

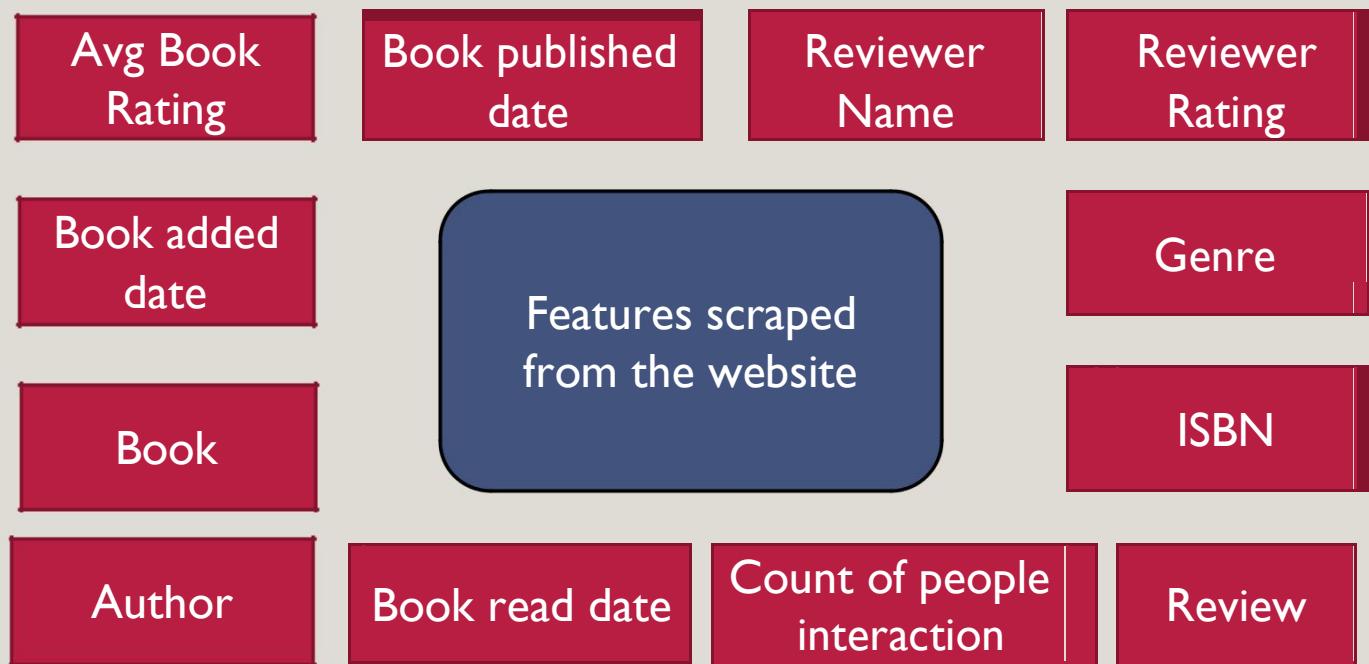
- Ratings given by reviewers for books, their genre, book authors, past history of the reviewers and the authors are related to them liking a particular book.



DATA COLLECTION



- The website that we have crawled for doing the prediction is <https://www.goodreads.com>



Features used for prediction



Genre

Reviewer
Ratings

Past Count
Genre

Past Rating
Genre

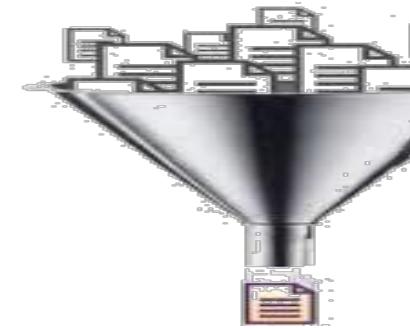
Book
Length

Past
Author
Book
count

TOOLS

- Python
- BeautifulSoup
- Chrome Web Driver
- Matplotlib and Pandas for plotting graphs





FILTERING AND CLEANING DATA

- We removed the books which were not in English language. (As no reviews were found)
- We removed the reviews which were not in English Language.
- We removed the records which did not have any reviewer rating
- We made the date format uniform

Book read date
Feb25,2015
Feb27,2016
Nov20,2014
Sep19,2015
Apr06,2014

WEBSITE USED FOR SCRAPING

goodreads Home My Books Browse ▾ Community ▾ Search books Search

Most popular 100 reviewers this week in The United States

People who wrote reviews that got the most votes on Goodreads this week (generated Jun 08, 2016 01:19PM)

this week | this month | in the last 12 months | all time

Rank	User Profile Picture	User Name	Location	Books Read	Friends	Votes This Week	Action
1.		karen	Woodside, NY (8 mi)	7,576 books	4,999 friends	961 votes	The United States The World Follow Reviews
2.		Jesse (JesseTheReader)	The United States	350 books	4,997 friends	949 votes	this week Follow Reviews
3.		Val⁺Shameless, Bitchy, Skanky, & Not Sorry⁺	ClayDanvers/RemingtonTate Manhattan Beach, CA (2455 mi)	3,604 books	700 friends	932 votes	this week Follow Reviews
4.		Candace	Edgewood, NM (1784 mi)	4,974 books	1,041 friends	810 votes	this week Follow Reviews
5.		Bill Kerwin	Columbus, OH (473 mi)	2,609 books	2,296 friends	729 votes	this week Follow Reviews
6.		Patrick	Freedom, WI (765 mi)	1,306 books	4,997 friends	690 votes	this week Follow Reviews
7.		Christy	Fairfield, OH (563 mi)	3,040 books	3,062 friends	669 votes	this week Follow Reviews
8.		Shelby *trains flying monkeys*	The United States	5,348 books	267 friends	651 votes	this week Follow Reviews

MEET PEOPLE

[top users](#)
[top readers](#)
[top reviewers](#)
[most popular reviewers](#)
[best reviews](#)
[online now](#)
[most followed](#)
[top librarians](#)
[recent statuses](#)
[recent reviews](#)

bookshelves
 all (7576)
 read (4055)
 currently-reading (3)
 to-read (1695)
 amd (470)
 ask-greg (288)
 at-my-desk (74)
 books-about-greg-s-mom (4)
 books-i-hate-more-than-most-other-t (1)
 buy-for-me-thanks (103)
 ceci-n-est-ce-pas-un-compte-rendu (18)
 continuing-ed-cookbooks (3)
 do-i-own-you (20)
 holy-grail-unicorn-tamerlane (90)
 i-am-the-one-percent (1)
 in-the-pipeline (15)
 nook-tbr (171)
 oh-dear (8)
 released-into-the-wild (127)
 romance-covers-i-have-loved (30)
 soon-to-buy (175)
 to-get-from-liberry (225)

 aaaahhrrrtt (128)
 americas-favorite-president (7)
 and-so-this-is-grad-school (63)
 animal-butts (3)
 appalachian-noir-southern-gothic (14)
 awkward-age-books-forgotten-til-now (113)
 babys-first-manga (15)

