

Ad Analysis using Machine Learning

Yashas C N¹, Kishan V², Pranay Reddy Juturu³ and Prof. Pavithra H⁴

^{1,2,3}Department of Computer Science and Engineering, R.V College of Engineering, Bengaluru, Karnataka, India

⁴Professor, Department of Computer Science and Engineering, R.V College of Engineering, Bengaluru, Karnataka, India

Abstract - In this paper we understand the concept of recommendation system and how beneficial it is to understand and analyze the advertisements as per user. As the audience keep evolving and there are different generation of users who are interested in different type of content. So, it is really important to recommend the right type of ad's in this evolving system so that the situation can be a win-win situation for the user and for the ad agencies as well. By using right type of recommendation system, this provides a solution of dynamically providing the information to the user and helps in promoting the product to the user and increasing their sales and business. In this paper, we'll discuss and analyze the type of approach used in solving this problem.

1. INTRODUCTION

Machine Learning is a field of computer science which uses statistical models to classify or predict the data. We are separating the data as test data and train data, where we use the train data to train the model and test data to test and see how that works and check the result meets the expectations and to improve the accuracy of the model.

In this paper, we have briefly explained a model which could recommend ads to the user in the most optimal way and recommends the content in the social media according to the user's interest and has a higher percentage of accuracy.

According to the recent study it is shown that, an average millennial spends around 8 hours a day on social media such (Facebook, Instagram etc..) and in YouTube around 4 billion videos are being watched every day. So social media, is a great source of obtaining the data. So, there are techniques incorporated by large corporations such as Facebook, Google who recommendation system to display the content according to the user's interest. But also, there are some websites by using the right type of recommendation system will be able to provide the right type of content to the right users.

2. RELATED WORK

There are lot of related work with respect to ad analysis using machine learning. One of the most used work is by video and sound analysis. Where the sound and video is segregated into frames and hashing is used so that there are no duplicates and naïve Bayes and sentiment analyzer is used for further recommendation. This method is not optimal as this method introduces to a lot of noise and accuracy will be low. The

approach used in this journal is feasible as the accuracy will be high and solves the above problem.

3. ANALYSIS METHODOLOGY

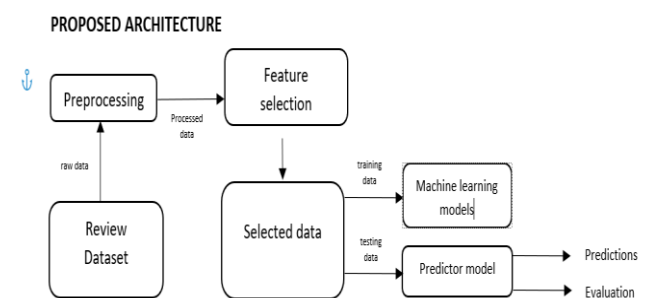


Fig 1.1

In the fig 1.1 the architecture of the proposed model has been explained. The following steps have been explained to achieve this process. They are: -

3.1 Data Visualization

In this method, the data is obtained from different e-commerce website. The data 100k rows and 25 columns and has different fields. The dataset which is chosen 80% test data and 20% train data. The data visualization is shown in the Fig. 1.2 and Fig 1.3.

