

Importing Required Libraries

```
In [1]: 1 import pandas as pd
        2 import matplotlib.pyplot as plt
        3 import numpy as np
        4 import sklearn
        5 import seaborn as sns
```

```
In [2]: 1 import nltk
        2 import re
        3 import string
        4 from string import digits
        5 from nltk.tokenize import word_tokenize
        6 from nltk.stem import WordNetLemmatizer, PorterStemmer
        7
        8 from sklearn.decomposition import NMF
        9 from sklearn.metrics import accuracy_score
       10 from sklearn.metrics import confusion_matrix
       11 from sklearn.model_selection import train_test_split, GridSearchCV
       12 from sklearn.naive_bayes import MultinomialNB
       13 from sklearn.ensemble import RandomForestClassifier
       14 from sklearn.pipeline import make_pipeline
       15 from matplotlib.pyplot import figure
       16
```

```
In [3]: 1 from sklearn.feature_extraction.text import CountVectorizer , TfidfTransform
```

```
In [4]: 1 from sklearn.utils._testing import ignore_warnings
        2 from sklearn.exceptions import ConvergenceWarning
        3
```

```
In [43]: 1 nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\muham\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
```

Out[43]: True

About the Dataset and Project Objective

The data has been taken from UCI repository and is regarding the use of a drug type for some medical condition alongwith the reviews and ratings from the user. The data was collected from various pharmaceutical company websites. There are 2 files which had been compiled for analysis as train and test data. The data percentages of train and test was set to 75% and 25% respectively.

The objective of the project is whether we can predict the rating a user is likely going to give depending on the review submitted. This is essentially a classification problem which would be addressed by using unsupervised machine learning approach. A comparison would also be done against some machine learning models to assess the effectiveness of the unsupervised model vs supervised model. The target variable is rating with the possible ratings between '1' and '10'.

The reason for taking this project is to practically deploy whatever has been learnt in this course to augment the concepts and approaches to different kind of problems. Additionally, augmenting NLP would also be great in identifying the sentiment of a review which in turn would determine the final rating given by a user/patient/carer.

Citation Felix Gräßer, Surya Kallumadi, Hagen Malberg, and Sebastian Zaunseder. 2018. Aspect-Based Sentiment Analysis of Drug Reviews Applying Cross-Domain and Cross-Data Learning.

Sequence of Work

The project will be carried out in the following sequence:

1. Data Import
2. Data Overview
3. Data Cleaning
4. EDA
5. Models Design
6. Models Testing
7. Models Comparison
8. Results and Analysis
9. Conclusion

1. Importing Data

```
In [87]: 1 df_train=pd.read_csv('drugsComTrain_raw.tsv',sep='\t')
          2 df_test=pd.read_csv('drugsComTest_raw.tsv',sep='\t')
```

2. Data Overview

In [7]: 1 df_train.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 161297 entries, 0 to 161296
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Unnamed: 0      161297 non-null int64
1   drugName        161297 non-null object
2   condition       160398 non-null object
3   review          161297 non-null object
4   rating          161297 non-null float64
5   date            161297 non-null object
6   usefulCount     161297 non-null int64
dtypes: float64(1), int64(2), object(4)
memory usage: 8.6+ MB
```

In [8]: 1 df_test.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53766 entries, 0 to 53765
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Unnamed: 0      53766 non-null int64
1   drugName        53766 non-null object
2   condition       53471 non-null object
3   review          53766 non-null object
4   rating          53766 non-null float64
5   date            53766 non-null object
6   usefulCount     53766 non-null int64
dtypes: float64(1), int64(2), object(4)
memory usage: 2.9+ MB
```

1. The **training dataset** has **161297 observations** against **7 features** whereas the **test dataset** has **53766 observations** against the same **7 features**.

2. The size of **training** is **8.6mb** and **test** **2.9mb**.

3. The dataset contains 4 object types, 2 int types and 1 float i.e. the ratings.

Feature Info

1. drugName (categorical): name of drug
2. condition (categorical): name of condition
3. review (text): patient review
4. rating (numerical): 10 star patient rating
5. date (date): date of review entry
6. usefulCount (numerical): number of users who found review useful

The first feature 'Unnamed:0' is simply the id the person giving the review.

3. Data Cleaning

Since, our objective for the project is to be able to predict the ratings based on the user review, we will be getting rid of the remaining columns.

```
In [88]: 1 df_train.drop(['Unnamed: 0', 'drugName', 'condition', 'date', 'usefulCount'], i
          2 df_test.drop(['Unnamed: 0', 'drugName', 'condition', 'date', 'usefulCount'], in
```

Checking for NA values

```
In [8]: 1 df_train.isna().sum()
```

```
Out[8]: review    0
        rating    0
        dtype: int64
```

```
In [9]: 1 df_test.isna().sum()
```

```
Out[9]: review    0
        rating    0
        dtype: int64
```

The number of observations remains the same, however, after getting rid of the redundant features as per our project objective, we have only 2 columns remaining namely 'review' and 'rating'.

Processing Text

There are two main steps involved in the processing of text, i.e. removing unnecessary alphanumeric, punctuations, numbers etc. and then vectorizing the sentences into words for further study and modelling.

We undertook the following steps in the model training:

1. Word Tokenization
2. Removal of Alpha numerics
3. Removal of Punctuations
4. Removal of Stop Words
5. Stemming the words
6. Lemmatizing the words
7. rejoining the words to return the modified sentence/text

Text Cleaning (review column)

```
In [7]: 1 df_train['review'] = df_train['review'].astype(str)
```

```

In [8]: 1 stop = set(nltk.corpus.stopwords.words('english'))
2 stop_words = stop.union({'said', 'mr', 'known', 'll', 'ms', 've', 'you', 'your',
3 'yours', 'ha', 'thi', 'now', 'onli', 'im', 'beca
4 "even", "go", "realli", "didnt", "abl",
5 'year'}) ###IMPORTANT#####
6
7
8 def cleaning(x):
9
10     # converting to words
11     tokens = word_tokenize(x)
12
13     # convert to lower case
14     words = [w.lower() for w in tokens]
15
16     # removing alphanumerics
17     words = [word for word in words if word.isalpha()]
18
19     # removing punctuations
20     table = str.maketrans('', '', string.punctuation)
21     stripped = [w.translate(table) for w in words]
22     # stripped = data['review1'].str.replace('[^\w\s]', '')
23
24
25     # removing stopwords after modification
26     words_mod = [w for w in stripped if not w in stop_words]
27
28     # word stemming
29     stemmer = PorterStemmer()
30     # stemmed_words = " ".join([stemmer.stem(w) for w in words_mod])
31     stemmed_words = [stemmer.stem(w) for w in words_mod]
32
33     # word lemmatization
34     lemmatizer = WordNetLemmatizer()
35     sentence = " ".join([lemmatizer.lemmatize(w) for w in stemmed_words])
36
37     sen = " ".join(sentence.split())
38
39     # sentence = " ".join(words)
40
41     return (sen)

```

```

In [9]: 1 train_text = df_train.review.apply(lambda x: cleaning(x))

```

```

In [10]: 1 vector = TfidfVectorizer( min_df = 50, ngram_range = (1,1), stop_words='eng
2 word_vec = vector.fit_transform(train_text)

```

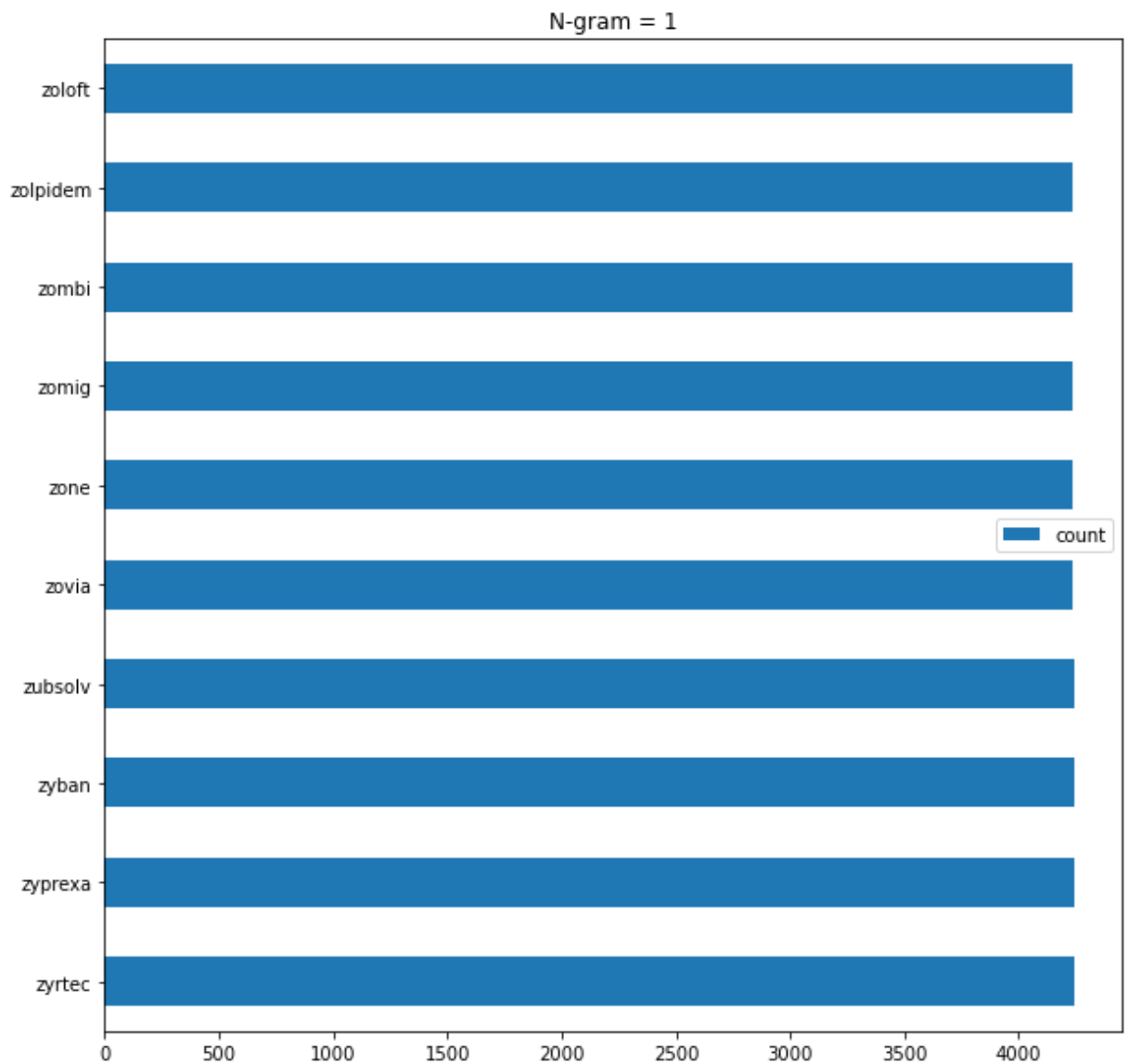
```

In [11]: 1 word_df = pd.DataFrame([vector.vocabulary_]).T
2 word_df.columns = ['count']
3 word_df.sort_values(by = 'count' , inplace = True ,ascending=False)