< previous 1 2 3 4 5 6 7 8 9 ... 378 379

cover	title	author	avg rating	rating	my rating	review ▾	date read	date added
	Divergent (Divergent, #1)	Roth, Veronica *	4.28	★★★★★	★★★★★	i need to make something perfectly clear. i am well aware that i gave 4 stars to Daughter of Smoke and Bone, and i am giving 5 stars to this one. the ...more	Jun 18, 2011	Jun 17, 2011
	Life and Death: Twilight Reimagined	Meyer, Stephenie	3.34	★★★★★	★★★★★	e.l. james runs out of ideas and desperately rewrites Fifty Shades of Grey from christian's perspective. stephenie meyers "celebrates" Twilight by swap ...more	not set	Oct 06, 2015
	Infinite Jest	Wallace, David Foster	4.32	★★★★★	★★★★★	this book... i think it is time to write a proper review for this book, as it is one of my all-time favorites and deserves way more than two words. bac ...more	not set	May 08, 2008
	Grey (Fifty Shades, #4)	James, E.L. *	3.75	★★★★★	★★★★★	my inner goddess is screaming out the safe word, but e.l. james is acting like we never set ground rules. rutabaga RUTABAGA!! RUTABAGA!!!!	not set	Jun 02, 2015

STATISTICS OF DATA SCRAPED

Attributes	Count
Total number of rows scraped	1760
After filtering(Removing rows with no rating)	1671
Original dataset number of features scraped	12
Final dataset number of features	57
Total number of reviewers scraped	98
Total number of reviews scraped per reviewer	Approx 20 (First page of each rev)

STATISTICS OF DATA SCRAPED

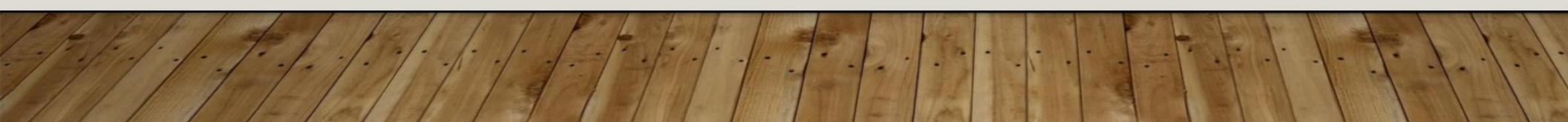
Groups	Features
Genre	Genre Type
Author	Author Name
Book	1)Book Name
	2)Book added date
	3)Book published date
	4)Book read date
	5)ISBN
Reviews	1)Review
	2)Reviewer Name
	3)Reviewer Rating
	4)Average book rating
	5)Count of people interactio with reviewer

Raw Features Processed Featu

Groups
Genre
Author
Book
Reviews

ORIGINAL DATASET

A	B	C	D	E	F	G	H	I	J
Author	Average book Rating	Book	Book added date	Book published date	Book read date	Count of people interaction	Genre	ISBN	Review
Fisher, Tarryn	4.18	F*ck Love	Jul08,2015	Dec31,2015	Jan-16	125	Romance	-	5 + STARS â€¢ 9
Kennedy, Elle	4.33	The Deal(Off-Campus, #1)	Feb24,2015	Feb24,2015	Feb25,2015	112	New Adult	-	5 WINNING S
Webster, K.	4.23	This is War, Baby(War & Peace, #1)	Feb26,2016	Feb29,2016	Feb27,2016	116	Dark	-	5 + I WAR S
Harmon, Amy	4.39	The Law of Moses(The Law of Moses, #1)	Sep24,2014	Nov18,2014	Nov20,2014	119	Romance	-	5 *Greats* 9
Cherry, Brittainy C.	4.29	The Air He Breathes(Elements, #1)	Jul31,2015	Sep25,2015	Sep19,2015	102	Romance	-	5 + STARS "H
Fisher, Tarryn	4.09	Mud Vein	Jul28,2013	Mar08,2014	Apr06,2014	110	Dark	-	The writing i:
Zapata, Mariana	4.36	The Wall of Winnipeg and Me	Dec26,2015	Feb29,2016	Mar-16	112	Romance	-	5 "TEAM GR/9
Williams, Nicole	4.08	Collared	Feb19,2016	Mar22,2016	Mar03,2016	118	Romance	-	5 + STARS â€¢ 9
Cole, Tillie	4.39	A Thousand Boy Kisses	Jul16,2015	Mar15,2016	Mar15,2016	102	Young Adult	-	This one was 9
Walters, A. Meredith	4.11	One Day Soon(One Day Soon, #1)	Nov04,2015	Feb18,2016	Feb15,2016	99	Romance	-	5 + STARS!â€¢ 9
Hoover, Colleen	4.27	Confess	Jan29,2015	Mar10,2015	Jan 31,2015	84	Romance	-	3.75 ~ 4 OM(9
Holden, Kim	4.25	So Much More	Feb27,2016	Mar29,2016	Mar29,2016	113	Romance	-	5 â€œWEâ€¢ 9
Torre, Alessandra	4.14	Hollywood Dirt(Hollywood Dirt, #1)	Jun27,2015	Sep07,2015	Sep03,2015	89	Romance	-	4.5 Coca Cola 9
Frazier, T.M.	4.25	King(King, #1)	Jan22,2014	Jun15,2015	Jun15,2015	101	Dark	-	I have been 9
Ward, Penelope	4.03	Sins of Sevin	Mar23,2015	Sep14,2015	Sep21,2015	112	Romance	-	2.75 Stars "H
Martinez, Aly	4.1	The Fall Up(The Fall Up, #1)	Aug12,2015	Oct26,2015	Oct13,2015	95	Romance	1518711391	5 "And falling 9
Lilley, R.K.	4.28	Breaking Him(Love is War, #1)	Aug25,2015	Oct03,2015	Oct14,2015	84	Romance	-	4.5 STARS â€¢ 9
Harmon, Amy	4.4	Making Faces	Oct19,2013	Oct12,2013	Oct21,2013	69	Romance	-	This isnâ€™t 9
Sheridan, Mia	3.76	Midnight Lily	Mar02,2016	Mar01,2016	Mar02,2016	101	Romance	-	3.5 StarsSom 9
Shen, L.J.	4.14	Sparrow	Feb27,2016	Mar01,2016	Mar09,2016	98	Romance	-	4.5 Stars â€¢ 9
Moriarty, Liane	4.16	Big Little Lies	May28,2014	Jul29,2014	Jul20,2014	20	Fiction	399167064	This one was 9
Kalanithi, Paul	4.37	When Breath Becomes Air	Jan14,2016	Jan12,2016	Jan31,2016	43	Nonfiction	-	As I finished 9
Doerr, Anthony	4.31	All the Light We Cannot See	Dec11,2013	May-14	Apr03,2014	66	Historical	1476746583	For me, this 9



