

```
In [1]: 1 %matplotlib inline
        2
        3 import sqlite3
        4 import pandas as pd
        5 import numpy as np
        6 import nltk
        7 import string
        8 import matplotlib.pyplot as plt
        9 import seaborn as sns
       10 from sklearn.feature_extraction.text import TfidfTransformer
       11 from sklearn.feature_extraction.text import TfidfVectorizer
       12
       13 from sklearn.feature_extraction.text import CountVectorizer
       14 from sklearn.metrics import confusion_matrix
       15 from sklearn import metrics
       16 from sklearn.metrics import roc_curve, auc
       17 from nltk.stem.porter import PorterStemmer
```

```
In [2]: 1 import json
        2 import pandas as pd
        3 final = pd.read_json('reviews_Movies_and_TV_5.json.gz', lines=True)
```

In [3]:

1 final

Out[3]:

	reviewerID	asin	reviewerName	helpful	reviewText	overall	summary	unixReviewTime	reviewTime
0	ADZPIG9QOCDG5	0005019281	Alice L. Larson "alice-loves-books"	[0, 0]	This is a charming version of the classic Dick...	4	good version of a classic	1203984000	02 26, 2008
1	A35947ZP82G7JH	0005019281	Amarah Strack	[0, 0]	It was good but not as emotionally moving as t...	3	Good but not as moving	1388361600	12 30, 2013
2	A3UORV8A9D5L2E	0005019281	Amazon Customer	[0, 0]	Don't get me wrong, Winkler is a wonderful cha...	3	Winkler's Performance was ok at best!	1388361600	12 30, 2013
3	A1VKW06X1O2X7V	0005019281	Amazon Customer "Softmill"	[0, 0]	Henry Winkler is very good in this twist on th...	5	It's an enjoyable twist on the classic story	1202860800	02 13, 2008
4	A3R27T4HADWFFJ	0005019281	BABE	[0, 0]	This is one of the best Scrooge movies out. H...	4	Best Scrooge yet	1387670400	12 22, 2013
...
1697528	AV657BUYHHXZ2	B00LT1JHLW	Mike Rules "Mike"	[1, 14]	wow \$269.99 for the entire series on Blu Ray??...	1	Way to Expensive!! WB = GREED	1406073600	07 23, 2014
1697529	A17W587EH23J0Q	B00LT1JHLW	Ron2900 "Ron"	[32, 48]	Finally, the holy grail of tv-on-dvd boxsets i...	5	HOLY BAT-BOXSET, BATMAN... I never thought thi...	1405641600	07 18, 2014
1697530	A3DE438TF1A958	B00LT1JHLW	thomas henry	[3, 10]	Could this be a true or I'm i dreaming batman ...	5	prayers have been answered because batman 60s ...	1405728000	07 19, 2014
1697531	A2RWCXDMANY0LW	B00LT1JHLW	wheev	[0, 4]	I've been a fan of the series since I was a yo...	5	can't Wait!	1405987200	07 22, 2014
1697532	A3ROPC55BE2OM9	B00LT1JHLW	WingLT	[11, 23]	People seriously need to wake up and realize t...	5	The Price is Insane? People Really Need to Wak...	1405728000	07 19, 2014

1697533 rows × 9 columns

In [4]:

```
1 # find sentences containing HTML tags
2 import re
3 i=0;
4 for sent in final['reviewText'].values:
5     if (len(re.findall('<.*?>', sent))):
6         print(i)
7         print(sent)
8         break;
9     i += 1;
```

In []:

1 final

Out[5]:

	reviewerID	asin	reviewerName	helpful	reviewText	overall	summary	unixReviewTime	reviewTime
0	ADZPIG9QOCDG5	0005019281	Alice L. Larson "alice-loves-books"	[0, 0]	This is a charming version of the classic Dick...	4	good version of a classic	1203984000	02 26, 2008
1	A35947ZP82G7JH	0005019281	Amarah Strack	[0, 0]	It was good but not as emotionally moving as t...	3	Good but not as moving	1388361600	12 30, 2013
2	A3UORV8A9D5L2E	0005019281	Amazon Customer	[0, 0]	Don't get me wrong, Winkler is a wonderful cha...	3	Winkler's Performance was ok at best!	1388361600	12 30, 2013
3	A1VKW06X1O2X7V	0005019281	Amazon Customer "Softmill"	[0, 0]	Henry Winkler is very good in this twist on th...	5	It's an enjoyable twist on the classic story	1202860800	02 13, 2008
4	A3R27T4HADWFFJ	0005019281	BABE	[0, 0]	This is one of the best Scrooge movies out. H...	4	Best Scrooge yet	1387670400	12 22, 2013
...
1697528	AV657BUYHHXZ2	B00LT1JHLW	Mike Rules "Mike"	[1, 14]	wow \$269.99 for the entire series on Blu Ray??...	1	Way to Expensive!! WB = GREED	1406073600	07 23, 2014
1697529	A17W587EH23J0Q	B00LT1JHLW	Ron2900 "Ron"	[32, 48]	Finally, the holy grail of tv-on-dvd boxsets i...	5	HOLY BAT-BOXSET, BATMAN... I never thought thi...	1405641600	07 18, 2014
1697530	A3DE438TF1A958	B00LT1JHLW	thomas henry	[3, 10]	Could this be a true or I'm i dreaming batman ...	5	prayers have been answered because batman 60s ...	1405728000	07 19, 2014
1697531	A2RWCXDMANY0LW	B00LT1JHLW	wheev	[0, 4]	I've been a fan of the series since I was a yo...	5	can't Wait!	1405987200	07 22, 2014
1697532	A3ROPC55BE2OM9	B00LT1JHLW	WingLT	[11, 23]	People seriously need to wake up and realize t...	5	The Price is Insane? People Really Need to Wak...	1405728000	07 19, 2014