```

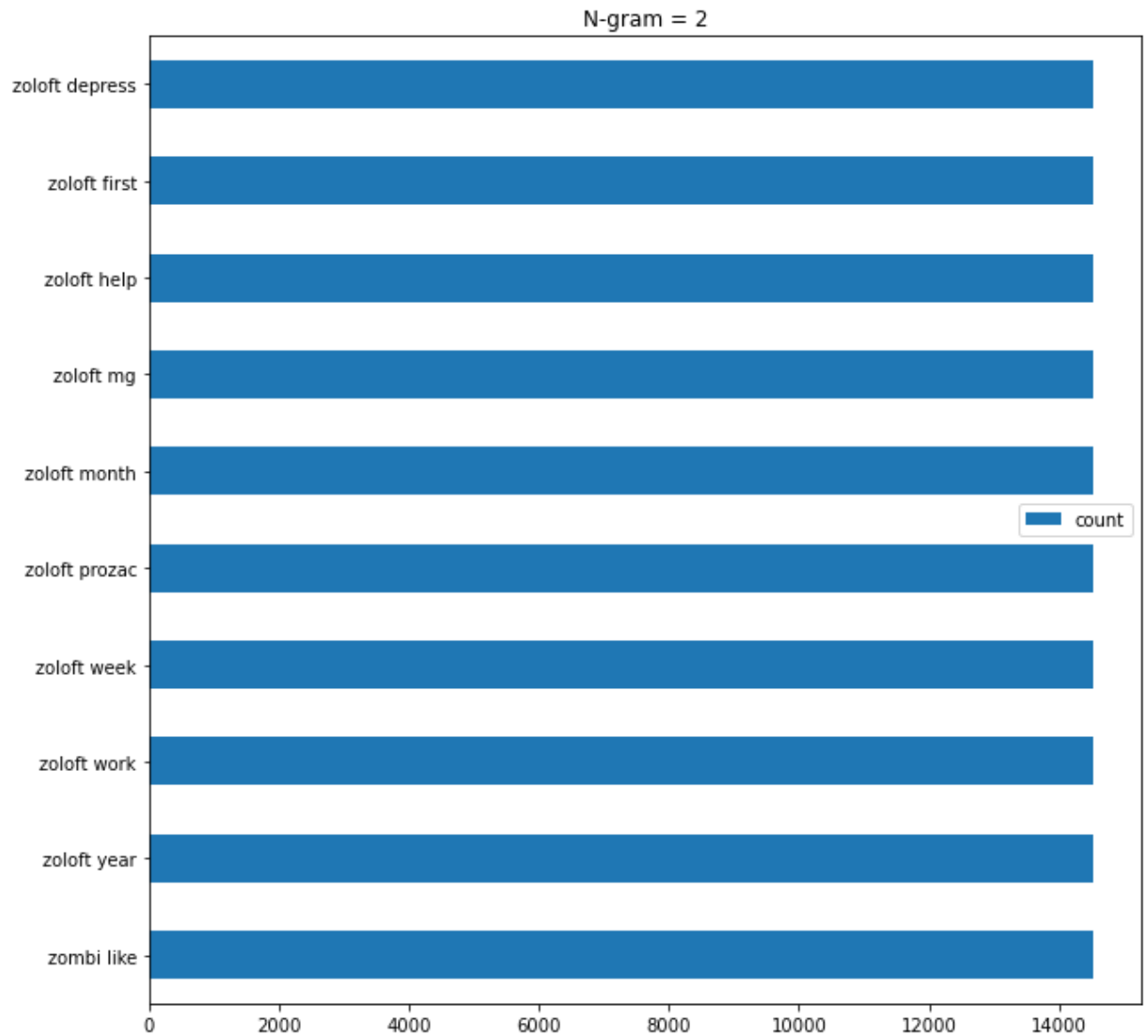
```
In [12]: 1 word_df.head(10).plot(kind = 'barh', figsize = (10,10))  
        2 plt.title('N-gram = 1')
```

Out[12]: Text(0.5, 1.0, 'N-gram = 1')



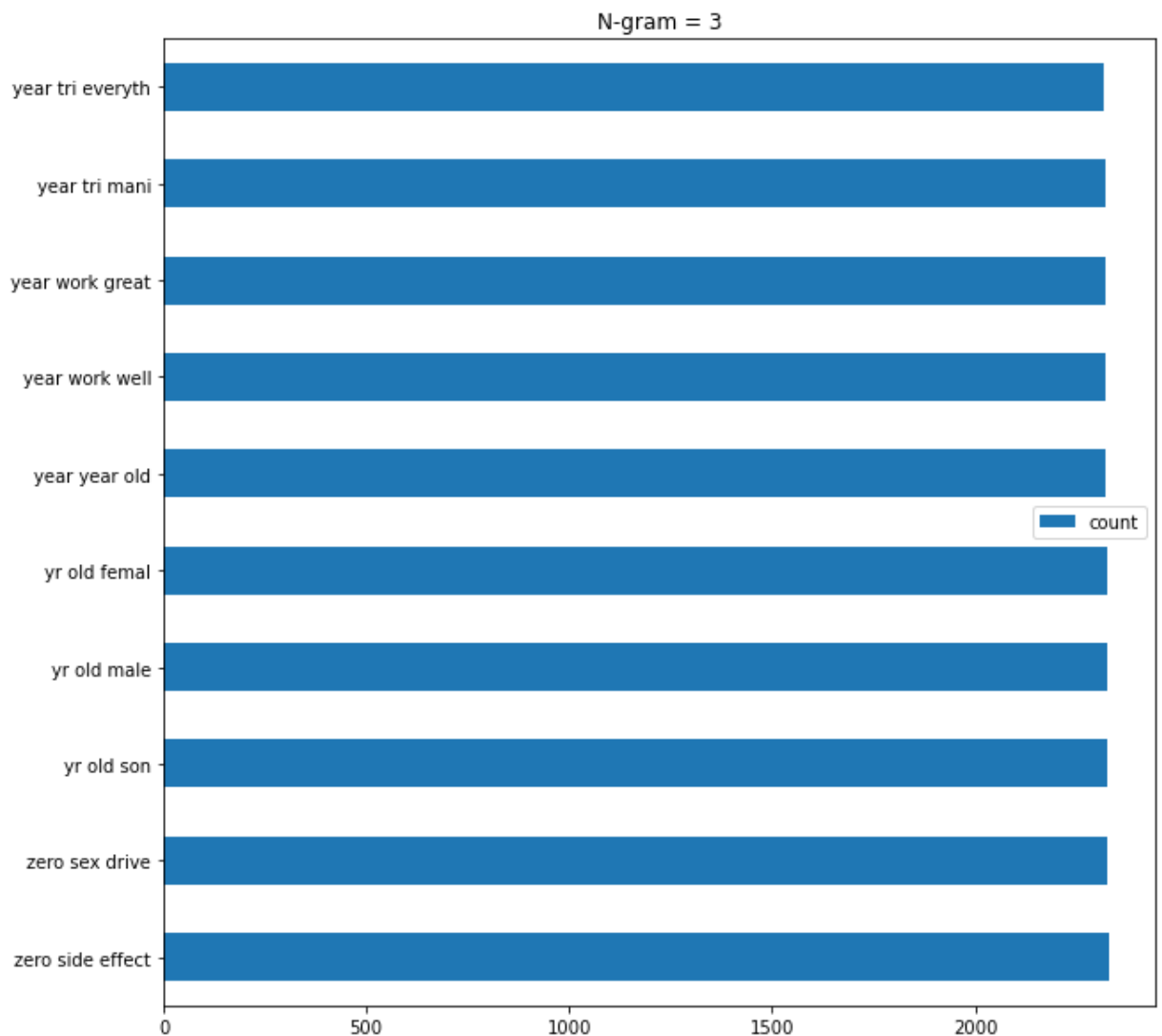
```
In [13]: 1 vector = TfidfVectorizer( min_df = 50, ngram_range = (2,2))
2 word_vec = vector.fit_transform(train_text)
3
4 word_df = pd.DataFrame([vector.vocabulary_]).T
5 word_df.columns = ['count']
6 word_df.sort_values(by = 'count' , inplace = True ,ascending=False)
7
8 word_df.head(10).plot(kind = 'barh', figsize = (10,10))
9 plt.title('N-gram = 2')
```