FINAL DATASET USED FOR PREDICTION

AV	AW	AX	AY	AZ	BA	BB	BC	BD
War	Western	Young Adult	Length_Of_Review	Reviewer_Ratings	Past_Count_Genre	Past_Rating_Genre	Book_Length	Past_Author_Book_Count
0	0	0	150	1	0	0	28	
0	0	0	150	1	1	1	27	
0	0	0	150	3	2	1	45	
0	0	0	152	1	3	1.67	36	
0	0	0	150	3	4	1.5	23	
0	0	0	150	3	5	1.8	21	
0	0	0	150	3	6	2	4	
0	0	0	154	3	7	2.14	15	
0	0	0	150	2	8	2.25	15	
0	0	0	156	3	9	2.22	24	
0	0	0	150	3	10	2.3	10	
0	0	0	150	3	11	2.36	14	
0	0	0	150	3	12	2.42	45	
0	0	0	150	1	13	2.46	55	
0	0	0	150	3	14	2.36	13	
0	0	0	150	3	15	2.4	18	
0	0	0	152	3	0	0	49	
0	0	0	150	3	0	0	16	
0	0	0	150	3	1	3	16	
0	0	0	150	3	0	0	6	
0	0	0	149	1	1	3	36	
0	0	0	150	1	2	2	36	
0	0	0	152	3	3	1.67	12	
0	0	0	150	1	4	2	17	
0	0	0	150	3	5	1.8	38	

FEATURE TYPE

- Author → Text
- Avg book Rating → Float
- Book → Alpha Numeric
- Book added date → Date
- Book published date → Date
- Book read date → Date
- Count of people interaction → Numeric
- ISBN → Alpha Numeric
- Review → Text
- Reviewer Name → Text

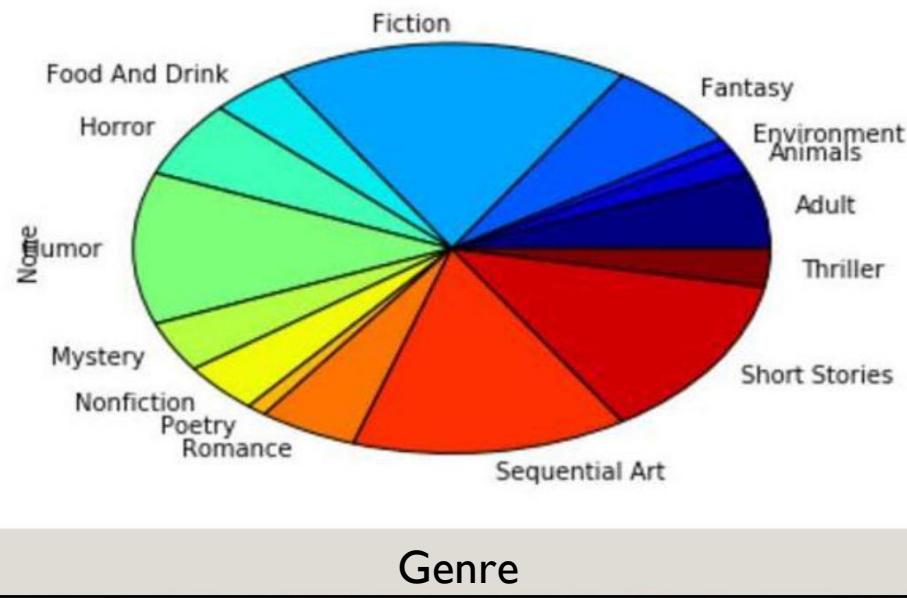
- Reviewer Rating → Numeric
- Genre → Categorical
- Past Count of Genre → Numeric
- Past Rating of Genre → Float
- Length Of Review → Numeric
- Book Length → Numeric
- Past Author Book Count → Numeric
- Past Author Book Rating → Float