1697533 rows × 9 columns

```
In [ ]: 1 df = final[~final['reviewText'].isnull()]
```

```
In [ ]: 1 from textblob import TextBlob
2 def preprocess(ReviewText):
3     ReviewText = ReviewText.str.replace("<br/>", "")
4     ReviewText = ReviewText.str.replace('<a.*>.*</a>', '')
5     ReviewText = ReviewText.str.replace('&', '')
6     ReviewText = ReviewText.str.replace('>', '')
7     ReviewText = ReviewText.str.replace('<', '')
8     ReviewText = ReviewText.str.replace('\xa0', ' ')
9     return ReviewText
10 df['reviewText'] = preprocess(df['reviewText'])
11
12 df['polarity'] = df['reviewText'].map(lambda text: TextBlob(text).sentiment.polarity)
13 df['review_len'] = df['reviewText'].astype(str).apply(len)
14 df['word_count'] = df['reviewText'].apply(lambda x: len(str(x).split()))
```

```
In [ ]: 1 print('5 random reviews with the highest positive sentiment polarity: \n')
2 cl = df.loc[df.polarity == 1, ['reviewText']].sample(5).values
3 for c in cl:
4     print(c[0])
```

5 random reviews with the highest positive sentiment polarity:

Dress up a video with 27 misses and you have one hit, including the 28th dress sung to the tune of "Here Comes the Bride." Delightful entertainment!

This man should be posthumously pardoned by the State of Montana. Excellent film. Forever one of the finest westerns made.

Every time I think about trading this DVD in, I put it on and am reminded of why UNDERWORLD has been the best show I've ever seen in my life.

Excellent music. Product was received in timely manner, and as advertised.

Lillies of the Field is one of the best movies by Sidney Poitier!! A thumbs up to ALL the nuns who played in it!!

```
In [ ]: 1 print('5 random reviews with the most neutral sentiment(zero) polarity: \n')
        2 cl = df.loc[df.polarity == 0, ['reviewText']].sample(5).values
        3 for c in cl:
        4     print(c[0])
```

5 random reviews with the most neutral sentiment(zero) polarity:

I had several of these horror flicks marked to watch and deleted them all from my watch list. Not much of a story. Bungled even by the likes of Gene Hackman, Michael Rosenbaum at last defines Lex Luthor. Kristen becomes a star. 10-year hit status assured.

This movie reminds me of the movie *Colors*, only updated to present day. There is a bit of *Adam 12*, *Colors*, *Training Day*, all mixed together and updated. Well done, and understated with the handheld camera technique. Keeps you glued to the screen, all the while developing the relationship between the two LAPD cops on patrol. Definitely a must-see.

This was not a movie for me. I actually could not finish watching it, just did not keep my interest.

This movie scared me to death 35 yrs ago and is still one of my favorites. Chilling is a word thrown around a lot in descriptions of movies of this genre but in this case it fits. It doesn't get any better than this for sci-fi/horror of the 50's or any decade.

```
In [ ]: 1 print('5 reviews with the most negative polarity: \n')
        2 cl = df.loc[df.polarity == -0.5, ['reviewText']].sample(5).values
        3 for c in cl:
        4     print(c[0])
```

5 reviews with the most negative polarity:

In this movie I thought the friend should have realized sooner what he did to his friends it was awful a person had to lose their life. But I would recommend this movie cause there is a lesson to be learned.

Not sure why people like this movie. This movie is not working. Just a sequence of scenes between dream / reality. Looks like the director didn't know what to do. I got bored and disappointed.