Out[13]: Text(0.5, 1.0, 'N-gram = 2')



```
In [14]: 1 vector = TfidfVectorizer( min_df = 50, ngram_range = (3,3))
2 word_vec = vector.fit_transform(train_text)
3
4 word_df = pd.DataFrame([vector.vocabulary_]).T
5 word_df.columns = ['count']
6 word_df.sort_values(by = 'count' , inplace = True ,ascending=False)
7
8 word_df.head(10).plot(kind = 'barh', figsize = (10,10))
9 plt.title('N-gram = 3')
```

Out[14]: Text(0.5, 1.0, 'N-gram = 3')



We can see from the graphs above that the name of the medicine is most occurring in the data when the min occurrence is set a 50 whereas the time and effect are mostly displayed when we are referring to the ngrams set to 3.

4. EDA

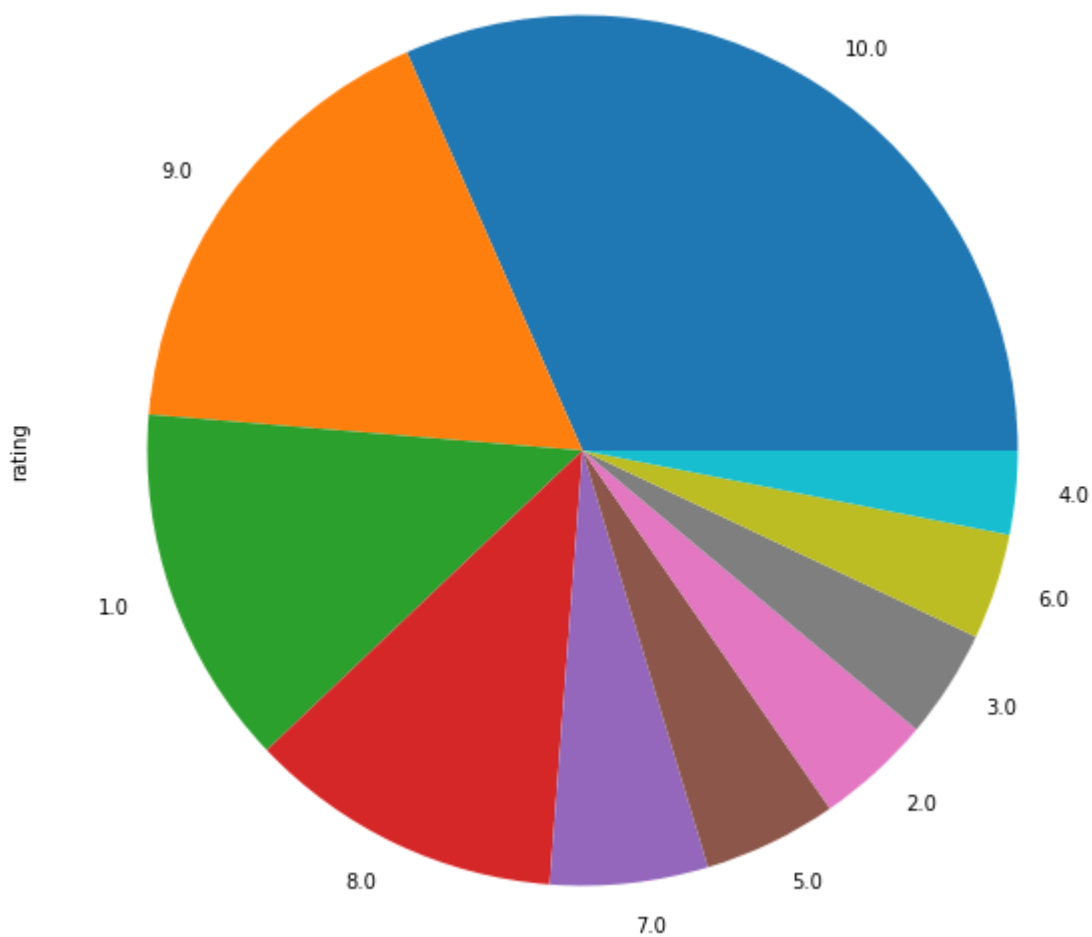

```
In [15]: 1 df_train.describe()
```

Out[15]:

	rating
count	161297.000000
mean	6.994377
std	3.272329
min	1.000000
25%	5.000000
50%	8.000000
75%	10.000000
max	10.000000

```
In [16]: 1 df_train.rating.value_counts().plot.pie(figsize = (10,10))
```