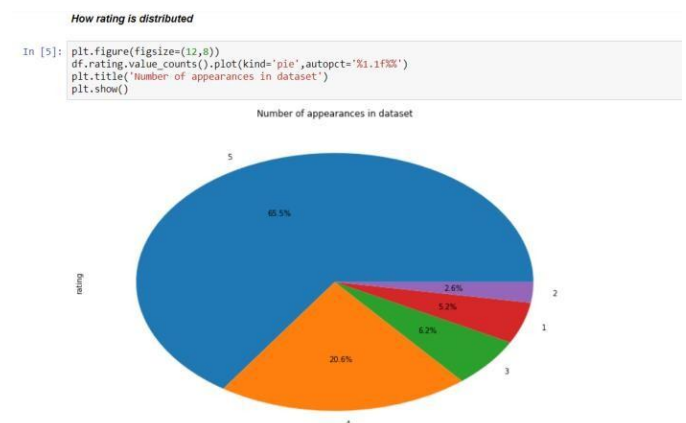


Fig 1.2. Distribution of Rating

The Year-wise distribution of products

Out[12]: <matplotlib.axes._subplots.AxesSubplot at 0x2454aab3b00>

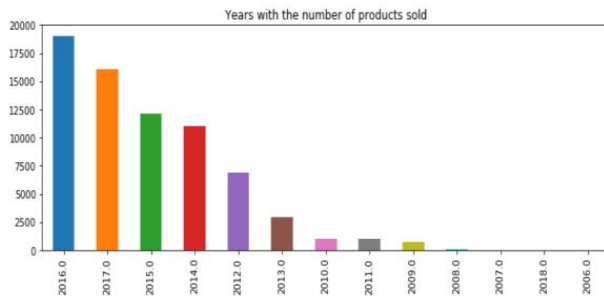


Fig 1.3. Year wise distribution of the product.

3.2 Data Preprocessing

Data pre-processing is the next step in this process as the data is obtained. There will be columns which have any kind of rows of content if there are rows with any type of non-relevant information, that particular row is dropped. And if there are NaN values, relevant information will be in those columns. The code snippet for this process is given below: -

Drop the columns which are not required and not useful for predictions

```
In [6]: drop_cols = ['Unnamed: 0', 'brand', 'categories', 'dateAdded', 'dateUpdated', 'keys', 'manufacturer', 'name', 'reviewsdate']
df = df.drop(drop_cols, axis=1)
df.head()
```

Fig 1.4 Dropping the values which are not relevant

Fill the NaNs with suitable values

```
In [17]: df['didPurchase'].fillna(True, inplace=True)
df['doRecommend'].fillna(True, inplace=True)
```

Fig 1.5 Fill Nan with suitable values.

3.3 Feature Selection

This is the next step in the process where various machine learning algorithms are incorporated to obtain the desired result: -

3.31 Gaussian Naïve Bayes

Naïve Bayes algorithm is based on Bayes' theorem with the assumption of independence with every pair of features. This requires extremely small amount of data to estimate the parameters but is fast compared to the more methods. However, it is a bad classifier but is used in real world application such as spam filtering and document classification.

Finding the optimal value using Gaussian Naive Bayes algorithm

```
In [71]: clf = GaussianNB()
clf.fit(X_train, y_train)
target_pred = clf.predict(X_test)
```

Find the accuracy score using Gaussian Naive Bayes Classifier

```
In [74]: accuracy_score(y_test, target_pred, normalize = True)
Out[74]: 0.9469755651010492
```

Fig 1.6 Accuracy using Gaussian Naïve Bayes

3.32 Natural language Processing (NLP)

Natural Language processing is a term which is used for automatic computation processing of human languages. This includes both algorithms that take human-produced text as input, and algorithms that produce natural looking text as input.

We have analyzed the review column for our dataset to understand customer review for each product. We have done preprocessing for removal special characters, digits, punctuations etc..

Removing Punctuation

Remove any punctuations present in the reviews for processing

```
In [111]: train['text'] = train['text'].str.replace("[^\w\s]","")
train['text'].head()
```

```
Out[111]: 0 i love this album its very good more to the hi...
1 good flavor this review was collected as part ...
2 good flavor
3 i read through the reviews on here before look...
4 my husband bought this gel for us the gel caus...
Name: text, dtype: object
```

Removal of Stop Words

Stop words are words which are filtered out before or after processing of natural language data (text). We remove some of the most common words—including lexical words, such as "and" in order to improve performance.

```
In [113]: from nltk.corpus import stopwords
stop = stopwords.words('english')
train['text'] = train['text'].apply(lambda x: " ".join(x for x in x.split() if x not in stop))
train['text'].head()
```

```
Out[113]: 0 love album good hip hop side current pop sound...
1 good flavor review collected part promotion
2 good flavor
3 read reviews looking buying one couples lubric...
4 husband bought gel us gel caused irritation fe...
Name: text, dtype: object
```

Fig 1.7 Removal punctuation and stop words.