DISTRIBUTION OF FEATURES

Distribution of Genre according to mid term data set

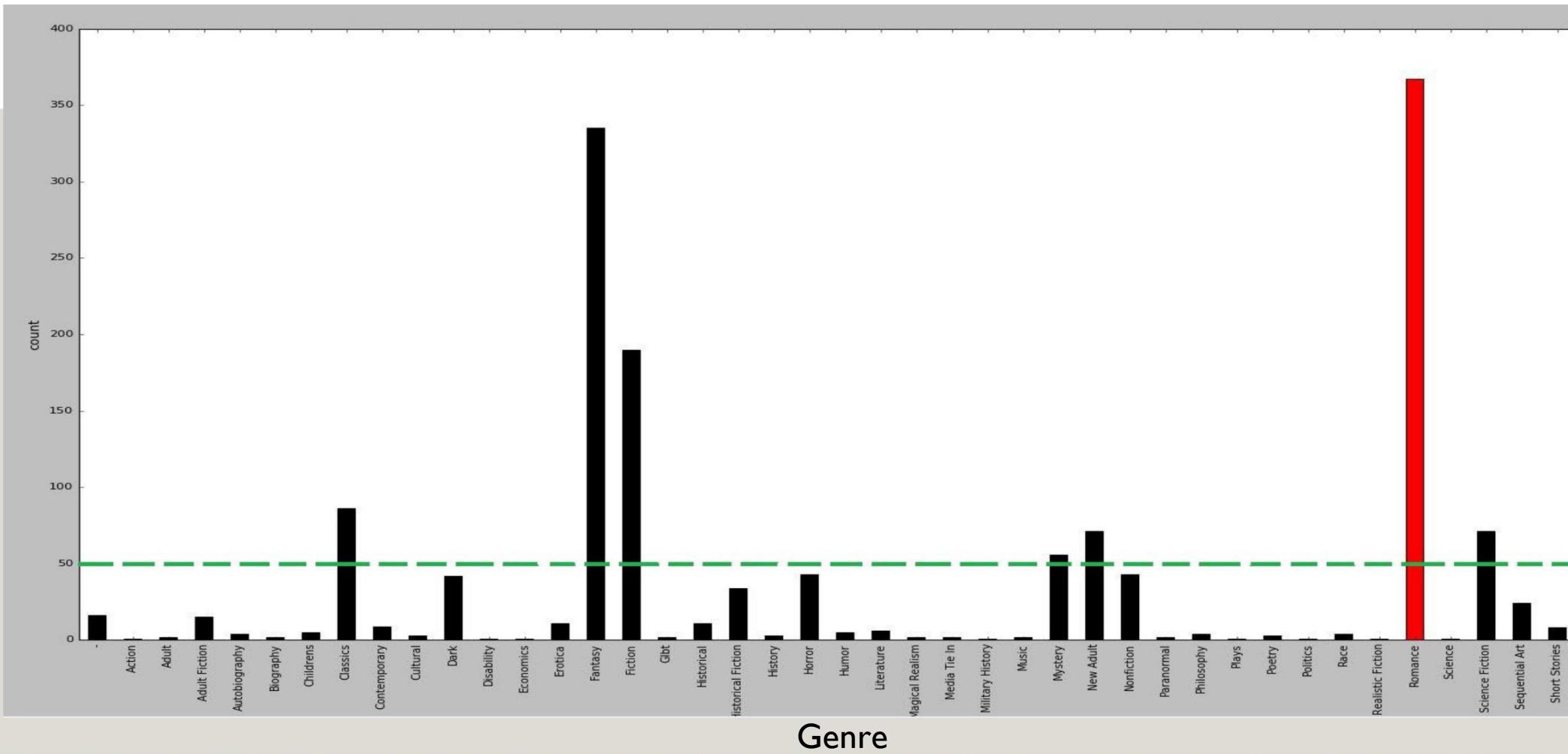
```
data.groupby('Genre').size().plot(kind='pie', colormap='jet' )
```

```
<matplotlib.axes._subplots.AxesSubplot at 0xe683cc0>
```

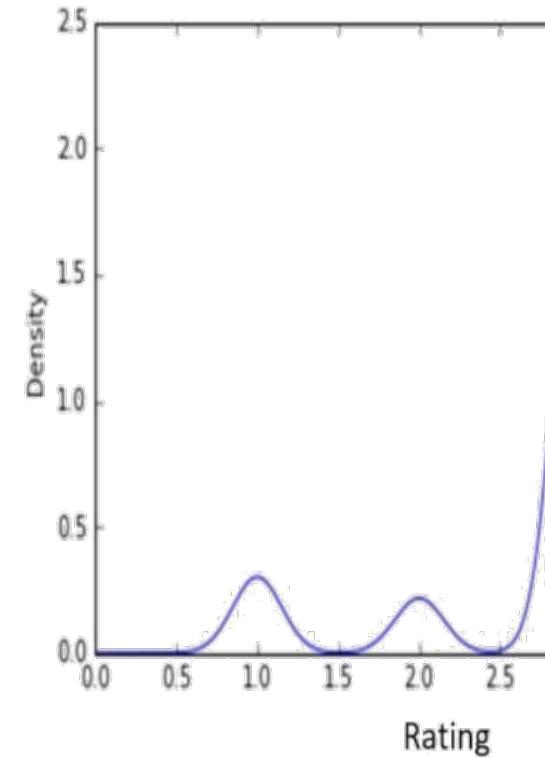
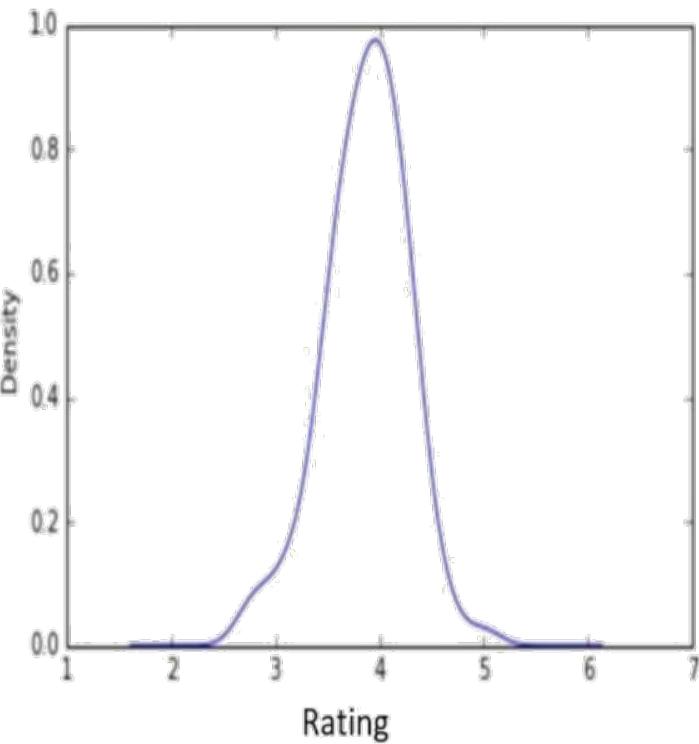


Distribution of Genre according to final data set

- The Genres: Romance, Fantasy and Fiction are the most prominent in our data set.