Out[16]: <AxesSubplot:ylabel='rating'>

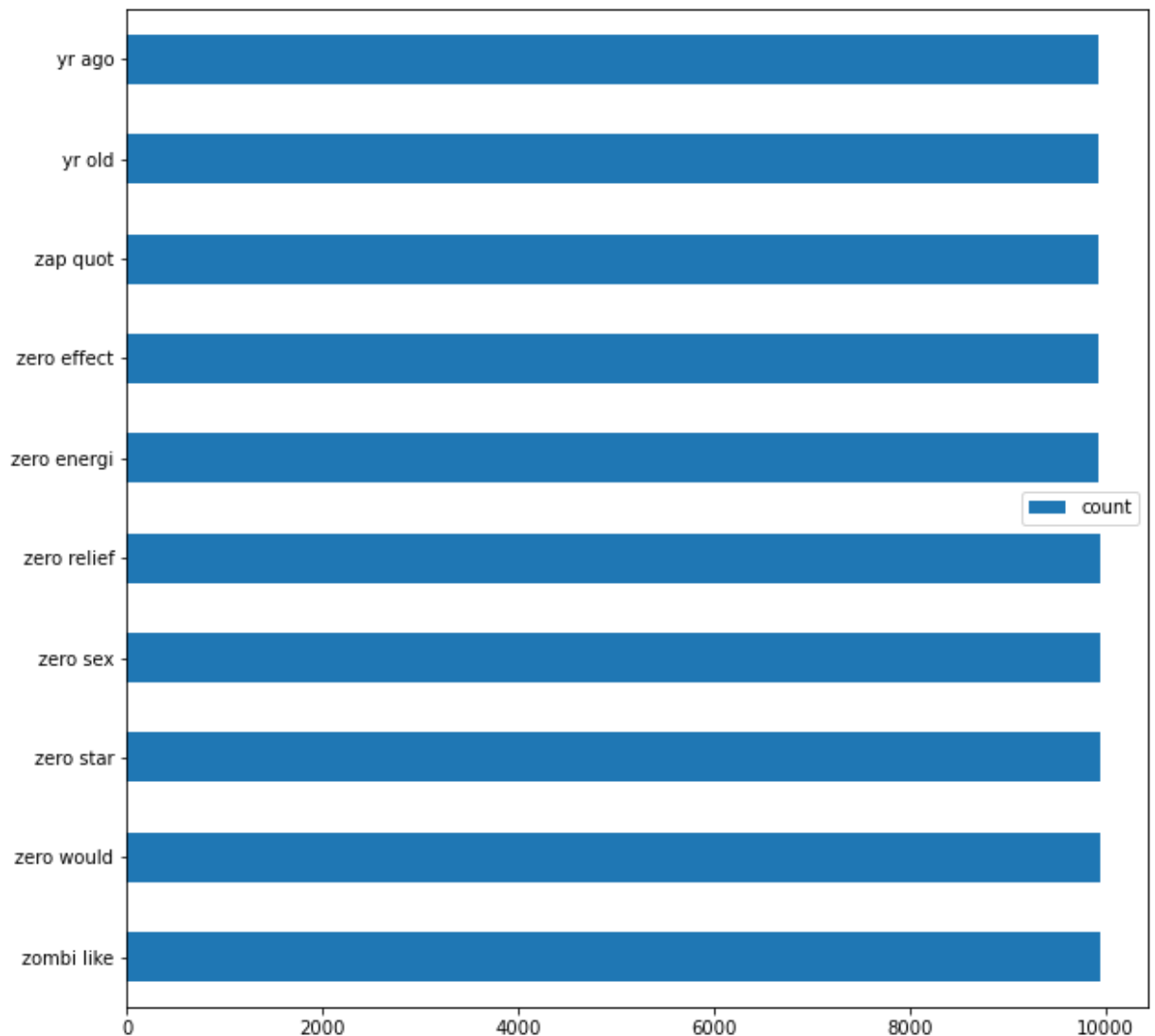


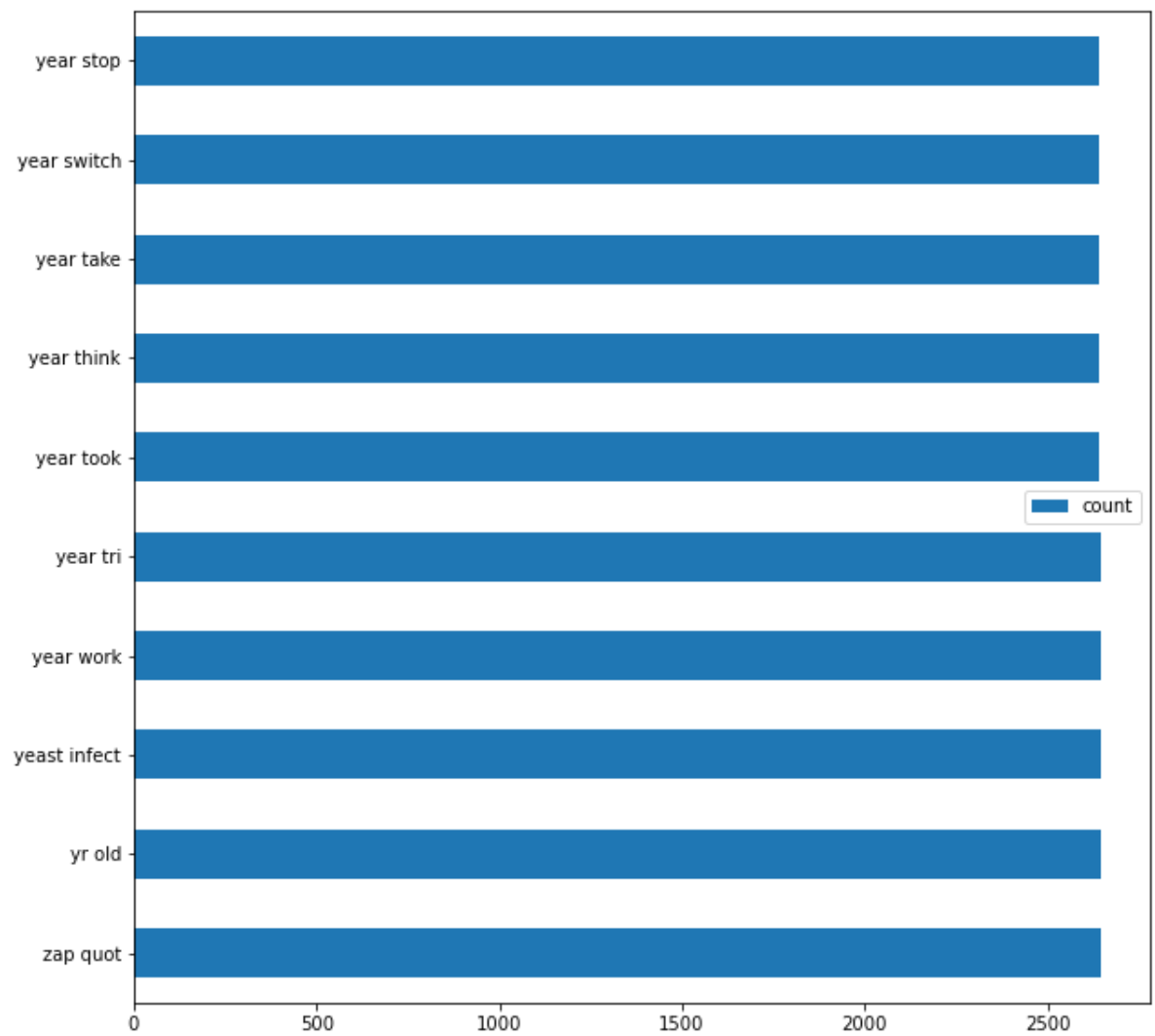
In [17]:

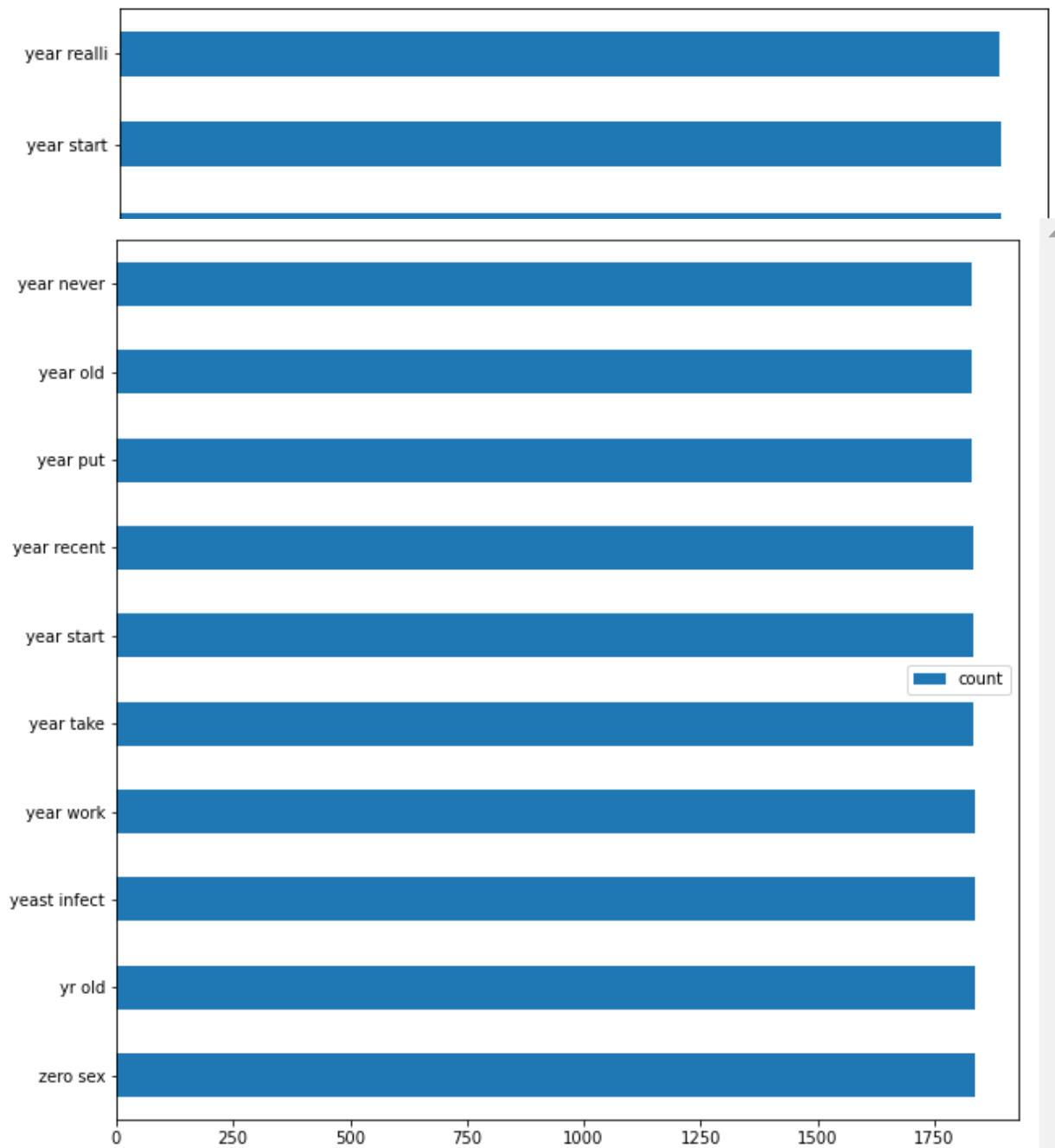
```
1 df_1 = df_train[df_train['rating'] == 1]
2 df_2 = df_train[df_train['rating'] == 2]
3 df_3 = df_train[df_train['rating'] == 3]
4 df_4 = df_train[df_train['rating'] == 4]
5 df_5 = df_train[df_train['rating'] == 5]
6 df_6 = df_train[df_train['rating'] == 6]
7 df_7 = df_train[df_train['rating'] == 7]
8 df_8 = df_train[df_train['rating'] == 8]
9 df_9 = df_train[df_train['rating'] == 9]
10 df_10 = df_train[df_train['rating'] == 10]
```

