



# Sentiment analysis of movie reviews using various algorithms

Mood Jyothi

Asst. Prof. Department of CSE, SVS Group of institutions, Bhemmaram, Warangal, Telangana

## To Cite this Article

Mood Jyothi. Sentiment analysis of movie reviews using various algorithms. International Journal for Modern Trends in Science and Technology 2022, 8, pp. 265-269. <https://doi.org/10.46501/IJMTST0802043>

## Article Info

Received: 18 January 2022; Accepted: 22 February 2022; Published: 25 February 2022

## ABSTRACT

*There are two major advancements in the state of viticulture technology research that are made by this effort. To begin, we take a comprehensive look at the evolution and current landscape of AV, IG, and ML in the wine business. By analyzing case studies from a variety of domains, including as crop yield estimation, vineyard management and monitoring, disease detection, quality evaluation, and grape phenology, we present a succinct review of current breakthroughs in vision systems and approaches. In this paper, we focus on how computer vision and machine learning might improve contemporary vineyard management and verification processes. In the second part of the study, we provide the brand new Grape CS-ML Database, which includes images of different grape varieties at different stages of ripeness along with the corresponding ground truth data (such as pH, Brix, etc.) gleaned via chemical analysis. The major purpose of this database is to motivate researchers in computer vision and machine learning to develop practical solutions for use in smart vines. By contrasting white and red cultivars across a variety of machine learning methods and color spaces and giving a set of assessment data, we demonstrate the database's potential for a color-based berry recognition application. The research concludes by discussing some of the challenges that must be overcome in the future to fully implement this technology in the viticulture sector.*

**Keywords:** Viticulture, computer vision, machine vision, visual computing, image processing, machine learning.

## 1. INTRODUCTION

The Internet is revolutionizing how people learn and work, but it also poses serious security risks as more and more of our lives move online. Recognizing network assaults, especially novel ones, is a pressing issue that has to be addressed immediately.

Computers, networks, programs, and data are all examples of digital assets, and cyber security is the practice of keeping them safe from attack. Within a network security system, both the network and the computers themselves are protected. These networks include several layers of protection, including firewalls, anti-virus software, and intrusion detection systems (IDS). Without an IDS, detecting unauthorized data

access, duplication, modification, or deletion might be challenging.

The term "security breach" can refer to both external and internal intrusions. The most frequent method of network analysis employed by IDSs is misuse-based, also known as signature-based, analysis, followed by anomaly-based and hybrid methods. The purpose of misuse-based detection systems is to spot malicious behavior by comparing it to known attack patterns [3]. When used for specific attacks, they won't produce too many false positives. The rules and signatures stored in the database, however, must be updated often by hand by administrators. Technology misuse makes it

impossible to monitor for new attacks or stop them in their tracks (zero days).

The normal functioning of a network or system is analyzed by using an anomaly detection method. One of their main selling points is that they can detect attacks that haven't been seen before. Additionally, normal-use profiles can be tailored by machine, program, or network, making it more difficult for attackers to predict what they will get away with. Signatures for abuse detectors can be defined using data on which anomaly-based strategies trigger an alert (new attacks). The primary disadvantage of anomaly-based techniques is that it is possible for previously unknown system behaviors to be labeled as abnormalities.

Hybrid detection combines misuse detection and anomaly detection[4]. By doing so, we can improve our detection of known incursions while decreasing the number of false positives caused by unknown attackers. The bulk of ML/DL systems are hybrid setups.

**There are 2 WORKS THAT CONNECT WITH THIS ONE.**

Sentiment analysis, with its inherent complexities, has drawn a huge number of academics because of the power of machine learning methods.

Sentiment analysis and a deep dive into machine learning techniques pave the way for a state-of-the-art intelligent system that can prove its AI prowess.

Yet, researchers may find it difficult to determine the best machine learning approach, and making the wrong choice might result in misleading findings and poor model performance.

This motivated us to investigate the relative merits of various machine learning techniques for sentiment analysis.

Only supervised machine learning approaches were considered, and we looked for ways to compare them in as many ways as possible.