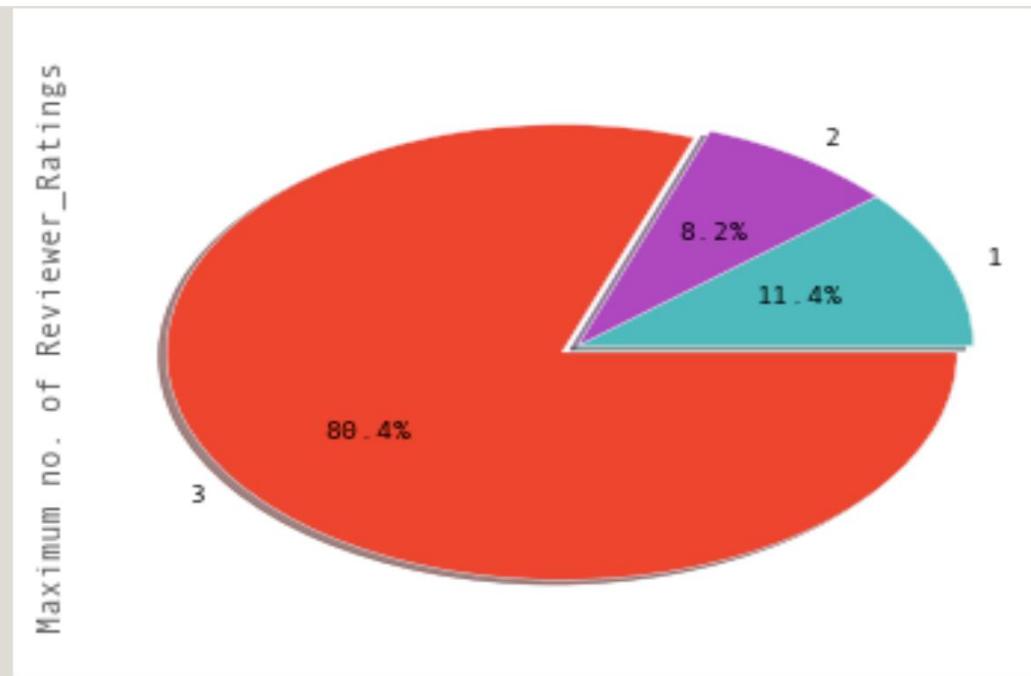


Reviewer Rating distribution and conversion



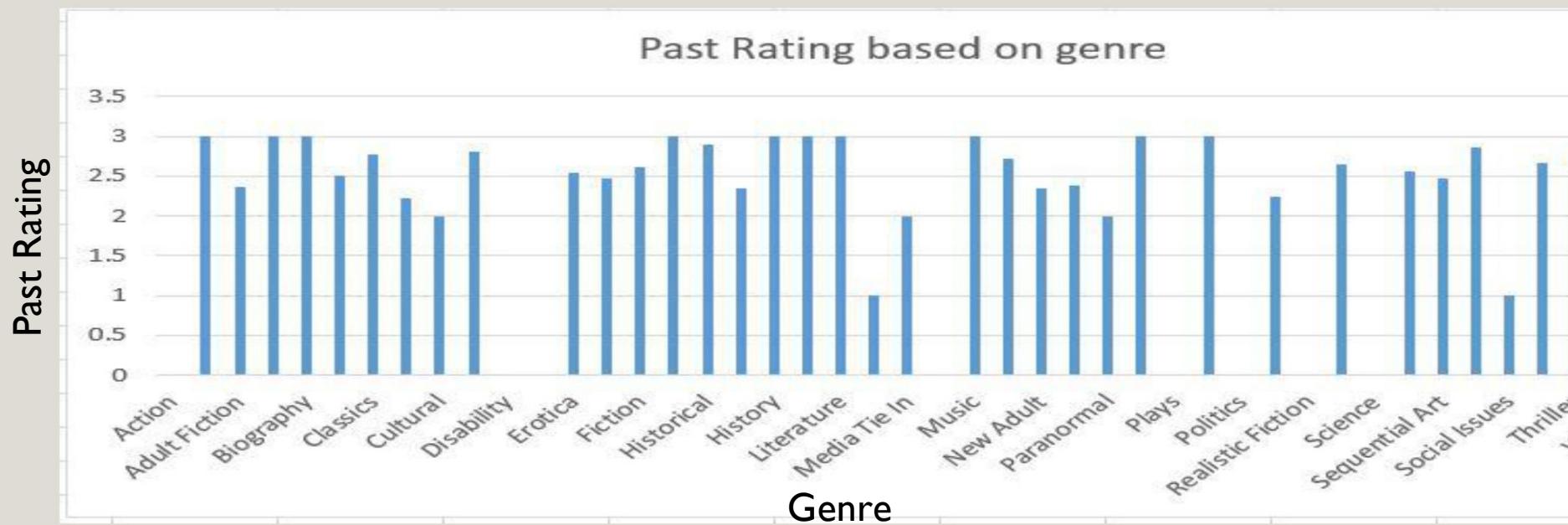
- The Reviewer Ratings were in the range 1 to 5. We converted them to 1,2,3.

Distribution of Reviewer Ratings



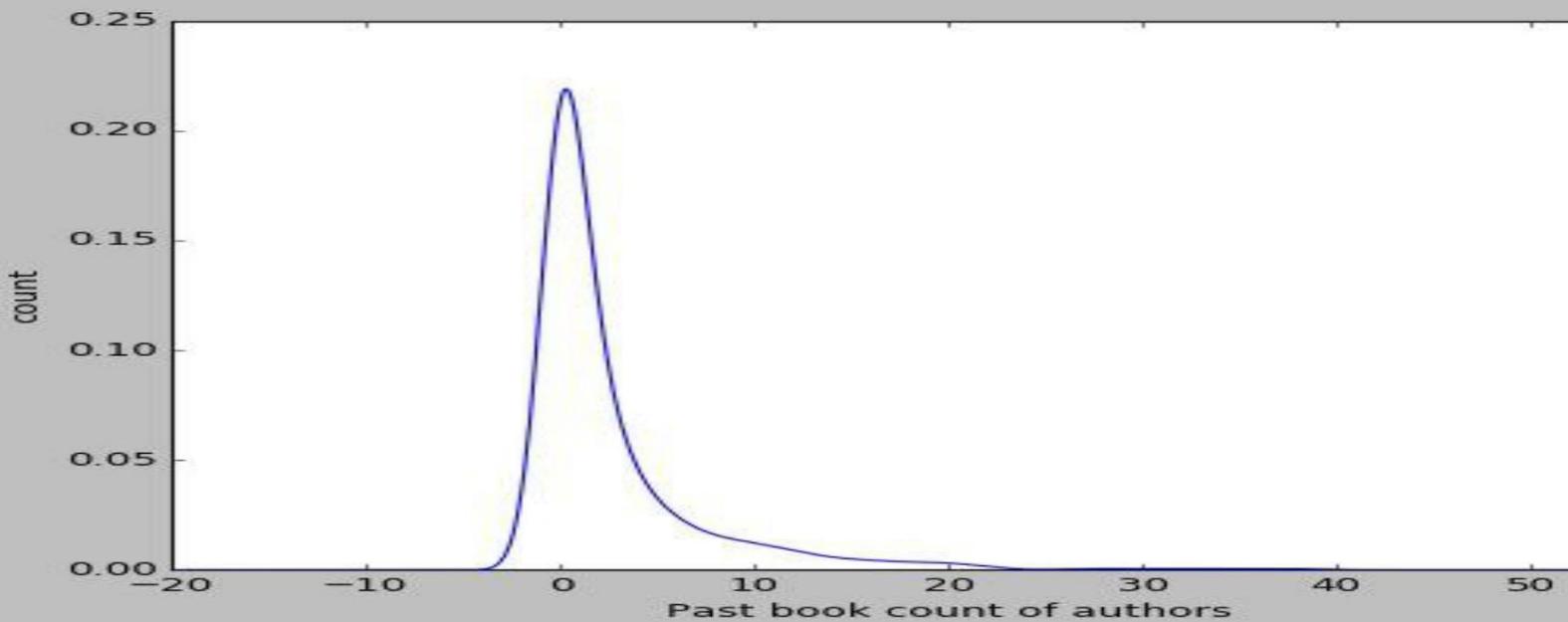
- Majority of the reviewer ratings were 3.