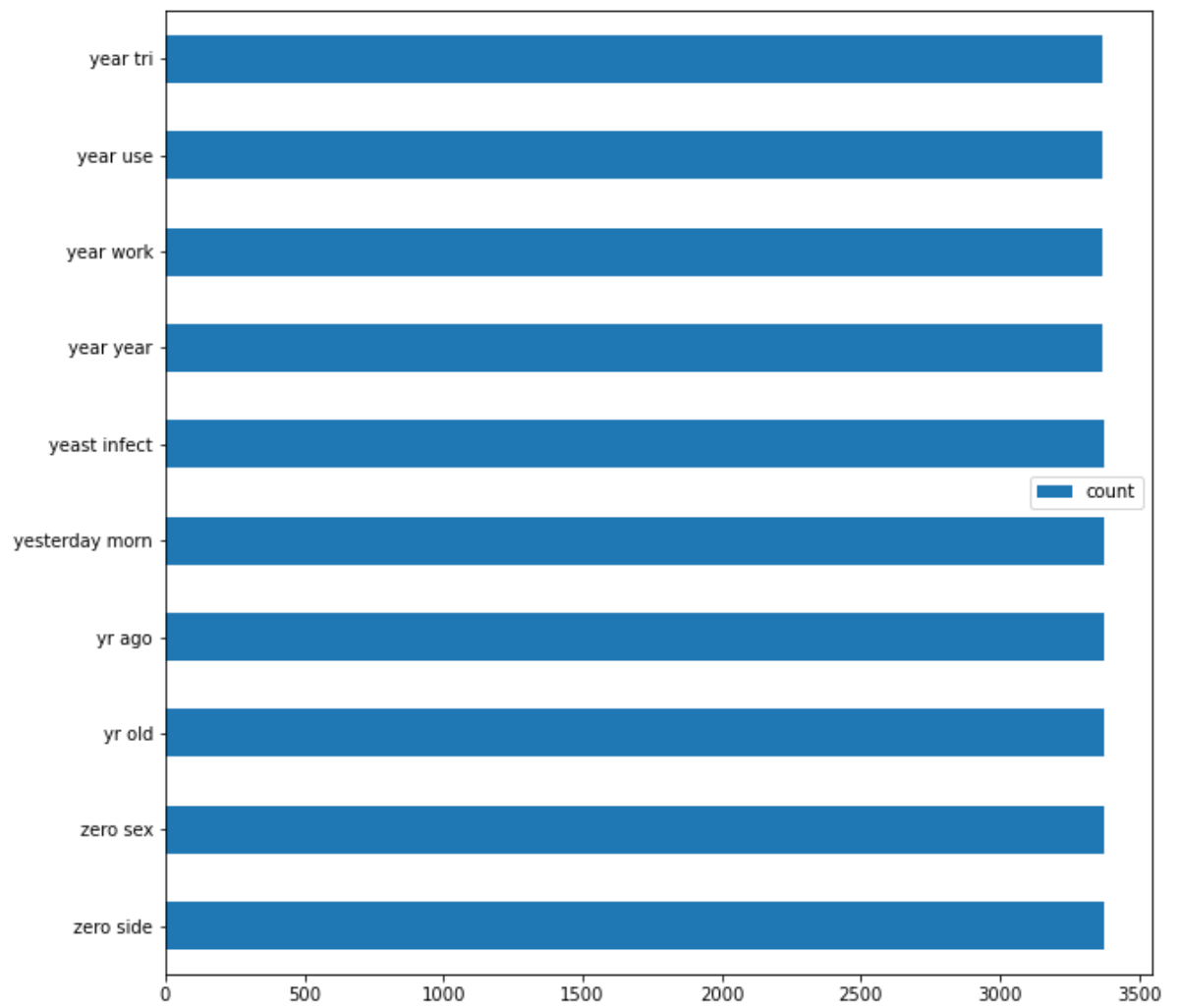
In [18]:

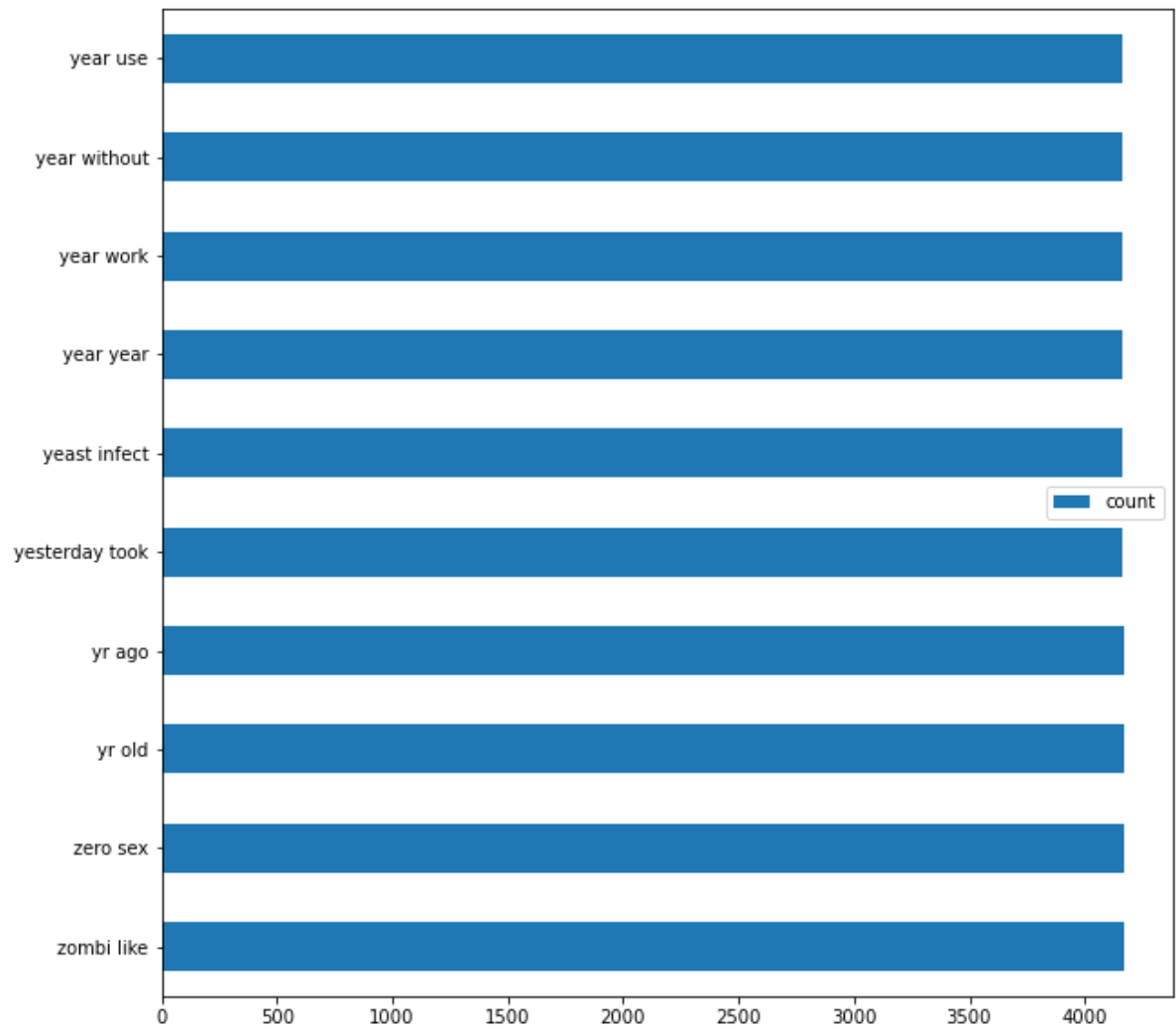
```
1 data= [df_1, df_2,df_3,df_4,df_5,df_6,df_7,df_8,df_9,df_10]
2 vector = TfidfVectorizer(min_df = 10 , ngram_range = (2,2))
3
4 for i in data:
5     clean_data = i.review.apply(lambda x: cleaning(x))
6     word_vec = vector.fit_transform(clean_data)
7
8     word_df = pd.DataFrame([vector.vocabulary_]).T
9     word_df.columns = ['count']
10    word_df.sort_values(by = 'count' , inplace = True ,ascending=False)
11
12    word_df.head(10).plot(kind = 'barh', figsize = (10,10))
13    # plt.title(i)
14    plt.show()
```

