

model considered has 16 input nodes with 10 nodes in hidden layers and 3 nodes in output layer. In case of hidden layer, different number of nodes are used for experiments.

There are two types of neural networks where the speech signals are considered such as the time delay network and length sequence which is processed at every step. The method followed is specified the inputs are shifted right for every unit delay line which links each set of inputs u units to the right adjacent and then the next input pattern is fed to the left position. perceptron for time sequence process that helps in changing the spatial sequence over corresponding units. The input layer will settled as man the recurrent neural networks. The time delay neural networks represents an effort to train a static multi-layer.

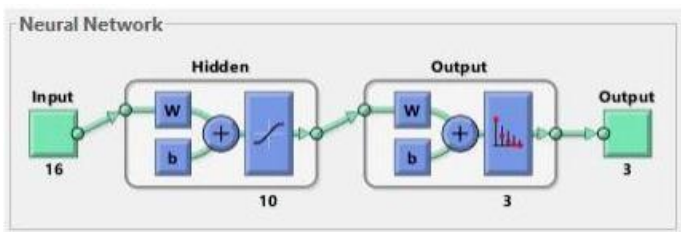


Figure 2. Neural network model

The recurrent neural networks help in providing a powerful extension of feed forward models by inserting the connections between the arbitrary pairs of units independently. The recurrent connections are reliable for an evolution of time of the internal state of the network.

3. PROPOSED ALGORITHM AND ITS WORKFLOW

In this research, the proposed a video classification method which uses audio as well as visual features of input video data for classification and the framework starts with the collection of sample video and is processing through pre-processing technique in which video and audio information is separated from the input video. The video data is converted into a set of frames followed by segmentation and extraction of critical features from both audio and video data. Later, feature level fusion process is initiated for concatenating amongst the audio, and video features and classification is achieved through a machine learning algorithm. The block diagram of the system being proposed is shown in Figure 3.

The workflow of the proposed system is illustrated in Fig. 4. Firstly, the video is extracted from the online website and is processed through the audio-video separation stage where the visual images and the audio file is extracted from the video. Initially, the video is converted into individual frames and processed through the key feature extraction stage, which is used to distinguish between the individual images and is processed through dimensionality reduction technique.

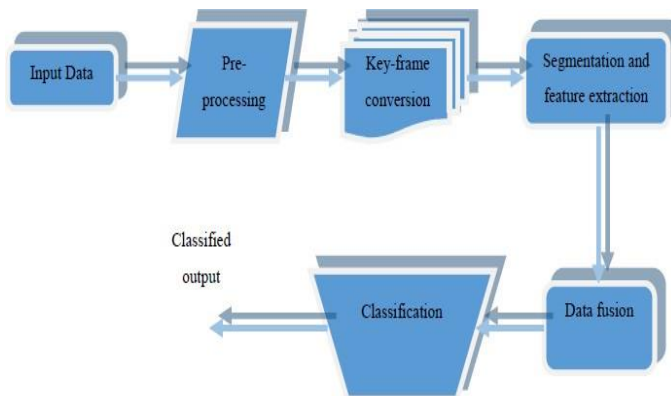


Figure 3. Block diagram of the proposed system

The histogram-based segmentation process is considered to enhance the quality of the image and to obtain the digitalized form.