### **2.1 LIMITATIONS OF EXISTING SYSTEM:**

When labels are absent from the training data, an unsupervised learning model must rely on its own intuition to correctly identify the data. The Teaching Strategy of Reinforcement In this case, the model is employed to guide official action. What to do in light of each fresh observation made in the external world. The model acquires knowledge by its experiences with the world, which are shaped by the data it collects about

the world. The Transduction Mechanism: It's quite similar to supervised machine learning; the main difference is that instead of building a function, it predicts future outcomes based on historical data.

### **3. PROPOSEDSYSTEM**

The methodology that was devised and implemented in this research is outlined below. Film Review Ratings Derived from Data: The collecting of data is the first step in any Sentiment analysis study, and there are already a number of publically available online. Information Cleaning The movie review data set contains characters, digits, special characters, and unrecognized characters. To protect our classifier from this potential risk, we employed a data set cleaning procedure following data collection. The method we used to eliminate unnecessary data from our databases. The next step, categorizing the reviews already in the databases, can now get under way.

Supervised machine learning methods can be applied to data that has already been classified into the available classes. So, it is necessary to classify reviews as either positive or negative based on their content. The effectiveness of a classifier is related to the time and effort spent training it and the quality of the data used to educate it, just as it is with any learning method. To keep things straightforward, it's common to employ a 70% training/30% testing split of the data set.

Training Data Sets: The Method of Putting the Model Through It The success of sentiment analysis relies heavily on the techniques employed during the training phase, as here is where the classifier is developed. So, it is crucial to provide the correctly labelled data as its input while training the model, and we must also take care that the training process is not stymied by too much data. This is because if the classifier is trained for too long, its accuracy may degrade. The training data set we constructed in the previous step contains 70% of all available data sets, therefore we'll use those.

### **3.2 ADVANTAGES OF PROPOSEDSYSTEM:**

Sentiment analysis of movie reviews is performed in a novel and easy approach using seven promising supervised machine learning techniques.

In the future, it is hoped that machines will be so intelligent that they can solve any problem without human intervention.

Hence, we expect robots to perform, and we teach them using machine learning techniques. Machine learning (ML) is a subfield of computer science concerned with enabling machines to learn new tasks without being explicitly programmed to do so. The goal is to examine the writing styles of movie critics to see if they have a generally positive or negative outlook.

This is the case because audiences will only spend their time and money on movies that satisfy their expectations.

People used to look to reviews written by people who had already watched the film in question to help them make a decision about whether or not to watch it.

Because viewers may not read all of the comments if there are too numerous, it is important to show or recommend the percentage of negative and positive feedbacks about the relevant movie. Which will help the viewer save time and make a more informed decision? The objective of this research is to organize the many forms of film criticism that have been written.

It may also show a breakdown of the number of positive and negative comments for the user's current movie option (s) in order to help them make an informed decision.

### 3.3 ARCHITECTURE

A CNN starts with a convolution layer and a subsampling layer, and then it can progress to a fully connected layer. Images in the form  $(m, m, r)$ , where  $(m, m)$  and  $(r, r)$  are the height and width of the picture and the number of channels, respectively; for instance, an RGB image has  $r=3$  and  $(r, r) = 3$ . The filters (or kernels) used in the convolutional layer, denoted by the notation  $kk$ , will have the dimensions  $n$  by  $n$  by  $q$  by  $q$ , where  $n$  is less than the size of the picture and  $q$  is equal to or fewer than the number of channels,  $r$ . Using convolution, each filter of size  $mn+1mn+1$  processes the image to produce a feature map of size  $kk$ . For subsampling the input maps, the standard practice is to use mean or max pooling over  $p \times p \times p$  continuous portions, where  $p$  is between 2 and 5 for small images (like MNIST) and 2 to 10 for large inputs. Before or after

the subsampling layer, each feature map is subjected to an additive bias and sigmoidal nonlinearity. The convolutional and subsampling sublayers that make up a CNN layer are seen in the image below. Units of the same color are interchangeable with one another.