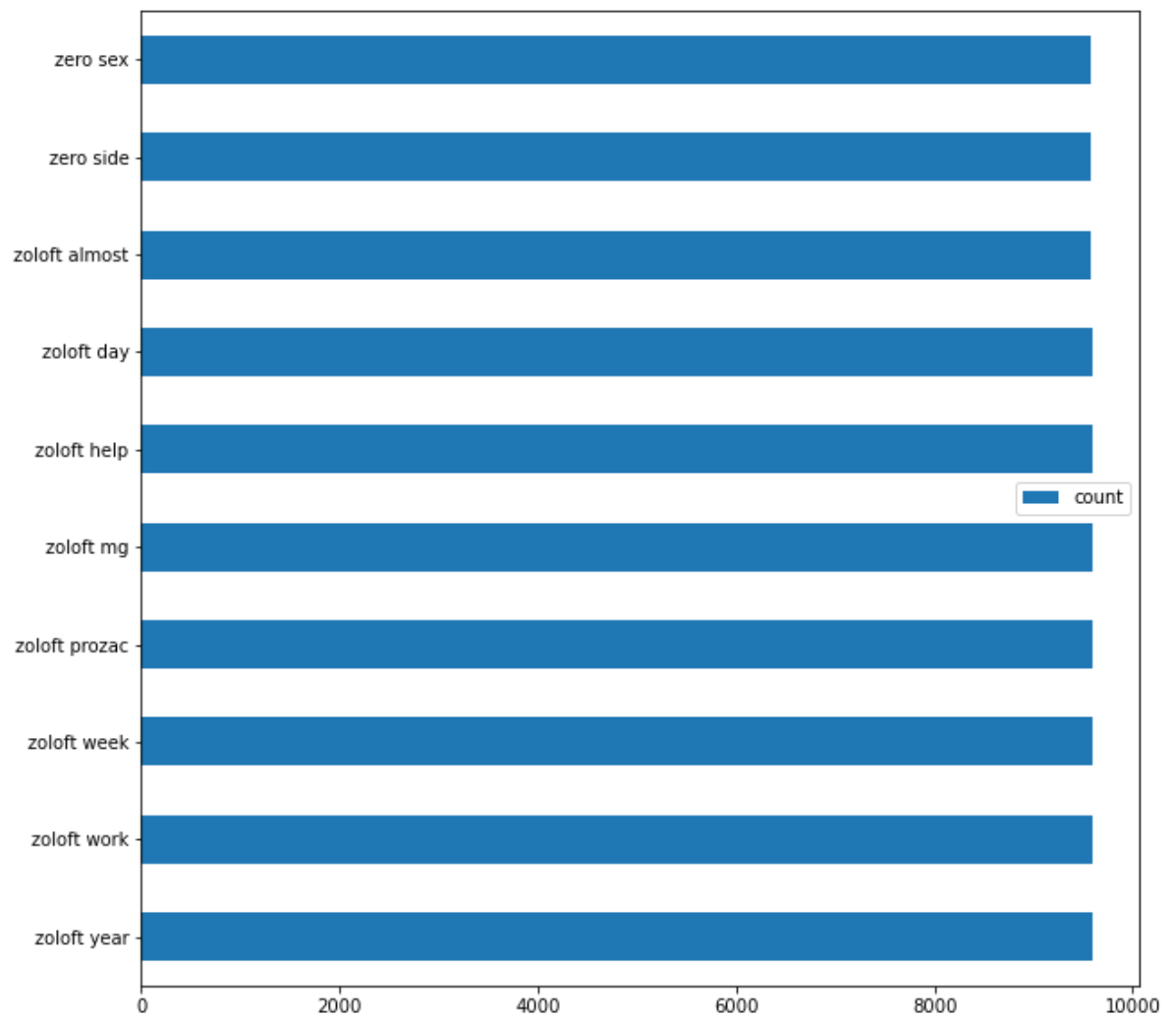


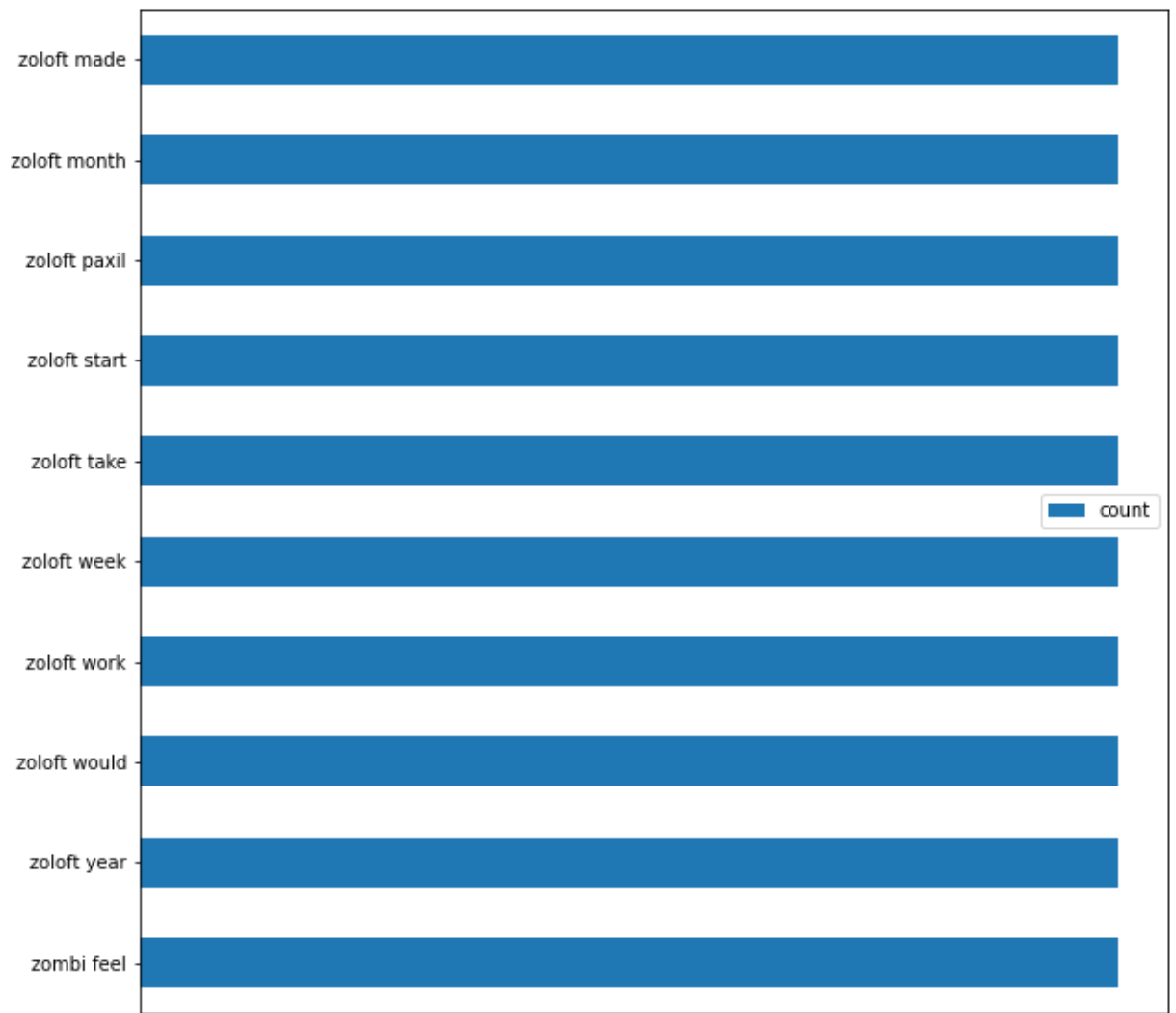


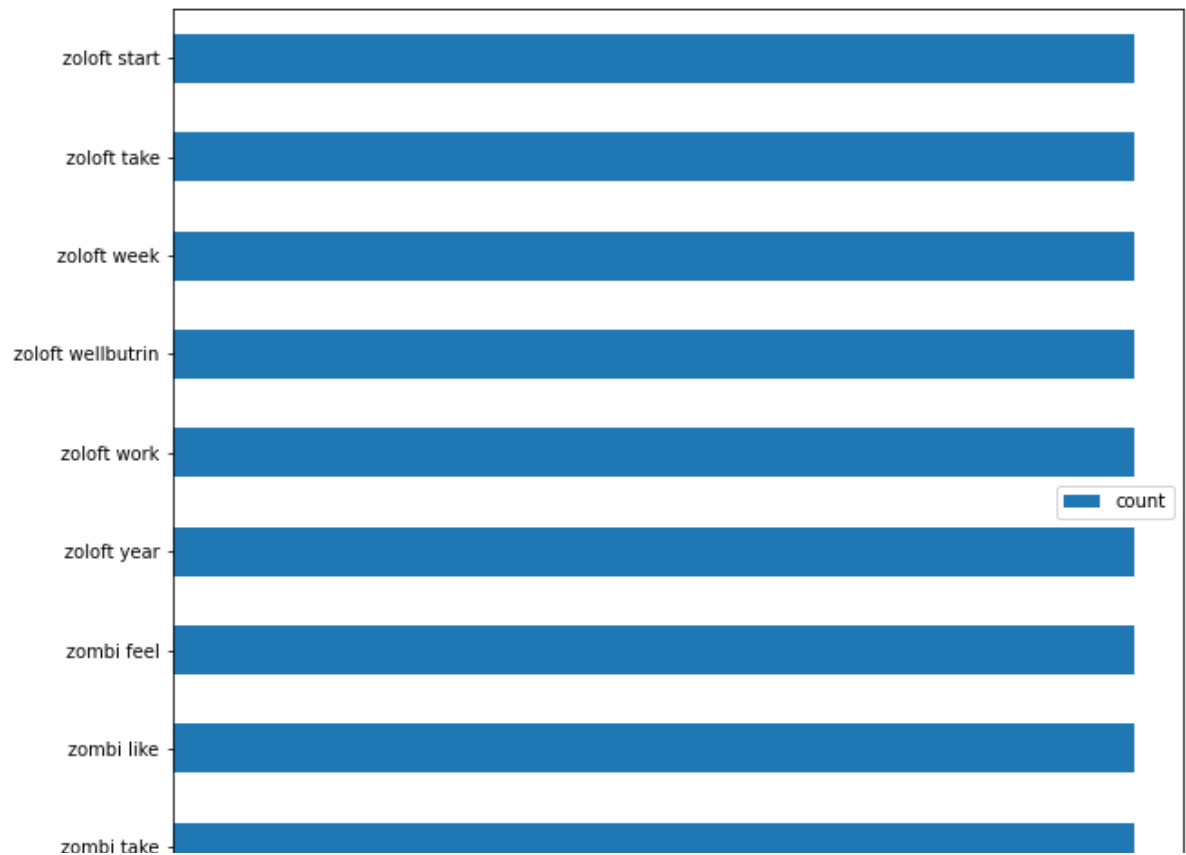












Here, we can take two things, firstly the ratings are not divided equally and a vast majority is skewed towards the higher side i.e. rating '9' and '10'. Secondly, the words among the various ratings seem to be quite similar. Due to this, our model accuracy will be compromised. Hence, it would be prudent to modify the ratings. **We will set the ratings into 3 categories, if rating is less than '5', we'll set it to '0', if it is within 5 to 8, we'll set it to '1' and higher than '8' to 2.**

'0' implies 'Bad'.

'1' implies 'average'.