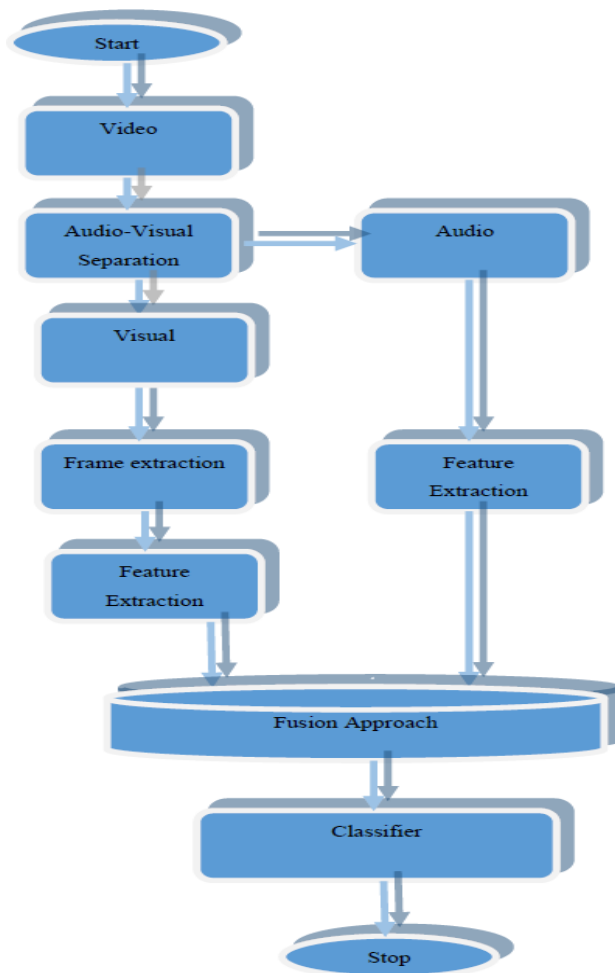


Figure 4. Workflow of the proposed system

Later, GLCM based features extraction technique is employed for extracting the useful features from video and zero forcing equalizer with Mel- Frequency Cepstral Coefficient (MFCC) for extracting the useful features. These features are later combined in the fusion stage and then processed through ANN_HMM stage to assess the performance of system. The in-depth description of the individual algorithm will be explained in the following stages as follows.

- In this research, the input data is processed through the keyframe extraction stage which is done on the basis of feature information. Color is the frequently accepted and essential feature for analyzing the input video.
- Features are stated as a function of individual or more number of measurement units in which quantifiable amount of data about the object is specified along with its significant characteristics of the video.

4. RESULTS AND DISCUSSION

The input videos considered are music, Sports and traffic. The input video is divided into 750 frames and those images are extracted using the feature extraction. The software used is MATLAB and following are the results parameters:

4.1 Confusion matrix

The confusion matrix is used in the machine learning to specify the complications of the statistical classification which is also called as error matrix and helps in visualizing the performance of a classification model on a set of test data for which the true values are known. In this the rows of the matrix represents the instances in a class while the columns represent the instances in an actual class. The

confusion matrix model for the proposed classification technique is shown on next page as fig 5.

The plot shows the variation of training and testing data. This is a simple technique for the summarization of the performance of the classification algorithm. By calculating the confusion matrix, the idea of right types of errors appearing in the model is identified. The basic terms used in the confusion matrix are:

- True positives (TP): The cases which predicted yes and they do.
- True negatives (TN): The case has predicted no and they don't.
- False positives (FP): The case was predicted yes but they actually don't.
- False negative (FN): The case was predicted no, but they actually do.

From the confusion matrix the list of rates that are predicted are,

- Accuracy
- Misclassification rate
- True positive rate
- False positive rate
- Specificity
- Precision
- Prevalence
- Null error rate
- Positive prediction value
- ROC curve

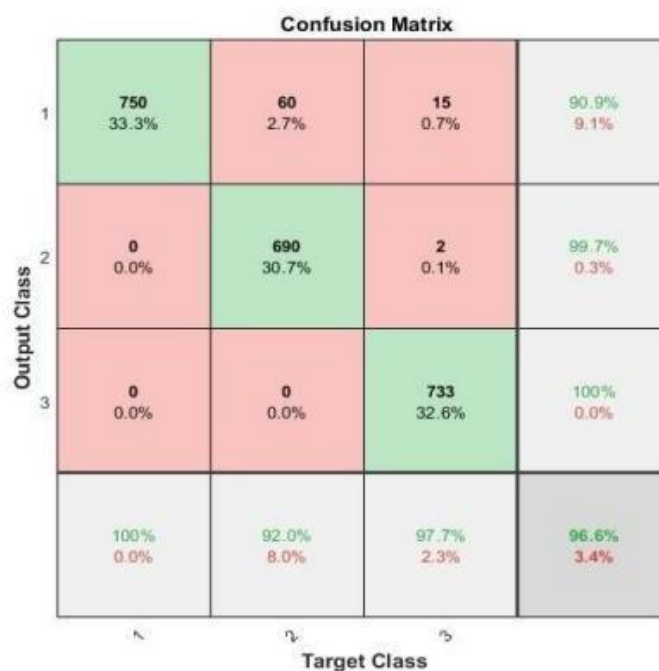


Figure 5. Confusion matrix for the proposed mode

4.2 ROC curve

ROC curve gives the ability to assess the performance of the classifier over the entire operating range. The neural networks are the classifier where these curves are used as shown in fig 6..