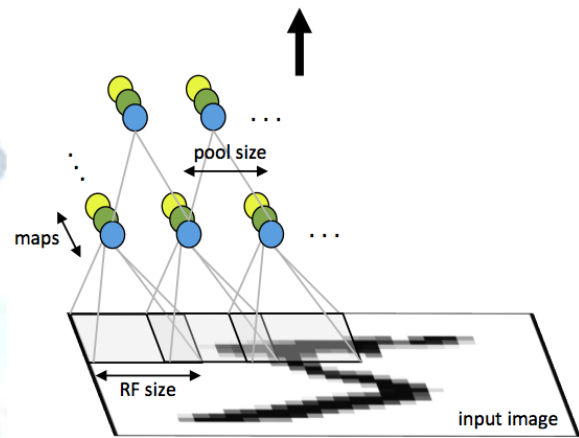


Figure 1:CNN

In Fig. 1 we see the first layer of a neural network with pooled convolutions. Adjacent units' weights are linked together, and the filter maps are shown by the units' colors.

### 3.4 Image processing

As has been previously discussed, the most general and potent feedforward neural network architecture is multilayer perceptrons, which are organized in layers such that each neuron within a layer receives a copy of all the outputs from the previous layer as its input. This type of model is particularly well-suited to the problem of learning from a small set of (more or less) random parameters.

Nevertheless, consider what happens to the model's parameter count (weights) when it is given unprocessed image data. Each neuron in the first hidden layer adds about 3000 more parameters to the model if we regard each channel of each pixel as an independent input to an MLP; CIFAR-10, for example, contains 32323 colored images. Long before the images reach a scale that most people would like to deal with, the situation gets uncontrolled as their size grows.

A common remedy, when using MLP, is to downsample images to a more manageable size. It would be great if we could apply some useful processing on the image before downsampling it, without needing to drastically expand the number of parameters, because doing so would prevent us from losing too much detail.

### 3.5 Convolutions

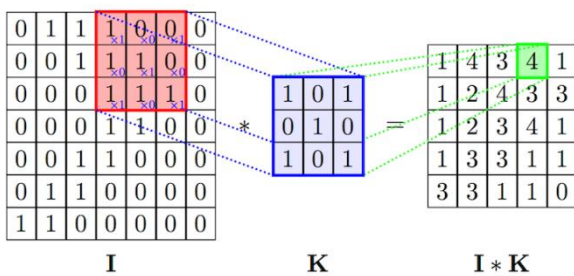
The structure of the information conveyed inside an image can be exploited in a highly efficient manner. Pixels that are physically near together are more likely to "cooperate" on forming an important aspect of the image than those that are farther apart. The same holds true for minor features: it doesn't matter where in the image they are located if they are found to be essential in determining the image's label.

Here, we make use of the convolution operator. Specifically, given a two-dimensional image  $I$  and a tiny matrix  $K$  of size  $h \times w$  (known as a convolution kernel), we calculate the convolved image  $I * K$  by repeatedly superimposing the kernel on top of the image and recording the sum of elementwise products between the image and the kernel:

$$(I * K)_{xy} = \sum_{i=1}^h \sum_{j=1}^w w_{Kij} \cdot I_{x+i-1, y+j-1}$$

(The precise definition really calls for inverting the kernel matrix first, but this step is unnecessary for ML.)

To better understand how the preceding formula and convolution (using two different kernels) function as an edge detector, consider the accompanying diagrams:



**Figure 2: Convolutional and pooling layers**

The convolution operator is the foundation of the convolutional layer and hence one of the most important parts of a CNN. Convolution of the output images of the previous layer with the given number of kernels,  $K$  (along with additive biases,  $b$  per kernel), defines a layer's operation (one per each output image). The next stage is to give each pixel in the final images an activation function. With a convolutional layer with  $d$  input channels, the final formula for a single picture channel's output (using kernel  $K$  and bias  $b$ ) looks like this: (for example, red, green, and blue in the input layer).

$$\text{conv}(I, K)_{xy} = \sigma(b + \sum_{i=1}^h \sum_{j=1}^w \sum_{k=1}^d w_{Kijk} \cdot I_{x+i-1, y+j-1, k})$$

It's important to keep in mind that, just like the weights of an MLP, the kernels can be learned from a given training dataset through gradient descent, as all we're doing here is adding and scaling the input pixels. Although an MLP could replicate a convolutional layer, it would need significantly more time (and data) in training to learn to operate in a manner that is even remotely close to that of the original.