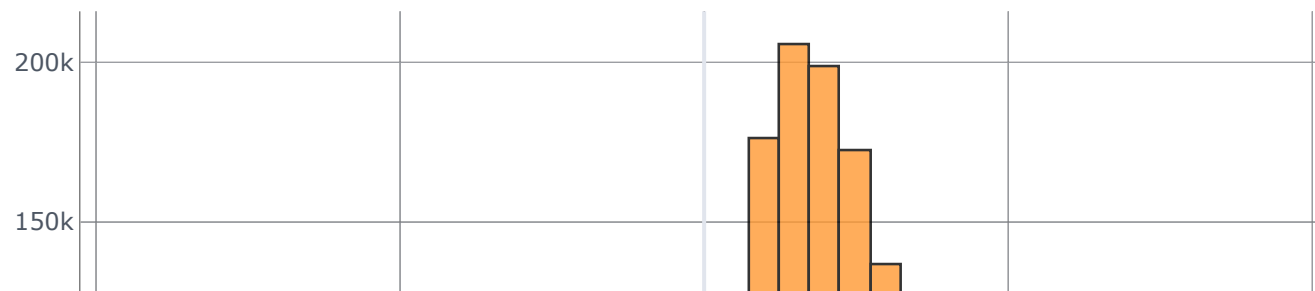
once again Hollywood fails to deliver on the sequel; why do they bother.

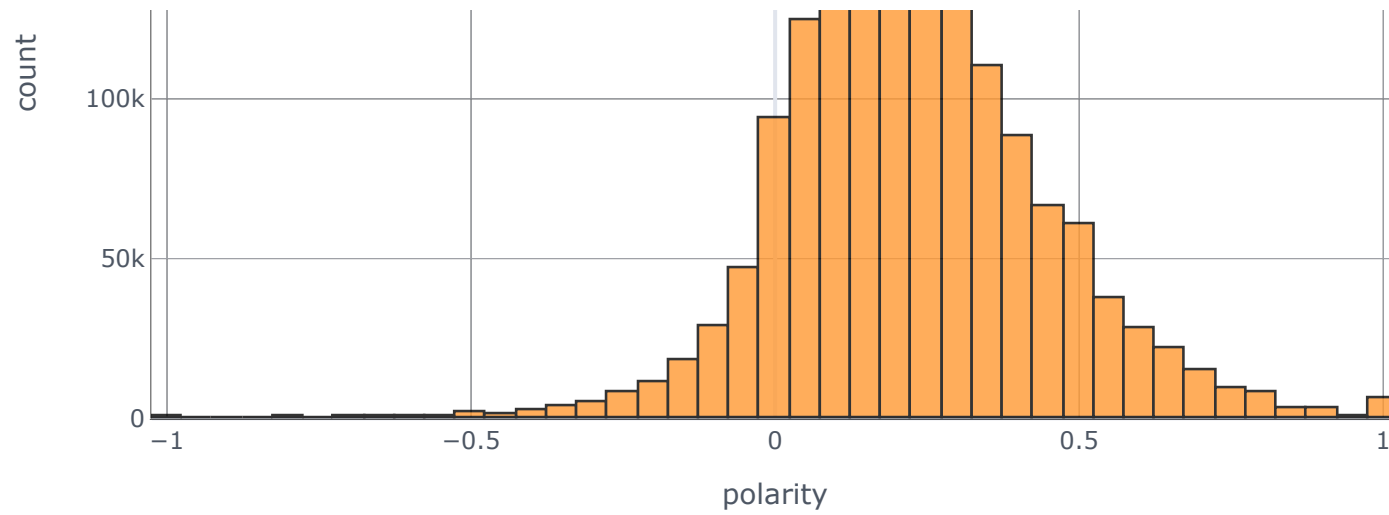
I tried watching it but just was not into this dvd I was so bored with this movie I just did not like it.

it's as scary as it looks....Very intense. And the clown is just horrifying. I wouldn't let no kids watch this. it scared me.

```
In [ ]: 1 #Import Libraries
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6
7 #Plotly Tools
8 from plotly.offline import init_notebook_mode, iplot
9 init_notebook_mode(connected=True)
10 import plotly.graph_objs as go
11 import plotly.offline as offline
12 offline.init_notebook_mode()
13 from plotly import tools
14 import plotly.tools as tls
15 init_notebook_mode(connected=True)
16 import cufflinks as cf
17 cf.go_offline()
18 cf.set_config_file(offline=False, world_readable=True)
19
20 df['polarity'].iplot(
21     kind='hist',
22     bins=50,
23     xTitle='polarity',
24     linecolor='black',
25     yTitle='count',
26     title='Sentiment Polarity Distribution')
```