'2' implies 'good'.

```
In [19]: 1 mod_rating = []
2 for i in df_train.rating:
3     if i < 5:
4         mod_rating.append(0)
5
6     elif i > 8:
7         mod_rating.append(2)
8
9     elif i >= 5 and i <=8:
10        mod_rating.append(1)
11
12 df_train['mod_rating'] = mod_rating
```

```
In [20]: 1 df_train.head()
```

```
Out[20]:
```

	review	rating	mod_rating
0	"It has no side effect, I take it in combinati...	9.0	2
1	"My son is halfway through his fourth week of ...	8.0	1
2	"I used to take another oral contraceptive, wh...	5.0	1
3	"This is my first time using any form of birth...	8.0	1
4	"Suboxone has completely turned my life around...	9.0	2

EDA Summary / Discussion

In the initial phase, we assessed the division of the ratings across the various markings. We found that ratings are not equally distributed and this would in turn would compromise the accuracy of any model. Moreover, we also found that the words were quite similar to each other in various brackets of the ratings. Hence, in order to counter skewness and redundancy of words in the various rating brackets, we grouped them together into 3 categories. From, this we would have ample amount of data on each rating to train an acceptable model.

Dividing into Train and Text Data

```
In [21]: 1 x = df_train.review
2 y = df_train.mod_rating
3
4 X = x.apply(lambda a: cleaning(a))
5
6
```

```
In [22]: 1 X_train , X_test, y_train, y_test = train_test_split(X , y, test_size =0.2,
```

4. Initial Model Design

We will design the following models in sequence:

1. Non-Zero Matrix Factorization (Unsupervised Learning)
2. Random Forest (Supervised Learning)
3. Naive Bayes (Supervised Learning)

(1) Non-Zero Matrix factorization

For the NMF, we will use the entire data for design initially. Additionally, we will also combine the test data to it to increase its accuracy since it is not effected by overfitting.

```
In [46]: 1 clf = NMF(n_components=3, random_state=2)
2
3 vector = TfidfVectorizer( ngram_range = (1,1), stop_words='english')
4 word_vec = vector.fit_transform(df_train.review)
5
6 clf.fit(word_vec)
7 w = clf.transform(word_vec)
8 # h = clf.components_
9
10 tp=pd.DataFrame({'id':pd.DataFrame(w).idxmax(axis=1)})
11 ## Number of categorizations
12
13 tp.value_counts()
14
15 accuracy_score(tp['id'].to_numpy() , df_train['mod_rating'])
```

D:\Anaconda_Launcher\lib\site-packages\sklearn\decomposition_nmf.py:312: FutureWarning: The 'init' value, when 'init=None' and n_components is less than n_samples and n_features, will be changed from 'nndsvd' to 'nndsvda' in 1.1 (renaming of 0.26).

warnings.warn(("The 'init' value, when 'init=None' and "

Out[46]: 0.28423343273588475

(2) Random Forest

```
In [27]: 1 model_rf = make_pipeline(TfidfVectorizer(stop_words='english'), RandomForest
2
3 model_rf.fit(X_train, y_train)
4
5 labels_train = model_rf.predict(X_test)
6 accuracy_score(y_test , labels_train)
```

Out[27]: 0.7960012399256045

(3) Naive Bayes

```
In [28]: 1
2 model_nb = make_pipeline(TfidfVectorizer(stop_words='english'), MultinomialNB
3
4 model_nb.fit(X_train, y_train)
5
6 labels_train = model_nb.predict(X_test)
7 accuracy_score(y_test , labels_train)
```

Out[28]: 0.596342219466832

In the above modelling, the models have been trained on the test data and their accuracy is based on the test data. The training and test data was achieved after dividing the original training dataset.

The 3 models seem to be performing quite differently with the **NMF** giving an accuracy of only 28%, Naive Bayes at 60% and Random Forest at 80%.

Accessing Models on the Actual Test Data

```
In [29]: 1 df_test['review'] = df_test['review'].astype(str)
2 test_clean = df_test.review.apply(lambda a: cleaning(a))
3
4 mod_test_rating = []
5 for i in df_test.rating:
6     if i < 5:
7         mod_test_rating.append(0)
8
9     elif i > 8:
10        mod_test_rating.append(2)
11
12    elif i >= 5 and i <=8:
13        mod_test_rating.append(1)
14
15 df_test['mod_rating'] = mod_test_rating
```

```
In [30]: 1 acc = []
2 # NMF
3 word_vec_test = vector.transform(df_test.review)
4 w_t = clf.transform(word_vec_test)
5 tpt=pd.DataFrame({'id':pd.DataFrame(w_t).idxmax(axis=1)})
6 acc.append(accuracy_score(tpt['id'].to_numpy() , df_test['mod_rating']))
7
8 # RF
9 res_rf = model_rf.predict(test_clean)
10 acc.append(accuracy_score(df_test.mod_rating,res_rf))
11
12 # NB
13 res_nb = model_nb.predict(test_clean)
14 acc.append(accuracy_score(df_test.mod_rating , res_nb))
15
16 acc
```

Out[30]: [0.28136740691143103, 0.7926198712941264, 0.5967897928058624]

1

5. Revised Model

Optimizing NMF

```

In [61]: 1 train_acc = []
          2 for i in range(60):
          3     vector = TfidfVectorizer(stop_words = 'english',min_df = i)
          4     word_vec = vector.fit_transform(df_train.review)
          5     clf = NMF(init = 'nndsvda',n_components=3, random_state=5)
          6     clf.fit(word_vec)
          7     w = clf.transform(word_vec)
          8     # h = clf.components_
          9     tp=pd.DataFrame({'id':pd.DataFrame(w).idxmax(axis=1)})
         10     train_acc.append(accuracy_score(tp['id'].to_numpy() , df_train.mod_ratin
         11

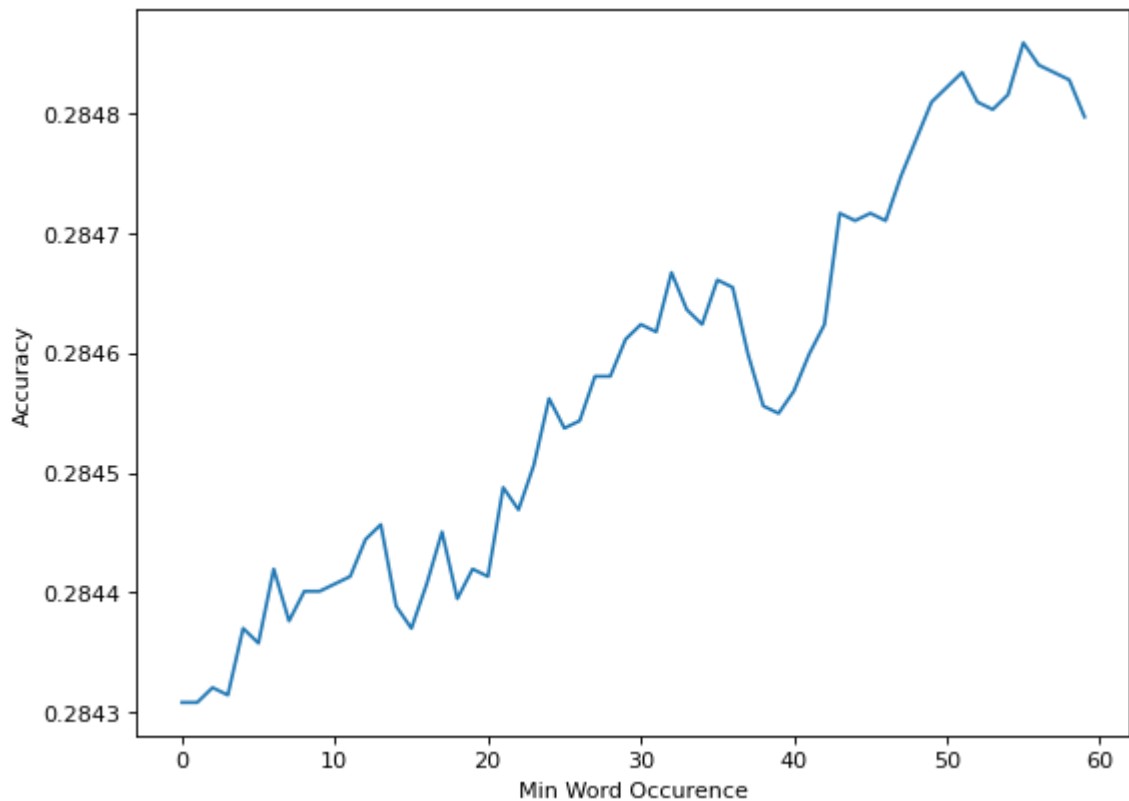
```

```

In [62]: 1 figure(figsize=(8, 6), dpi=80)
          2 plt.plot(train_acc)
          3 plt.xlabel("Min Word Occurence")
          4 plt.ylabel("Accuracy")

```

Out[62]: Text(0, 0.5, 'Accuracy')



```

In [63]: 1 train_acc.index(max(train_acc))

```

Out[63]: 55

```

In [73]: 1 vector = TfidfVectorizer( min_df=9, stop_words='english')
          2 word_vec = vector.fit_transform(train_text)