This step (natural language processing) is achieved by 3 algorithms: -

- A. Bag of Words.
- B. TF-IDF.(Term Frequency, Inverse Document Frequency)
- C. Hashing

The code snippet for each of these algorithm is given below: -

Bag of Words

Bag of Words (BoW) refers to the representation of text which describes the presence of words within the text data. The intuition behind this is that two similar text fields will contain similar kind of words, and will therefore have a similar bag of words. Further, that from the text alone we can learn something about the meaning of the document.

CountVectorizer

The CountVectorizer provides a simple way to both tokenize a collection of text documents and build a vocabulary of known words, but also to encode new documents using that vocabulary.

```
[146]: from sklearn.feature_extraction.text import CountVectorizer
bow = CountVectorizer(max_features=1000, lowercase=True, ngram_range=(1,1), analyzer = "word")
train_bow = bow.fit_transform(train['text'])
train_bow
Out[146]: <70967x1000 sparse matrix of type '<class 'numpy.int64''
with 852741 stored elements in Compressed Sparse Row format>
```

Fig 1.8 Bag of Words using CountVectorizer

Term Frequency – Inverse Document Frequency (TF-IDF)

TF-IDF is the multiplication of the TF and IDF which we calculated above.

```
In [142]: tfidf = TfidfVectorizer(max_features=1000, lowercase=True, analyzer="word",
stop_words='english', ngram_range=(1,1))
train_vect = tfidf.fit_transform(train['text'])
train_vect
Out[142]: <70967x1000 sparse matrix of type '<class 'numpy.float64''
with 721460 stored elements in Compressed Sparse Row format>
```

Hashing with HashingVectorizer

The HashingVectorizer class implements this approach that can be used to consistently hash words, then tokenize and encode documents as needed.

```
In [162]: from sklearn.feature_extraction.text import HashingVectorizer
# create the transform
vectorizer = HashingVectorizer(n_features=20)
# encode document
vector = vectorizer.transform(train['text'])
vector
Out[162]: <70967x20 sparse matrix of type '<class 'numpy.float64''
with 648047 stored elements in Compressed Sparse Row format>
```

Fig 1.9 TF-IDF using TfidfVectorizer

```
In [163]: train['text'][:5].apply(lambda x: TextBlob(x).sentiment)
Out[163]: 0 (-0.09999999999999999, 0.575)
1 (0.0, 0.0)
2 (0.0, 0.0)
3 (0.014090909090909093, 0.6594444444444445)
4 (0.0, 0.0)
Name: text, dtype: object
```

Fig 1.10 Hashing using HashingVectorizer

3.33 Sentiment Analysis

Sentiment is like a combination of tone of voice, word choice, and writing style all rolled into one. Natural language with labels about positivity or negativity, we can develop agents that can learn to understand any sentiment. The code snippet for sentiment analysis is given below.

Sentiment Analysis

It is a process of computationally identifying and categorizing opinions expressed in a piece of text, especially in order to determine whether the writer's attitude towards a particular topic, product, etc., is positive, negative, or neutral.

```
In [148]: train['text'][:5].apply(lambda x: TextBlob(x).sentiment)
Out[148]: 0 (-0.09999999999999999, 0.575)
1 (0.0, 0.0)
2 (0.0, 0.0)
3 (0.014090909090909093, 0.6594444444444445)
4 (0.0, 0.0)
Name: text, dtype: object
```