Sentiment Polarity Distribution

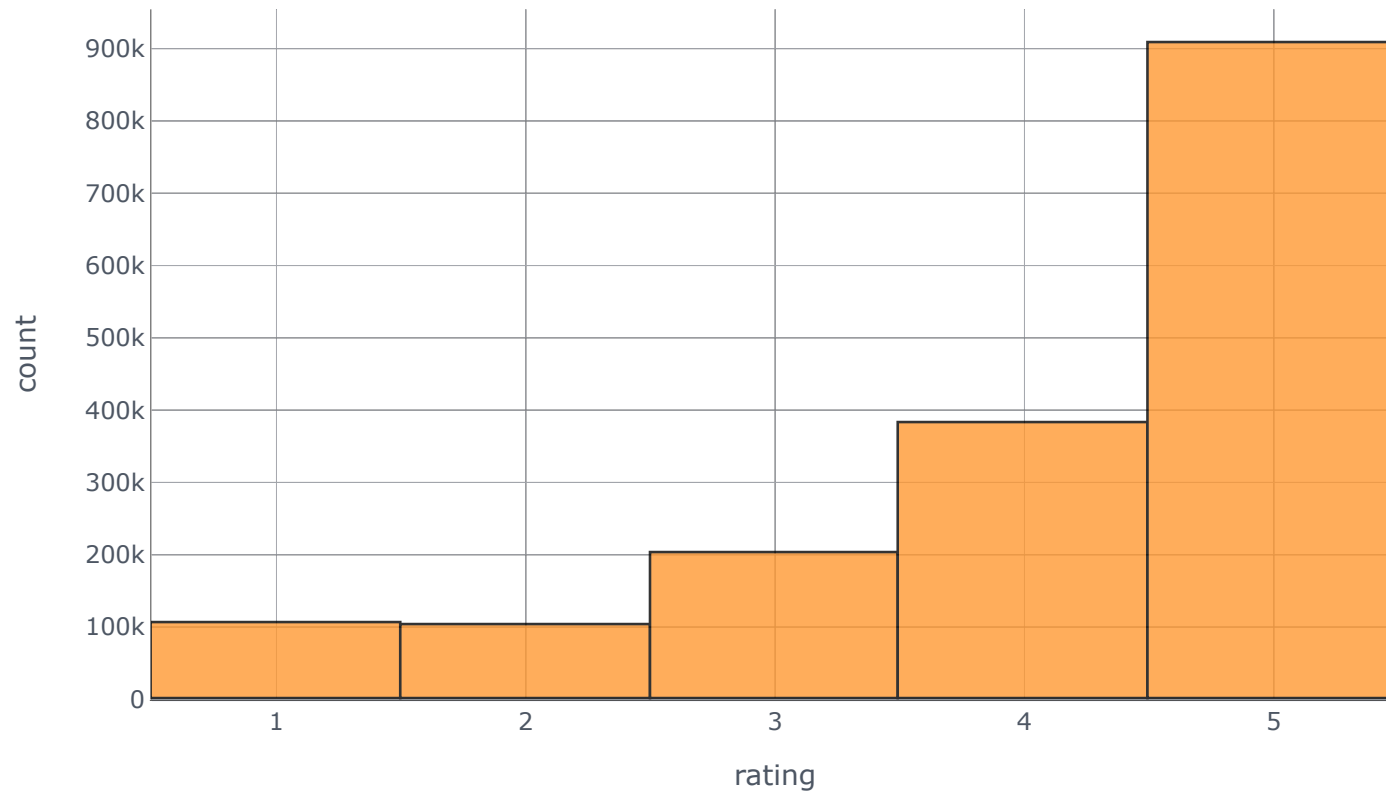




[Export to plot.ly »](#)


```
In [ ]: 1 df['overall'].ipplot(  
2     kind='hist',  
3     xTitle='rating',  
4     linecolor='black',  
5     yTitle='count',  
6     title='Review Rating Distribution')
```