Distribution of Past Rating according to Genre



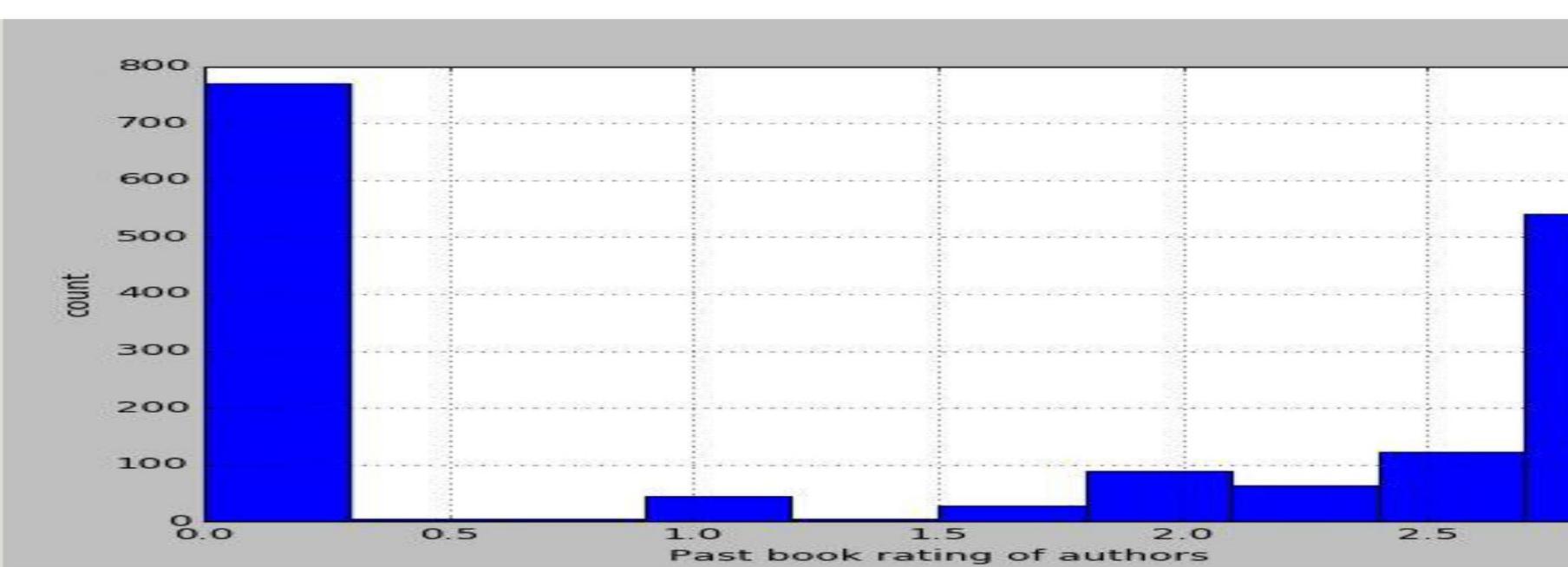
- We calculated the past rating for each record based on the ratings of the books read before and belonging to the same genre.

Distribution of Past book count of authors



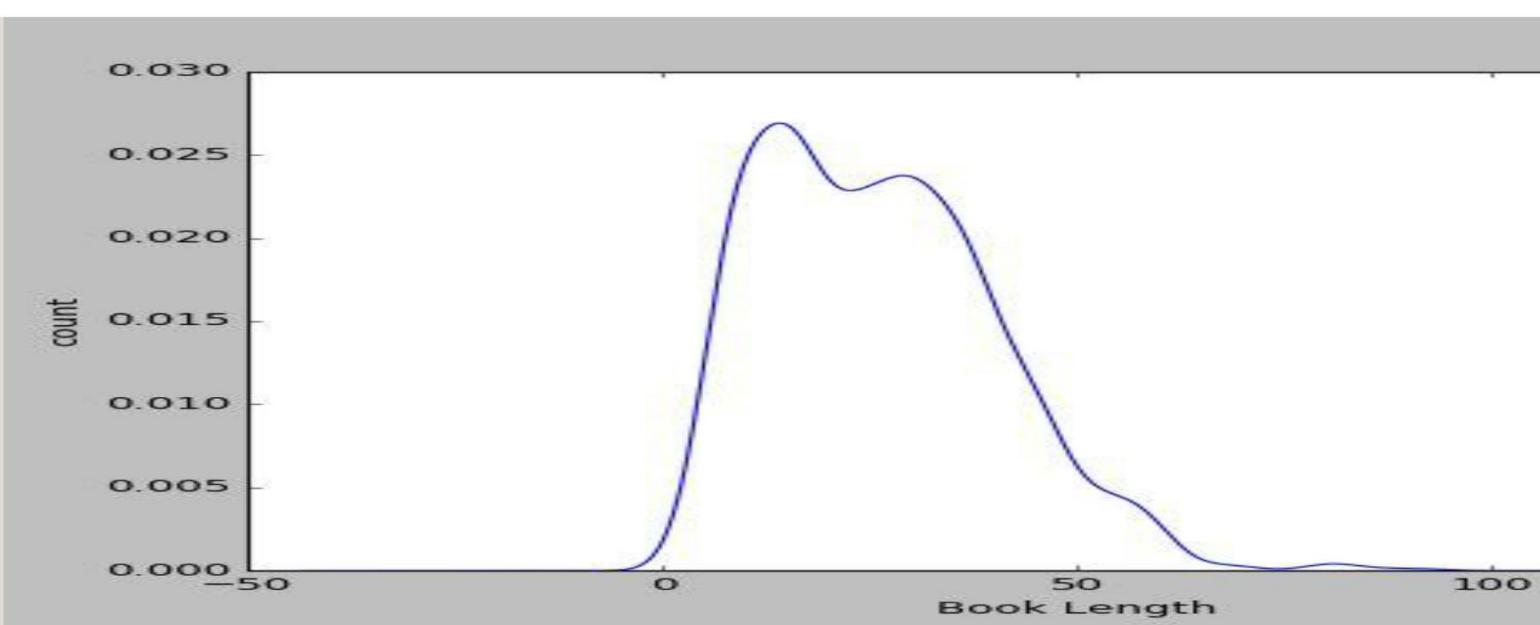
- We calculated the past book count for each record based on the books read before that and belonging to the same author.

