



Semantic Content Recommendation System Using AWS Sagemaker

This project tries to develop a semantic content recommendation system using AWS Sage maker . Data has been trained and tested on a 20 newsgroups dataset and comparison of two supervised algorithms namely Multinomial Naive Bayes Classifier and Support Vector Machine(SVM) has been done on the basis of it.

AUTHORS

Eeshan Garg (11918)

AFFILIATIONS

Cluster Innovation Centre,
University Of Delhi, 110007, India
Mentor - Dr. Nirmal Yadav

INTRODUCTION

Information retrieval is the science of searching for information in a document, searching for documents themselves, or searching for metadata that describe data. This project combines the techniques of naive bayes and support vector machine for information retrieval. This approach uses bayes formulae for text document preprocessing and then uses support vector machine for classification

DATASET

List of the 20 newsgroups, partitioned respective to their subject matter is as follows:

- alt.atheism
- comp.graphics
- comp.os.ms-windows.misc
- comp.sys.ibm.pc.hardware
- comp.sys.mac.hardware
- comp.windows.x
- misc.forsale
- rec.autos
- rec.motorcycles
- rec.sport.baseball
- rec.sport.hockey
- sci.crypt
- sci.electronics
- sci.med
- sci.space
- soc.religion.christian
- talk.politics.guns
- talk.politics.mideast
- talk.politics.misc

From: lerxst@wam.umd.edu (where's my thing)
Subject: WHAT car is this?
Nntp-Posting-Host: rac3.wam.umd.edu
Organization: University of Maryland, College Park
Lines: 15

I was wondering if anyone out there could enlighten me on this car I saw the other day. It was a 2-door sports car, looked to be from the late 60s/early 70s. It was called a Bricklin. The doors were really small. In addition, the front bumper was separate from the rest of the body. This is all I know. If anyone can tell me a model name, engine specs, years of production, where this car is made, history, or whatever info you have on this funky looking car, please e-mail.

Thanks,
- IL

----- brought to you by your neighborhood Lerxst -----

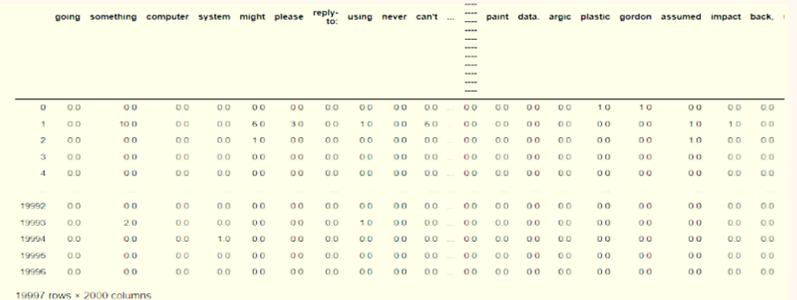
Example of a Newsgroup Document

DATA PREPROCESSING

Stripping the headers, footers and quotations

"A fair number of brave souls who upgraded their SI clock oscillator have\shared their experiences for this poll. Please send a brief message detailing\your experiences with the procedure. Top speed attained, CPU rated speed,\nadd on cards and adapters, heat sinks, hour of usage per day, floppy disk\functionality with 800 and 1.4 m floppies are especially requested.\n\nI will be summarizing in the next two days, so please add to the network\knowledge base if you have done the clock upgrade and haven't answered this\poll. Thanks."

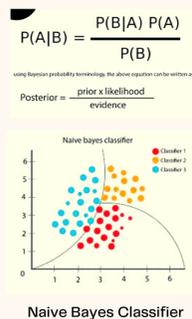
Data frame after Tokenization



MODELS

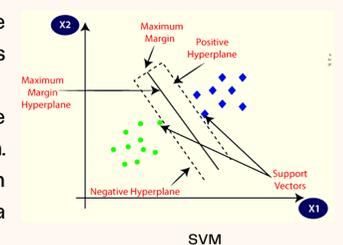
Multinomial Naive Bayes

This algorithm is a probabilistic learning method that is mostly used in Natural Language Processing. It predicts the tag of a text such as a piece of email or newspaper article. It calculates the probability of each tag for a given sample and then gives the tag with the highest probability as output.



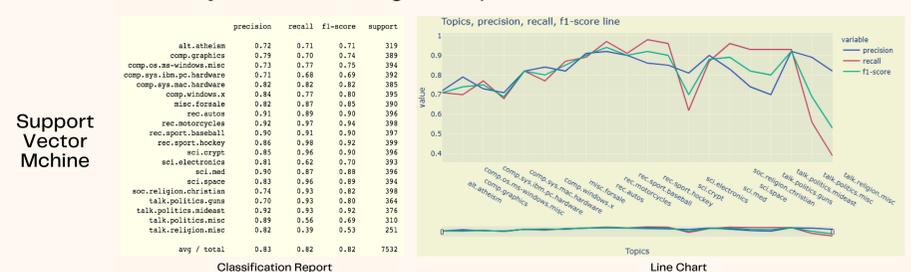
Support Vector Machine

The objective of this algorithm is to find a hyperplane in an N-dimensional space that distinctly classifies the data points. To separate the two classes of data points, there are many possible hyper planes that could be chosen. Our objective is to find a plane that has the maximum margin, i.e. the maximum distance between data points of both classes.



RESULTS

The following classification report and line graph show performance of each classifier Multinomial Naive Bayes and SVM, in regards to precision, recall scores and F1-Score.



CONCLUSION

From this study, we found that Naive Bayes does very well in text classification. It should be noted that SVM does achieve very good scores as well. In short, given a text classification problem, any of the two could be used without a big compromise on the classification accuracy. But if we want to do it in less training time, then it would be advised to use Naive Bayes instead of SVM.

FUTURE DIRECTIONS

It would be fascinating to explore what makes every one of the classifiers proceed as they do, may it be from the scores every classifier accomplishes in characterizing the text, or the assets it consumes; time, and computer memory during training and testing. Research can be done on the effect of component extraction and portrayal on the performance of every classifier. For example Multinomial Naive Bayes in this study performs far superior to SVM on account of the BOW approach.

REFERENCES

- <https://aws.amazon.com/getting-started/hands-on/semantic-content-recommendation-system-amazon-sagemaker/>
- https://medium.com/@siyao_sui/nlp-with-the-20-newsgroups-dataset-ab35cd0ea902
- <https://github.com/tanishq9/Text-Classification-20-Newsgroups>
- <http://archive.ics.uci.edu/ml/datasets/Twenty+Newsgroups>