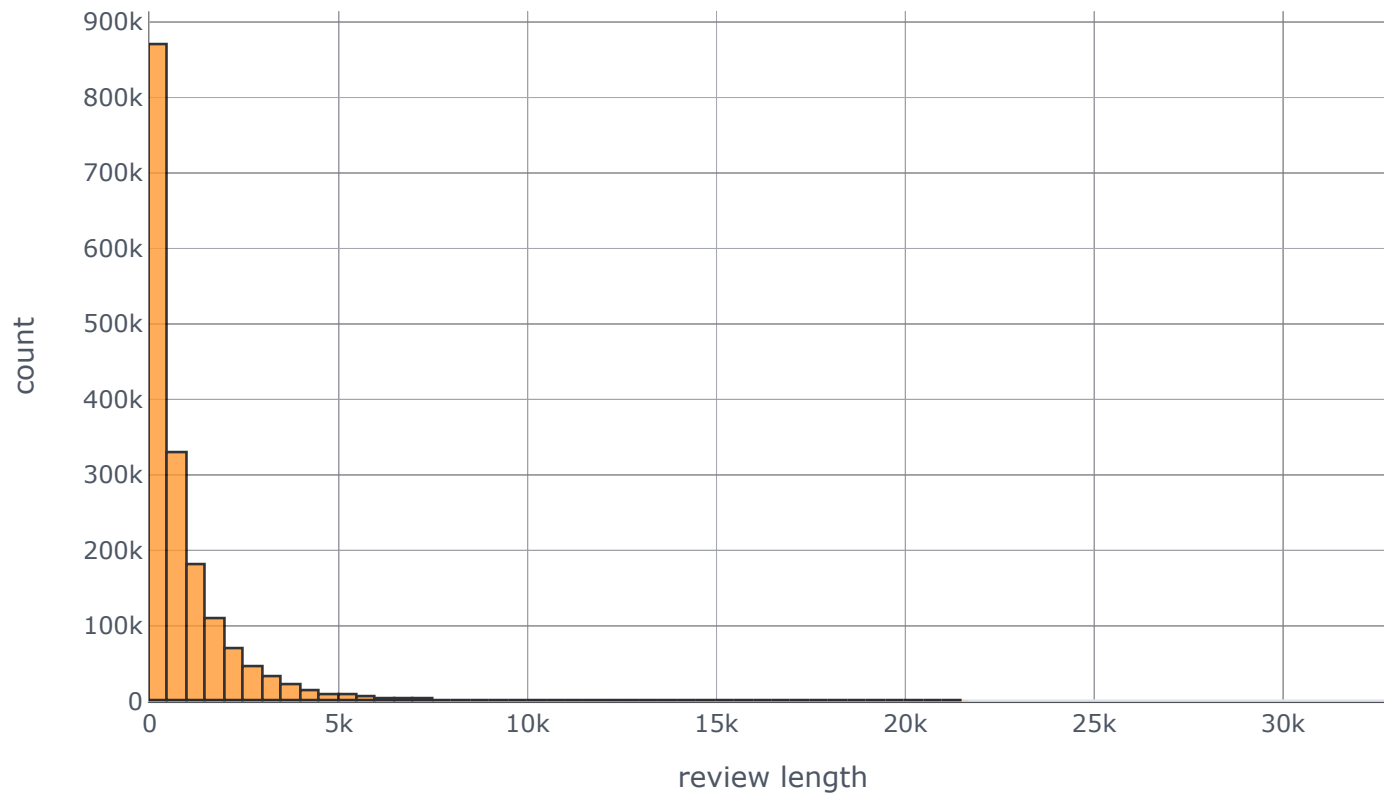
Review Rating Distribution



[Export to plot.ly »](#)

```
In [ ]: 1 df['review_len'].ipplot(  
2     kind='hist',  
3     bins=100,  
4     xTitle='review length',  
5     linecolor='black',  
6     yTitle='count',  
7     title='Review Text Length Distribution')
```