Distribution of Past book rating of authors



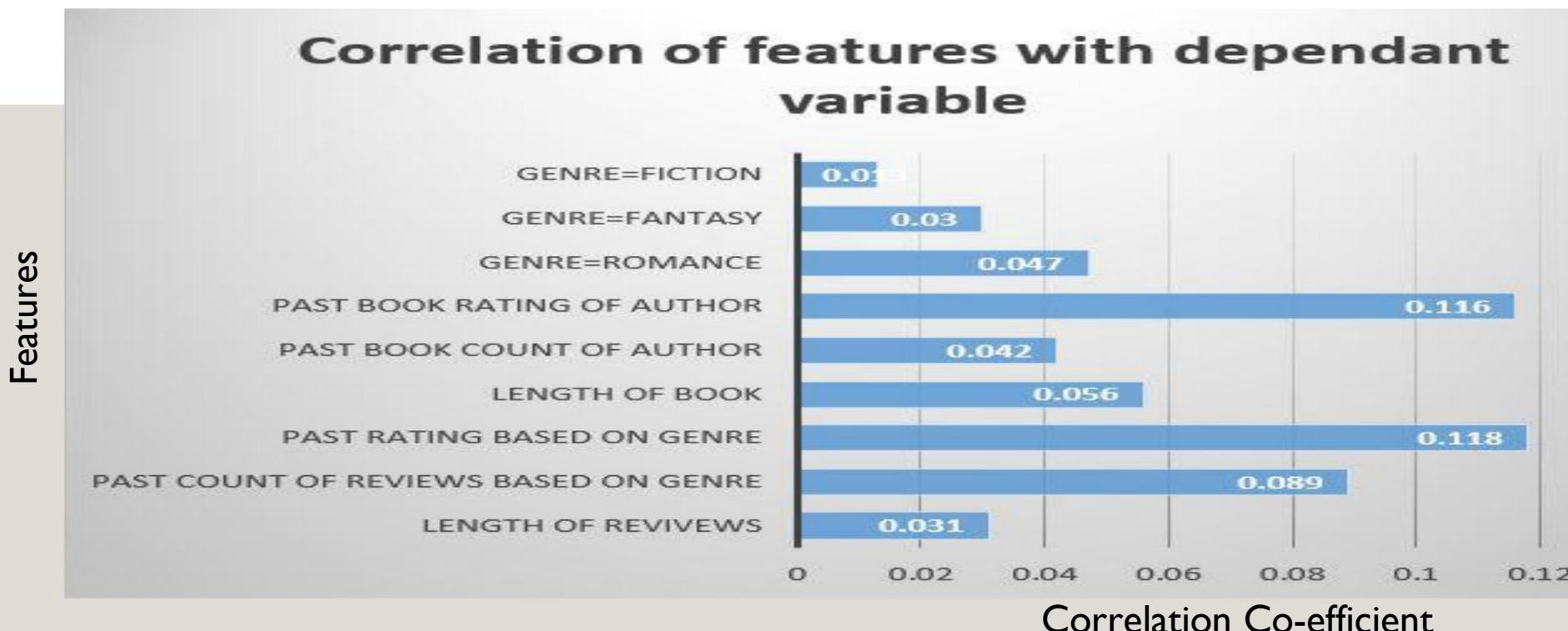
- We calculated the past rating for each record based on the ratings of the books read before book and belonging to the same author. Most of the past book ratings were 0 as the authors not repeat much in our dataset.

Distribution of Book Title Length



- We calculated the book title length for each record. We considered this feature on the assumption that sequels would have the same book title length.

FEATURES CORRELATION



- The 'Past book rating of author' and 'Past rating based on genre' had a high correlation with our dependent variable 'Reviewer Ratings'.

METHODS TO ACHIEVE OUR GOAL

After website scraping and data cleaning we used the following:

- KNN
- Logistic regression
- Naive Bayes
- Random Forest

We also used Voting Classifier on our 3 best algorithms.



PRE-PROCESSING

Making the dataset random

```
df1=csv.reindex(np.random.permutation(csv.index))
```

Training and Test data set

```
df1
count=len(csv)
trainCount=int((70*count)/100)
testCount=(count-trainCount)
rev_train=df1[0:trainCount]
rev_test=df1[trainCount+1:count]
```

Normalisation

```
import pandas as pd
def norm(x):
    print x
    df_max_aa=x.max()
    df_mean_aa=x.mean()
    df_min_aa=x.min()
    x=(x-df_min_aa)/(df_max_aa)
    return x
```

ACCURACY OF MODELS

KNN

```
KNN=KNeighborsClassifier(5)
KNN.fit(rev_train_Norm.astype(int),labels_train.astype(int))
predicted=KNN.predict(rev_test_Norm)
print accuracy_score(predicted,labels_test)
```

0.802395209581

```
{'n_neighbors': 1, 'weights': 'uniform'} 0.514970
{'n_neighbors': 1, 'weights': 'distance'} 0.514970
{'n_neighbors': 3, 'weights': 'uniform'} 0.526946
{'n_neighbors': 3, 'weights': 'distance'} 0.526946
{'n_neighbors': 5, 'weights': 'uniform'} 0.800684
{'n_neighbors': 5, 'weights': 'distance'} 0.800684
{'n_neighbors': 7, 'weights': 'uniform'} 0.800684
{'n_neighbors': 7, 'weights': 'distance'} 0.800684
{'n_neighbors': 9, 'weights': 'uniform'} 0.800684
{'n_neighbors': 9, 'weights': 'distance'} 0.800684
{'n_neighbors': 11, 'weights': 'uniform'} 0.800684
{'n_neighbors': 11, 'weights': 'distance'} 0.800684
{'n_neighbors': 13, 'weights': 'uniform'} 0.800684
{'n_neighbors': 13, 'weights': 'distance'} 0.800684
{'n_neighbors': 15, 'weights': 'uniform'} 0.800684
{'n_neighbors': 15, 'weights': 'distance'} 0.800684
{'n_neighbors': 17, 'weights': 'uniform'} 0.800684
{'n_neighbors': 17, 'weights': 'distance'} 0.800684
{'n_neighbors': 19, 'weights': 'uniform'} 0.800684
{'n_neighbors': 19, 'weights': 'distance'} 0.800684
{'n_neighbors': 21, 'weights': 'uniform'} 0.800684
{'n_neighbors': 21, 'weights': 'distance'} 0.800684
{'n_neighbors': 23, 'weights': 'uniform'} 0.800684
{'n_neighbors': 23, 'weights': 'distance'} 0.800684
{'n_neighbors': 25, 'weights': 'uniform'} 0.800684
{'n_neighbors': 25, 'weights': 'distance'} 0.800684
{'n_neighbors': 27, 'weights': 'uniform'} 0.800684
{'n_neighbors': 27, 'weights': 'distance'} 0.800684
{'n_neighbors': 29, 'weights': 'uniform'} 0.800684
{'n_neighbors': 29, 'weights': 'distance'} 0.800684
{'n_neighbors': 31, 'weights': 'uniform'} 0.800684
{'n_neighbors': 31, 'weights': 'distance'} 0.800684
```