```

```
In [74]: 1 initial = ["random", "nndsvd", "nndsvda", "nndsvdar"]
2 train_score = [[],[],[],[]]
3 test_score = [[],[],[],[]]
4 beta = ['frobenius', 'kullback-leibler']
5 for i in range(len(initial)):
6     for j in beta:
7         if (j == 'kullback-leibler'):
8             clf = NMF(init = initial[i], n_components=3, random_state=5, beta=j)
9         else:
10            clf = NMF(init = initial[i], n_components=3, random_state=5, beta=j)
11            clf.fit(word_vec)
12            w = clf.transform(word_vec)
13            h = clf.components_
14
15            tp=pd.DataFrame({'id':pd.DataFrame(w).idxmax(axis=1)})
16
17            train_score[i].append(accuracy_score(tp['id'].to_numpy(), df_train[
18
19
```

D:\Anaconda_Launcher\lib\site-packages\sklearn\decomposition_nmf.py:226: UserWarning: The multiplicative update ('mu') solver cannot update zeros present in the initialization, and so leads to poorer results when used jointly with init='nndsvd'. You may try init='nndsvda' or init='nndsvdar' instead.

warnings.warn("The multiplicative update ('mu') solver cannot update ")

```
In [75]: 1 df = pd.DataFrame(columns = ["Frobenius", "Kullback_leibler"],
2 index = ['random', 'nndsvd', 'nndsvda', 'nndsvdar'])
3 df = df.T
4 df['random'] = [train_score[0][0], train_score[0][1]]
5 df['nndsvd'] = [train_score[1][0], train_score[1][1]]
6 df['nndsvda'] = [train_score[2][0], train_score[2][1]]
7 df['nndsvdar'] = [train_score[3][0], train_score[3][1]]
8
```

```
In [76]: 1 df
```

Out[76]:

	random	nndsvd	nndsvda	nndsvdar
Frobenius	0.297805	0.267941	0.268071	0.267928
Kullback_leibler	0.345146	0.276521	0.303347	0.300328

After identifying the optimal parameters for the NMF, let us now join both the training and test set to increase the library of the NMF for better accuracy.

```
In [77]: 1 joint_df = pd.concat([df_train, df_test])
2 joint_df['review'] = joint_df['review'].astype(str)
3
```

```

In [82]: 1 vector = TfidfVectorizer( min_df = 55 , stop_words='english')
          2 joint_word_vec = vector.fit_transform(joint_df.review)
          3
          4 clf = NMF(init = 'random' , n_components=3, random_state=2,solver = 'mu' ,
          5
          6 clf.fit(joint_word_vec)
          7 w = clf.transform(joint_word_vec)
          8 h = clf.components_
          9
          10 tp=pd.DataFrame({'id':pd.DataFrame(w).idxmax(axis=1)})
          11 ## Number of categorizations
          12
          13 tp.value_counts()
          14
          15 accuracy_score(tp['id'].to_numpy() , joint_df['mod_rating'])

```

Out[82]: 0.40917312601423766

Modelling Summary / Discussion

We started off by dividing the train dataset into further train and test sets with the division set at 0.2 for test set.

NMF: The first model that we trained was the Non-Zero Matrix Factorization of the complete train data. The reason we took the entire data was it is not effected by any kind of overfitting since the data will not be based on any initial set outcome. After training the data and testing it on the train data we got an **accuracy of a mere 28%**. We then tested it on the actual test data provided as well which also gave an accuracy of 28%. This tends us to believe that the nmf is not really a good choice here with the lack of accuracy it is showing in both the datasets.

We later tried to optimize the NMF by optimizing the minimum acceptable word occurrence which came out to be at '55'. We also went over the various init types and loss function to identify the best parameters. From these, we found that the best **init was set at 'random'** with the **loss function at 'kullback leibler'**. **Once we set these new parameters and joined both the train the test set for enhanced accuracy, we got an accuracy of 41%** which is mark difference from our initial model but still is not accurate enough.

Random Forest: The Random forest model seemed to give the best result with an **accuracy of 80% on the train set** and **79% on the test set**. These accuracies are quite near which indicates that the model is substantially good fit without any overfitting taking place.

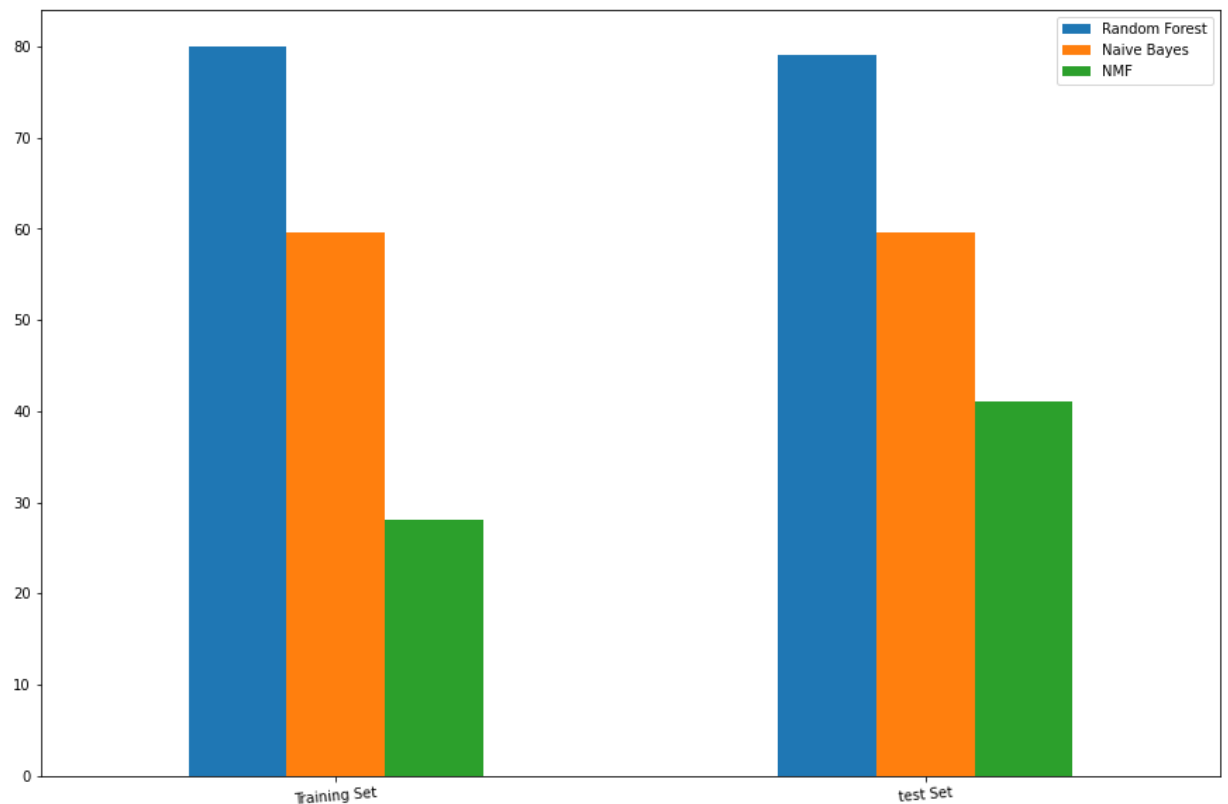
Naive bayes: The naive Bayes model seems to fall in between the Random Forest and NF in terms of performance with a **constant accuracy of 60% on both the train and test set**.

Results and Analysis


```
In [84]: 1 df_comp = pd.DataFrame(columns = ["Training Set" , "test Set" ],
2         index = ['Random Forest' , 'Naive Bayes' , 'NMF'])
3 df_comp = df_comp.T
4 df_comp['Random Forest'] = [80 , 79]
5 df_comp['Naive Bayes'] = [59.6 , 59.6]
6 df_comp['NMF'] = [28 , 41]
```

```
In [86]: 1 df_comp.plot(kind = 'bar', rot=5, fontsize=10, figsize=(15,10))
```

Out[86]: <AxesSubplot:>



The final selected model is Random Forest which seems to be performing the best in both the training data and the test data. The accuracy for each data is as follows:

Train Data : 80 %

Test Data : 79%

Lets try to optimize the Random Forest further,

```
In [89]: 1 max_feat = [ 'auto', 'sqrt' , 'log2' ]
2
3 rf_scores = []
4
5 for i in max_feat:
6     model = make_pipeline(TfidfVectorizer(stop_words='english'), RandomForest
7
8     model.fit(X_train , y_train)
9
10    pred = model.predict(X_test)
11
12    rf_scores.append(accuracy_score(y_test , pred))
```

```
In [90]: 1 rf_scores
```

```
Out[90]: [0.7953812771233726, 0.7964352138871668, 0.7830130192188469]
```

The default parameter setting of the Random Forest seems to be giving the best result and thus does not require any modification in terms of the features. There is one thing to note however in the above graph that NMF showed quite an improvement when the parameters were fine tuned. This if in conjunction with further improvement in the text processing can probably be improved to produce much better results.

Conclusion

In this project, multiple models were used to train the data both supervised and unsupervised learning methods. Both supervised learning models seemed to do better than the unsupervised learning method. However, it is important to note that the Matrix Factorization took a lot less time in model training and prediction as compared to the supervised learning method. The most important takeaway from this project was the need to properly clean the text. Many attempts were used to correctly identify the noise from the data to improve accuracy as well as time complexity. However, it is felt that the data still has a lot of room for improvement since many digits remained integrated with words which increased the time complexity and degraded the model overall accuracy.