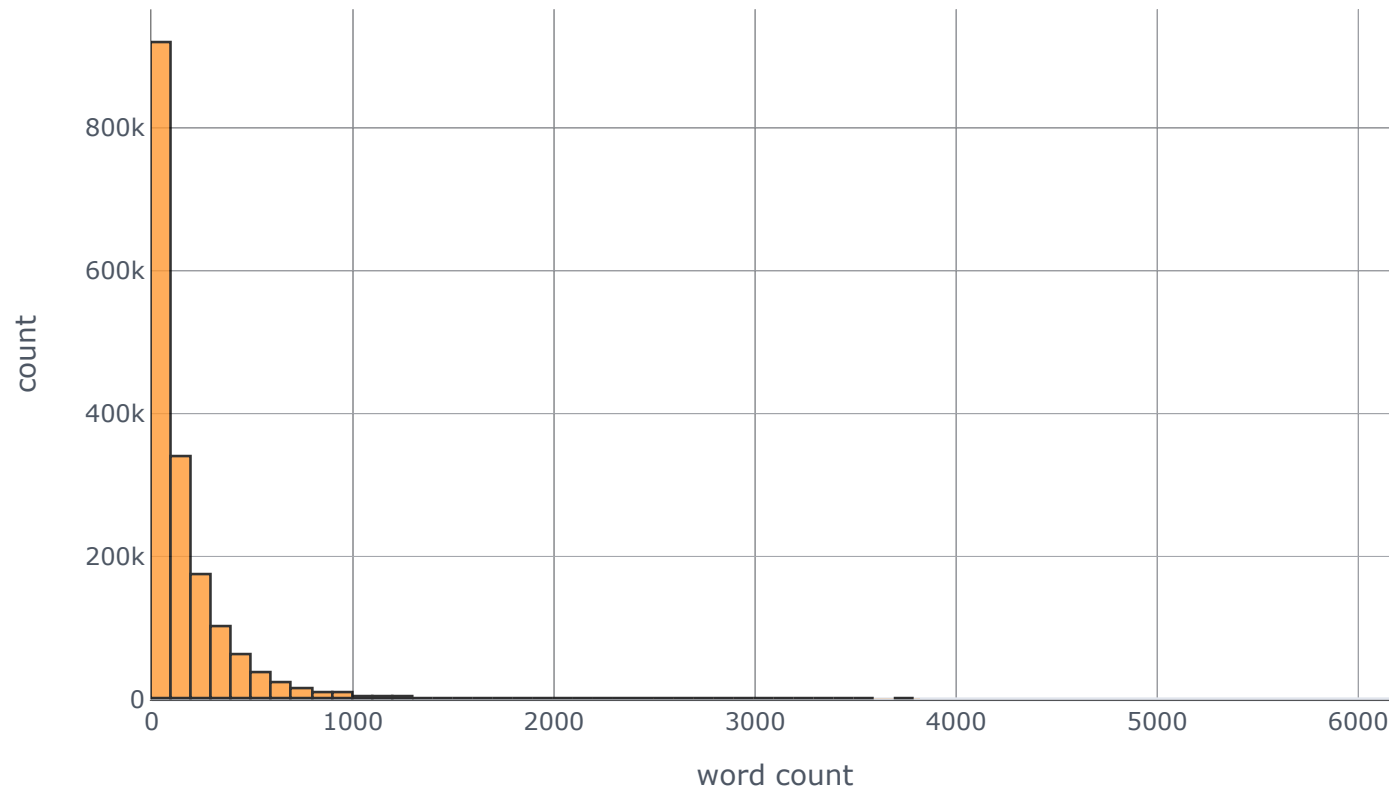
Review Text Length Distribution



[Export to plot.ly »](#)

```
In [ ]: 1 df['word_count'].iplot(  
2     kind='hist',  
3     bins=100,  
4     xTitle='word count',  
5     linecolor='black',  
6     yTitle='count',  
7     title='Review Text Word Count Distribution')
```

Review Text Word Count Distribution



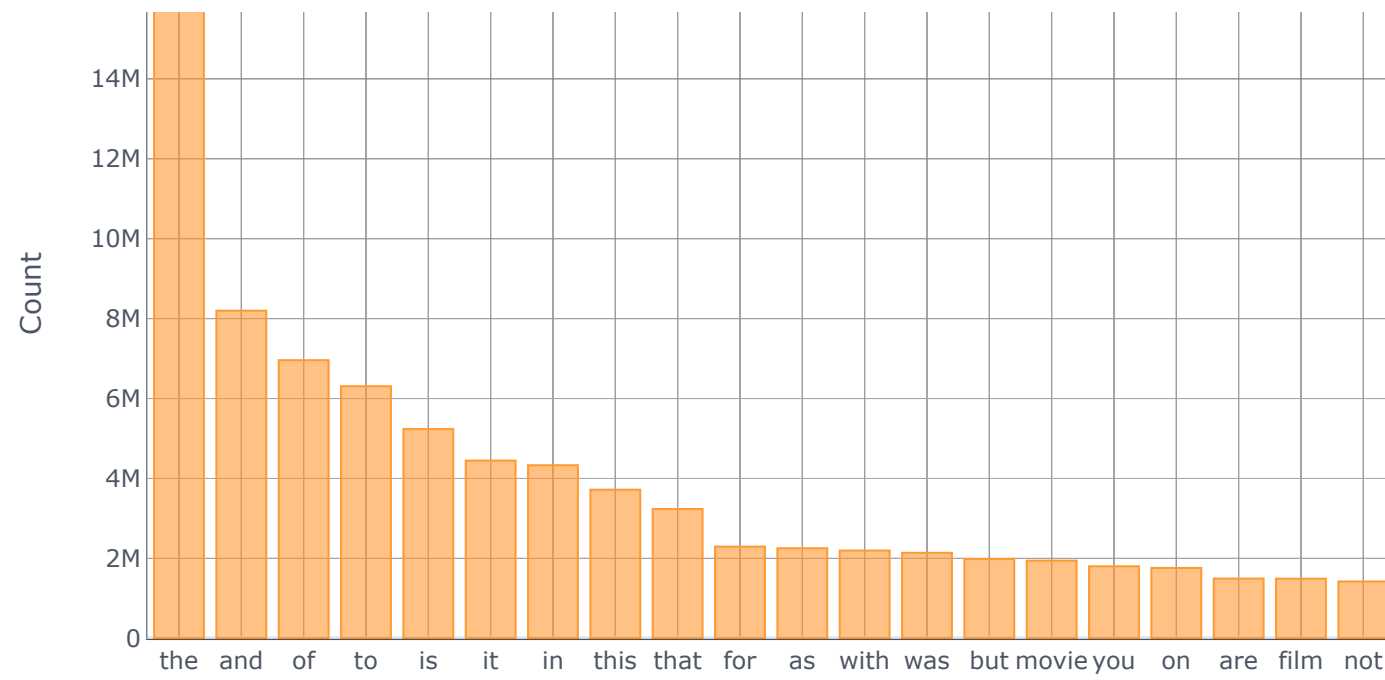
[Export to plot.ly »](#)

```
In [ ]: 1 def get_top_n_words(corpus, n=None):
2         vec = CountVectorizer().fit(corpus)
3         bag_of_words = vec.transform(corpus)
4         sum_words = bag_of_words.sum(axis=0)
5         words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary_.items()]
6         words_freq = sorted(words_freq, key = lambda x: x[1], reverse=True)
7         return words_freq[:n]
8 common_words = get_top_n_words(df['reviewText'], 20)
9 for word, freq in common_words:
10     print(word, freq)
11 df1 = pd.DataFrame(common_words, columns = ['reviewText' , 'count'])
12 df1.groupby('reviewText').sum()['count'].sort_values(ascending=False).plot(
13     kind='bar', yTitle='Count', linecolor='black', title='Top 20 words in review before removing stop words')
```

```
the 16399695
and 8198290
of 6959404
to 6307177
is 5235998
it 4447618
in 4330611
this 3720441
that 3235537
for 2295381
as 2254253
with 2196737
was 2140099
but 1984339
movie 1945059
you 1803016
on 1761735
are 1495803
film 1492831
not 1422246
```