The plotting of the true positive rates against the false positive rates are performed at various threshold values. The above figure shows the graph of the three video signals and is classified as class 1, class 2, class3. The comparison of the two operating characteristics is performed as the criterion changes. The above figure shows that the class 1 has accuracy of almost 100% and the class B has accuracy of 90% whereas class c has accuracy of 99%.

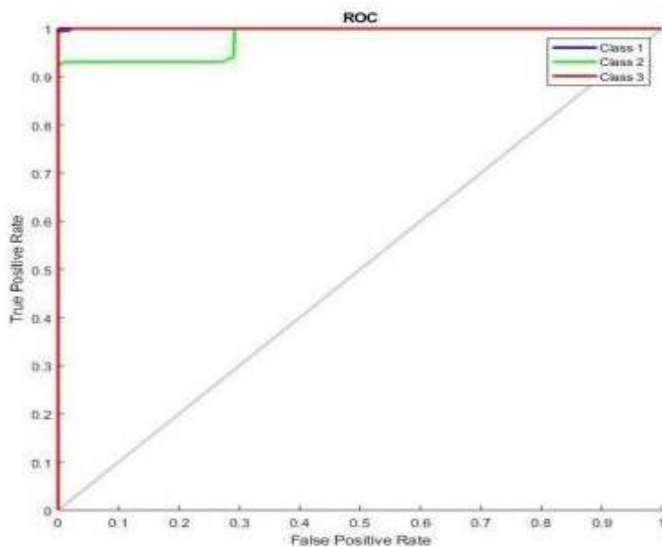


Figure 6. ROC output curve

4.3 HMM output

The below figure7 shows Hidden Markov model graph where the parameters such as emission and transition are compared. As the transition rate increases the emission rate increases gradually and shows a sudden increase at some point and remains constant after few transitions.

The table 1.1 shows the complete values which is obtained after the simulation process and gives the accuracy achieved. In the above table, considering Music Video-2, Sports video-2, and Traffic Video-2 we achieved 96.5777 % accuracy. So we conclude that 96.5777 5% data is correctly classified and the Sensitivity is 100%.

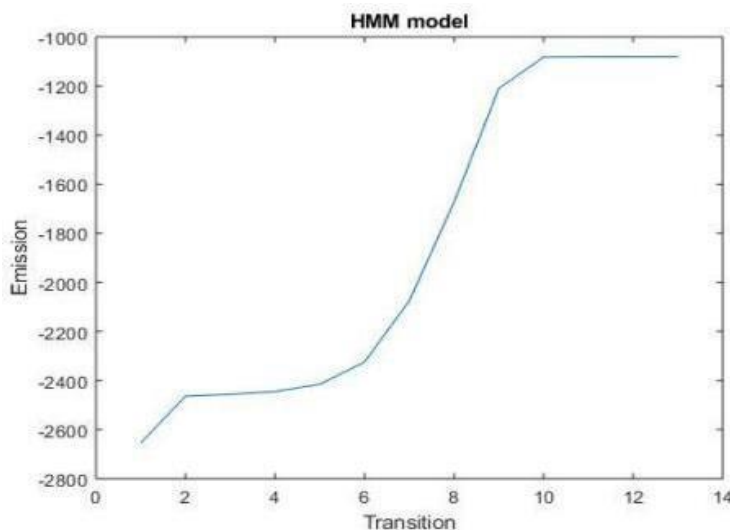


Figure 7. Output of HMM

Table 1. Performance Evaluation Table

| Video Name | M2 S2 T2 (music video 2, Sport video2, Traffic Video 2) | M3 S3 T3 (music video 3, Sport video3, Traffic Video 3) |
|-----------------|--|--|
| Accuracy | 96.5777 | 88.5777 |
| Sensitivity | 100 | 87.6000 |
| Specificity | 94.8666 | 89.0666 |
| Precision | 90.6892 | 80.0243 |
| Recall | 100 | 87.6000 |
| Freq_Measurment | 95.1173 | 83.6409 |
| Gmean | 97.3995 | 88.330 |