Best parameters {'n_neighbors': 5, 'weights': 'uniform'} 0.812375249501

Logistic Regression

```
LR=LogisticRegression ()  
LR.fit(rev_train,labels_train)  
#use the classifier to predict  
predicted=LR.predict(rev_test)  
#print the accuracy  
print accuracy_score(predicted,labels_test)
```

0.818363273453

Naïve Bayes

```
rev_train_Norm = sklearn.preprocessing.normalize(rev_train, norm='l2')  
rev_test_Norm=sklearn.preprocessing.normalize(rev_test, norm='l2')  
clf3 = MultinomialNB()  
clf3.fit(rev_train_Norm,labels_train)  
pred=clf3.predict(rev_test_Norm)  
print accuracy_score(pred,labels_test)
```

0.796407185629

Random F

```
clf3 = RandomForestClassi  
clf3.fit(rev_train_Norm,  
pred=clf3.predict(rev_te  
print accuracy_score(prec
```

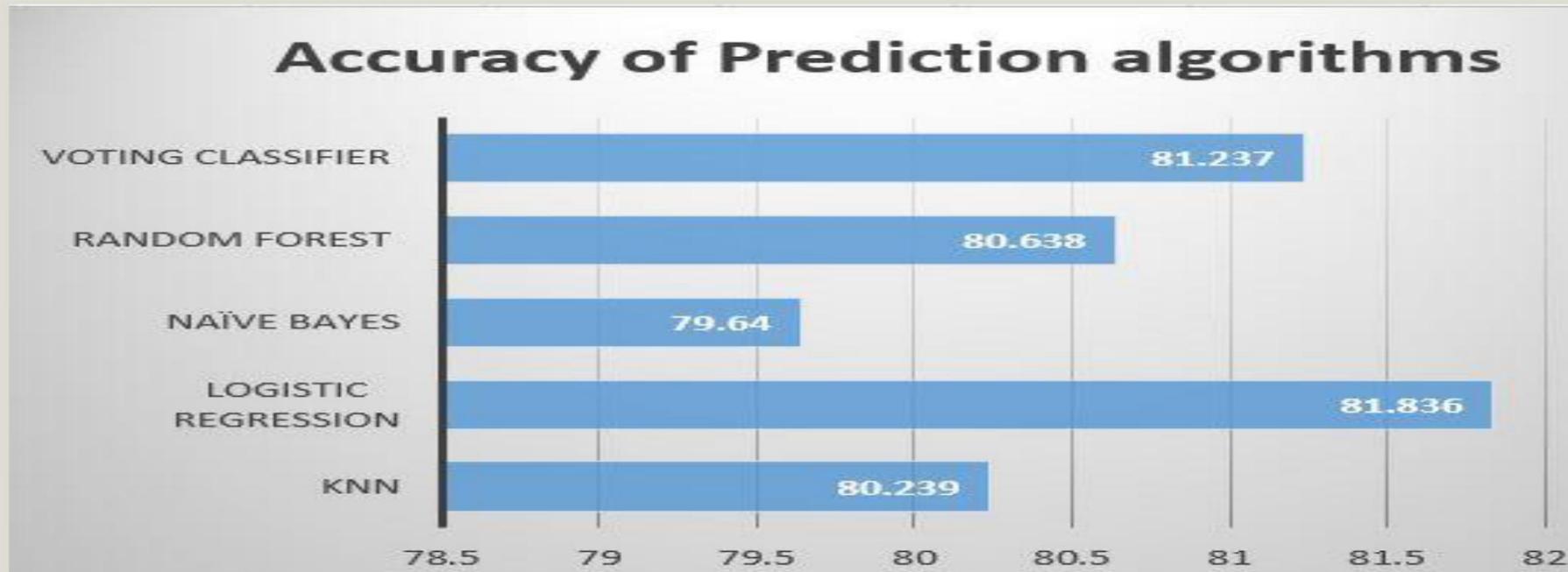
0.806387225549

Voting Classifier

```
clf1 = LogisticRegression()
clf2 = KNeighborsClassifier(5)
clf3 = RandomForestClassifier(n_estimators=100)
eclf = VotingClassifier(estimators=[('lr', clf1), ('knn', clf2), ('rf',clf3)], vo
eclf.fit(rev_train_Norm,labels_train)
pred=eclf.predict(rev_test_Norm)
print accuracy_score(pred,labels_test)
```

0.812375249501

SUMMARY



- The 'Logistic Regression' model proved to be the best for our prediction.

DETAILED SUMMARY

With length of review				
Algos	Accuracy without normalisation and reviewer ratings=1,2,3	Accuracy with normalisation & reviewer ratings =1,2,3	Accuracy without normalisation and reviewer ratings=1,2,3,4,5	Accuracy with & reviewer rating
KNN	77.45	80.39	45.5	
Logistic Regression	79.44	80.83	55.28	
NB	32.5	79.24	17.77	
Random Forest	78.04	80.03	52.89	
Combination of 3(KNN,Random Forest,Logistic)	78.64	81.83	51.49	
Without length of review				
Algos	Accuracy without normalisation and reviewer ratings=1,2,3	Accuracy with normalisation & reviewer ratings =1,2,3	Accuracy without normalisation and reviewer ratings=1,2,3,4,5	Accuracy with & reviewer rating
KNN	76.44	81.43	46.03	
Logistic Regression	79.44	80.63	56.8	
NB	54.08	81.63	33.93	
Random Forest	78.84	79.24	49.9	
Combination of 3(KNN,Random Forest,Logistic)	80.83	80.638	52.09	

- Length of review feature will not be used for the prediction.