Top 20 words in review before removing stop words



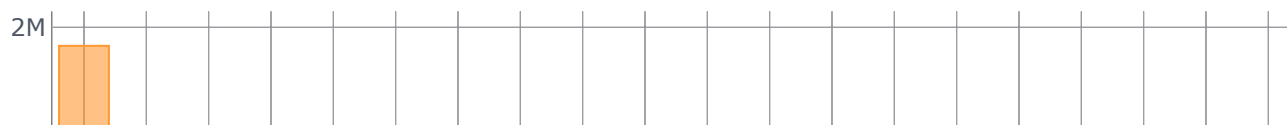


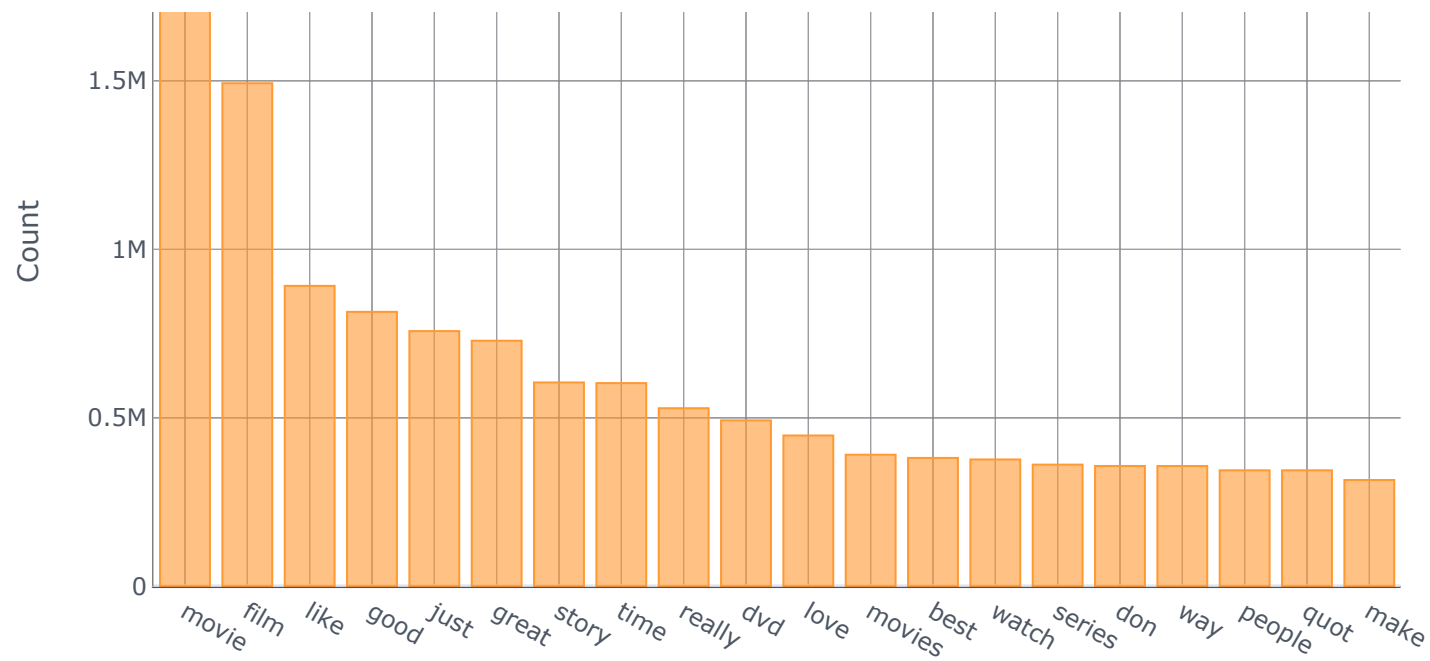
[Export to plot.ly »](#)

```
In [ ]: 1 def get_top_n_words(corpus, n=None):
2         vec = CountVectorizer(stop_words = 'english').fit(corpus)
3         bag_of_words = vec.transform(corpus)
4         sum_words = bag_of_words.sum(axis=0)
5         words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary_.items()]
6         words_freq = sorted(words_freq, key = lambda x: x[1], reverse=True)
7         return words_freq[:n]
8 common_words = get_top_n_words(df['reviewText'], 20)
9 for word, freq in common_words:
10     print(word, freq)
11 df2 = pd.DataFrame(common_words, columns = ['reviewText' , 'count'])
12 df2.groupby('reviewText').sum()['count'].sort_values(ascending=False).plot(
13     kind='bar', yTitle='Count', linecolor='black', title='Top 20 words in review after removing stop words')
```

```
movie 1945059
film 1492831
like 891611
good 814498
just 757456
great 728824
story 605037
time 603128
really 528601
dvd 492446
love 447648
movies 390588
best 381055
watch 376545
series 361202
don 357229
way 357008
people 344288
quot 344266
make 315549
```