Finally, it is worth noting that a convolutional operator is not limited to data that is structured in only two dimensions; in fact, most machine learning frameworks (Keras included) will supply you with pre-built layers for 1D and 3D convolutions.

Even though the number of parameters in a convolutional layer is much smaller than in a fully connected (FC) layer, more hyperparameters (parameters whose values must be determined before training begins) are included.

In specifically, the following hyperparameters should be chosen for use within a single convolutional layer: The height and width of each kernel; the number of kernels to convolve with the output of the previous layer (depth) (height and width). The "stride" refers to the amount by which the kernel is shifted between iterations in order to determine the next pixel in the output. To indicate the overlap between individual pixels in the output, this is often set to 1, which matches the previously provided formula. Keep in mind that as your stride length increases, your output sizes will decrease. Convolution by any kernel larger than 1111 will lower the size of the output image, even though it is often preferred that the input and output images be the same size. Here, the image has been suitably padded with zeros around the edges. This is often referred to as "same padding," in contrast to "valid" (no) padding. While any amount of padding is theoretically possible, the most popular choices are same padding and valid padding.

Although recent advances on all-convolutional networks are promising, in most cases a CNN's sole purpose is to downsample an image so that it may be processed by an MLP, and convolutions are employed for this purpose.

Since the pooling layer mixes multiple separate but similarly sized regions of the image, it is a popular method for downsampling photos (2222). Max-pooling is the most popular way for doing this aggregation, and

it entails picking the pixel with the highest value within a particular chunk. The 2x2 max-pooling algorithm is shown in diagram form below.

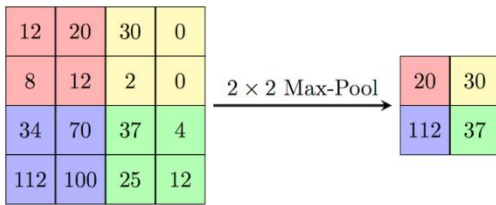


Fig 3:Max pooling

#### 4. CONCLUSION

Here, we present a novel, straightforward approach to sentiment analysis of film reviews using seven promising supervised machine learning algorithms. The data suggests that linear SVC/SVM is the best classifier for reliably classifying a huge volume of data, like movie reviews. We hope to investigate its effectiveness when used with huge datasets using unsupervised and semi-supervised machine learning techniques in the future.

#### Conflict of interest statement

Authors declare that they do not have any conflict of interest.

#### REFERENCES

- [1] Arturo Aquino, Maria P. Diago, Borja Millan, Javier Tard Aguila, "A new methodology for estimating the grapevine berry number per cluster using image analysis", Biosystems Engineering, Vol. 156, pp. 80-95, 2017.
- [2] Tardaguila, J., Diago, M.P., Millán, B., Blasco, J., Cubero, S. García-Navarrete, O.L., and Aleixos, N., "Automatic estimation of the size and weight of grapevine berries by image analysis", CIGR- AgEng, 2012.
- [3] [http://wcigr.ageng2012.org/images/fotosg/tabla\\_137\\_C1735.pdf](http://wcigr.ageng2012.org/images/fotosg/tabla_137_C1735.pdf)
- [4] G. M. Dunn, S. R. Martin, "Yield prediction from digital image analysis: a technique with potential for vineyard assessments prior to harvest.", Australian Journal of Grape and Wine Research, 10, pp.196-198, 2004.
- [5] Chamelat, R., Rosso, E., Choksuriwong, A., Rosenberger, C., Laurent, H., and Bro, P., "Grape Detection by Image Processing", Proceedings of IEEE 32nd Annual Conference on Industrial Electronics (IECON 2006), Paris, France. pp. 3697-3702, 2006.
- [6] Rabatel, G., &Guizard, C., "Grape berry calibration by computer vision using elliptical model fitting", 6th European Conference on Precision Agriculture (ECPA 2007), Skiathos, Greece. pp. 581-587, 2007.