If we consider Music Video-3, Sports video-3, and Traffic Video-3 we can achieve 88.5777 % accuracy. So we can conclude that 88.5777% data is correctly classified. Sensitivity is 87.6000. Thus the results obtained shows that the accuracy obtained in different types of video has maximum accuracy of 97% where the video is classified.

5. CONCLUSION

With an increasing number of users and modes of communications, users are spending an amount of time on videos. However, it is quite difficult for human beings to categorize and caption too many videos. So, videos which we are going to watch require advanced techniques to do video classification and captioning. Also, in a few cases, we need to block some contents. Like, parents don't want their kids to watch violent or abusive content on the internet and for this purpose, we need such an advanced video classification technique which can find and block that unwanted content. In this work, we have explored methods and architectures to understand videos. The automatically categorization and caption are helpful for users to have better experiences when watching videos. The aim of this work is to know how to categorize and caption the video automatically. We have proposed an enhanced HMM- ANN based classification of video recordings with the aid of audio-visual feature extraction. The results evaluated and analysed shows the better accuracy when compared with the traditional. The obtained accuracy of 97% by the ANN-HMM algorithm shows that the classifier used for classifying the type of video based on the categories used.

References

- [1] M. Dawood and M. Ghanbari, "Scene content classification from MPEG coded bit stream," in Proc. 1999 IEEE 3rd Workshop on Multimedia Signal Processing, pp. 253–258, Copenhagen, Denmark, September 1999.
- [2] Aizawa, K., Nakamura, Y., & Satoh, S. I. (Eds.). (2004). *Advances in Multimedia Information Processing-PCM 2004: 5th Pacific Rim Conference on Multimedia*, Tokyo, Japan, November 30-December 3, 2004, Proceedings (Vol. 3332). Springer.
- [3] Anjum, N., & Cavallaro, A. (2010). Trajectory clustering for scene context learning and outlier detection. In *Video search and mining* (pp. 33-51). Springer, Berlin, Heidelberg.
- [4] Araújo, C. S., Magno, G., Meira, W., Almeida, V., Hartung, P., & Doneda, D. (2017, September). Characterizing videos, audience, and advertising in Youtube channels for kids. In *International Conference on Social Informatics* (pp. 341-359). Springer, Cham.
- [5] Bahatti, L., Bouattane, O., Echibat, M. E., & Zaggaf, M. H. (2016). An Efficient Audio Classification Approach Based on Support Vector Machines. *international journal of advanced computer science and applications*, 7(5), 205-211.
- [6] Barbancho, A. M., Tardón, L. J., López-Carrasco, J., Eggink, J., & Barbancho, I. (2015). Automatic classification of personal video recordings based on audiovisual features. *Knowledge-Based Systems*, 89, 218-227.
- [7] Camastra, F., Vinciarelli, A., & Yu, J. (2009). Machine learning for audio, image and video analysis. *Journal of Electronic Imaging*, 18(2), 029901-029901. 50
- [8] Casey, M. A. (2002). *Sound Classification and Similarity. Introduction to MPEG-7: Multimedia Content Description Interface*, 309-317.
- [9] C. W. Ngo, T. C. Pong, H. J. Zhang, and R. T. Chin, "Motion characterization by temporal slice analysis," in Proc. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, vol. 2, pp. 768–773, HiltonHead Island, SC, USA, June 2000.
- [10] Chougule, S. V., & Chavan, M. S. (2014). Channel Robust MFCCs for Continuous Speech Speaker Recognition. In *Advances in Signal Processing and Intelligent Recognition Systems* (pp. 557-568). Springer, Cham.