Top 20 words in review after removing stop words





[Export to plot.ly »](#)

```
In [ ]: 1 def get_top_n_trigram(corpus, n=None):
2         vec = CountVectorizer(ngram_range=(3, 3), stop_words='english').fit(corpus)
3         bag_of_words = vec.transform(corpus)
4         sum_words = bag_of_words.sum(axis=0)
5         words_freq = [(word, sum_words[0, idx]) for word, idx in vec.vocabulary_.items()]
6         words_freq = sorted(words_freq, key = lambda x: x[1], reverse=True)
7         return words_freq[:n]
8     common_words = get_top_n_trigram(df['reviewText'], 20)
9     for word, freq in common_words:
10         print(word, freq)
11 df6 = pd.DataFrame(common_words, columns = ['reviewText', 'count'])
12 df6.groupby('reviewText').sum()['count'].sort_values(ascending=False).iplot(
13     kind='bar', yTitle='Count', linecolor='black', title='Top 20 trigrams in review after removing stop words')
```

```
In [ ]: 1 trace1 = go.Scatter(
2         x=df['polarity'], y=df['overall'], mode='markers', name='points',
3         marker=dict(color='rgb(102,0,0)', size=2, opacity=0.4)
4     )
5 trace2 = go.Histogram2dContour(
6         x=df['polarity'], y=df['overall'], name='density', ncontours=20,
7         colorscale='Hot', reversescale=True, showscale=False
8     )
9 trace3 = go.Histogram(
10        x=df['polarity'], name='Sentiment polarity density',
11        marker=dict(color='rgb(102,0,0)'),
12        yaxis='y2'
13    )
14 trace4 = go.Histogram(
15        y=df['overall'], name='Rating density', marker=dict(color='rgb(102,0,0)'),
16        xaxis='x2'
17    )
18 data = [trace1, trace2, trace3, trace4]
19
20 layout = go.Layout(
21     showlegend=False,
22     autosize=False,
23     width=600,
24     height=550,
25     xaxis=dict(
26         domain=[0, 0.85],
27         showgrid=False,
28         zeroline=False
29     ),
30     yaxis=dict(
31         domain=[0, 0.85],
32         showgrid=False,
33         zeroline=False
34     ),
35     margin=dict(
36         t=50
37     ),
38     hovermode='closest',
39     bargap=0,
40     xaxis2=dict(
41         domain=[0.85, 1],
42         showgrid=False,
43         zeroline=False
44     ),
45     yaxis2=dict(
```



```
46         domain=[0.85, 1],
47         showgrid=False,
48         zeroline=False
49     )
50 )
51
52 fig = go.Figure(data=data, layout=layout)
53 iplot(fig, filename='2dhistogram-2d-density-plot-subplots')
```

In []:

1