Above, you can see that it returns a tuple representing polarity and subjectivity of each tweet. Here, we only extract polarity as it indicates the sentiment as value nearer to 1 means a positive sentiment and values nearer to -1 means a negative sentiment. This can also work as a feature for building a machine learning model.

```
In [150]: train['sentiment'] = train['text'].apply(lambda x: TextBlob(x).sentiment[0])
train[['text', 'sentiment']].head()
Out[150]:
```

	text	sentiment
0	album hip hop side current pop sound type list...	-0.100000
1	favor	0.000000
2	favor	0.000000
3	read review looking buying one couple lubrican...	0.014091
4	husband bought get u get caused irritation fel...	0.000000

Fig 1.11 Sentiment Analysis for the text

3.34 Stack Ensemble Model

Ensemble modeling is the process of running two or more related but different analytical models and then synthesizing the results into a single score or spread in order to improve the accuracy of predictive analysis and data mining applications. Stacking is a way of combining multiple models, that introduces the concept of a meta learner. Applying stacked models to real-world big data problems can produce greater prediction accuracy and robustness than do individual models.

Voting Classifier Ensemble

```
In [4]: X=dff[['didPurchase', 'rating']]
y=dff['doRecommend']
clf1 = LogisticRegression(random_state=1)
clf2 = RandomForestClassifier(random_state=1)
clf3 = GaussianNB()
clf4 = KNeighborsClassifier(n_neighbors=10)
ecf = VotingClassifier(estimators=[('lr', clf1), ('rf', clf2), ('gnb', clf3), ('knn', clf4)], voting='hard')
ecf.fit(X_train, y_train)
In [5]: ecf.score(X_test, y_test)
Out[5]: VotingClassifier(estimators=[('lr', LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
penalty='l2', random_state=1, solver='liblinear', tol=0.0001,
verbose=0, warm_start=False)), ('rf', RandomForest...wskl',
metric_params=None, n_jobs=1, n_neighbors=10, p=2,
weights='uniform'))], flatten_transform=None, n_jobs=1, voting='hard', weights=None)
In [6]: ecf.score(X_test, y_test)
Out[6]: 0.9484531242664663
Voting classifier gives an accuracy of 95%
```

Fig 1.12 Accuracy using Stacked Ensemble for Voting Classifier.

4. Results

Keeping the independent and dependent variables same across various algorithms, we found acceptable results of various algorithms used for our E-commerce dataset.

Dependent variable: doRecommend Independent variable: ProductId, didPurchase, Username, Rating.

Since we take into consideration various independent variables and find the correlation between them, our aim to predict the likelihood for the Advertisement being legitimate to the user and percentage if correctness.

Gaussian Naïve Bayes	94.69%
Natural Language Processing for Sentiment Analysis	Bag of Words TF-IDF Hashing, and sentiment analysis
Stacked Ensemble model	Voting Classifier ~94%

5. Conclusion

By building this analysis system, we are able to find products which are able to find the products which are similar and can be recommended to a set of users who have similar buying pattern through a set of ads generating based on the current analysis. After using various algorithms, we concluded that, Gaussian Naïve Bayes, NLP technique are most suitable for our dataset and for performing the required sentiment analysis serves as a way to let us know products demand and the kind of Ad based on them.

REFERENCES

- [1] R Vinit Kaushik, Raghu R, Maheshwar Reddy, Ankita Prasad, Sai Prasanna M S, "Ad analysis using Machine learning", PESIT – Bangalore South Campus.
- [2] Anitha anandhan, Liyana Shubu, Maizatul Akmar Ismail and Ghulam Mujtaba "Social Media Recommender System: Review and Open Research Issues", IEEE Transaction on knowledge and data Engineering, date of publication February 27,2018.
- [3] Ashutosh Bansal, Saleena B, Prakash B, "Using Data mining techniques to analyze the customer reaction toward social media."
- [4] Bengio, Samy; Goodfellow, Ian J.; Kurakin, Alexey (2017)," Adversarial Learning at Scala" Google AI.