
MLG-42 Real Time Sentiment Analysis

1. Devashish N Mehta - dnmehta@iitk.ac.in
2. Nikita Awasthi - anikita@iitk.ac.in
3. Shivam Utreja - sutreja@iitk.ac.in
4. Mrinaal Dogra - mrinaald@iitk.ac.in
5. Mohd. Abbas Zaidi - mzaidi@iitk.ac.in

Problem Statement

1 We aim to develop a user-end application to get time-mapped viewer sentiment
2 while watching a video from the following set of 6 different emotions:

$$E = \{\text{neutral, happy, sad, angry, surprised, fearful}\}$$

3 The input will be the real time viewer's webcam feed while watching the video
4 and the output will be an emotion specified at regular time intervals throughout the
5 video.

6 1 Problem motivation

7 • Widespread need of genuine user feedback

8 Users are exposed to a lot of visual content online on a daily basis. This could be in the
9 form of product advertisements, youtube videos, movie trailers etc. The creators of these
10 will be benefited immensely if given the user's emotional feedback on being exposed to
11 this content. The recommendation systems will also be able to use this data in order to
12 make better recommendations in the future to the same user. This will hence, improve the
13 experience of the consumers, and also provide with constructive feedback to the creators.

14 • Drawbacks of the current user review methods

15 Presently, the only methods by which creators can get reviews of their content online
16 are via comments, reaction buttons, reviews, feedback forms etc.. There are two major
17 disadvantages of all these methods. Firstly, only a small percentage of users opt to fill
18 these out (which itself is a biased set of users in the first place), of which, an even smaller
19 percentage coincides well with the true feelings of the reviewer, since the reviews can
20 also be affected by external factors[1]. Secondly, these responses are not real-time, but an
21 overall summarized response of what the user thinks he felt of the content as a whole, after
22 watching it completely. This provides a very crude feedback to the creators as well as the
23 recommendation systems to be of any significant use.

24 • Advantages of Realtime Sentiment Analysis

25 Ideally, the recommendation systems and creators will be benefitted the most if they could
26 know the responses of the users exposed to a given piece of content. These are the pure
27 sentiments that each specific user had exhibited throughout the duration of the video. This
28 data set would be a lot richer, being the time-mapped review of the video. It would also be
29 a lot larger, and less biased as it would be inclusive of all the users watching that specific
30 content. It would also be a highly accurate estimate of the true user emotions, as opposed to
31 what the user chooses to mention in his/her written review.

32 2 Previous Work

- 33 • A number of features are used to analyse user sentiment including voting, rating, etc.
34 Comments are one of the most informative features and they try to analyse user comments
35 on video to attach a sentiment to the video. [2]
- 36 • A framework for classifying images according to high-level sentiment which subdivides the
37 task into three primary problems: emotion classification on faces, human pose estimation,
38 and 3D estimation and clustering of groups of people. [3]
- 39 • A project was made using deep learning techniques to predict the emotion depicted by an
40 image. [4]
- 41 • Currently, sentiment analysis on Youtube profiles of the person hosting the video or the
42 comments by users to identify possible polarisation by the video. However as mentioned in
43 [5], to gauge the extent of radicalization possible, the cue is basically text.
- 44 • For social media analysis using tweets for analysing user feedback and sentiment, the
45 existing work predicts the emotion labels as positive, negative or neutral. This can limit
46 the amount of information that can be extracted from the content for feedback to content
47 creators (videos in case of Youtube). This might be disadvantageous for content creators
48 who wish to maximise the reach of their content by analysing the highs and lows of the
49 video mapped to time frame. [6] [7]

50 2.1 Further Discussion and Problems

- 51 • Most of the existing methods rely on some sort of feedback received from a user. As we
52 discussed above it cannot be trusted as a reliable source of information about the actual
53 feelings of an individual. Therefore using tweets, comments or reacts as the feedback
54 may not actually serve the purpose when the analysis is based to gauge user feedback to a
55 particular commodity/ product.

56 3 Novel Contribution

57 With a greater participation of individuals on social media, coporates have also shifted to social media
58 for popularisation of their products. Revenue generation is directly linked to the reach of any product.
59 With advertisements shifting to a video based platform, it is important that the content creators receive
60 a more detailed feedback on their content. Our application proposes to bridge that gap between the
61 producers and the consumers. With a time-mapped sentiment feedback to provide better analysis, the
62 producers have an opportunity to identify their weak points. The current feedback systems lag behind
63 in this regard. Most of the review systems rely on cumulative feedback at the end of the video in the
64 form of comments or reactions. We create an end-to-end application that takes a viewers' video feed
65 as input and gives the time mapped graph of the sentiments.

66 4 Methodology

67 4.1 Data

- 68 • Cohn-Kanade face database was compiled by researchers at Carnegie Mellon University, this
69 database is easily available upon request using the EULA form. The dataset is available in
70 two versions CK and CK+, Version 1, includes 486 sequences from 97 posers and is referred
71 to as CK. It consists of a range of expression from neutral to peak. The peak expression is
72 EMACS coded. However the given labels correspond to the requested expression and not
73 the actual expression. Cohn Kannade 2 Version 2 of the data set includes both spontaneous
74 as well as requested set of images. It also provides standards for facial feature tracking and
75 sentiment analysis. The CK/ CK+ dataset provided a set of 120 images, we expected to
76 extract around 100-150 more images using FACS to emotions conversion. This has been
77 discussed in great detail later. [8]
- 78 • ISED Dataset or Indian Spontaneous Expression Database provides frontal images cor-
79 responding to video while watching emotion inducing clips. The labelling is done by 4
80 decoders and validated by the stimuli and the self-report of sentiments. It covers 4 emotions

81
82
83
84
85
86

which include include happiness, disgust, sadness, and surprise. The dataset comprises of around a set of 500 peak images.[9]

- FERC 2013 dataset has a collection of 28,709 low resolution images. Each of these images is a 48×48 pixel grayscale image of the face, with the face being centered and occupying roughly the same area in each. The labels for this dataset are not posed unlike CK+, making them hard to classify, however owing to the large size of our data set this comes to benefit since the classifier becomes robust.

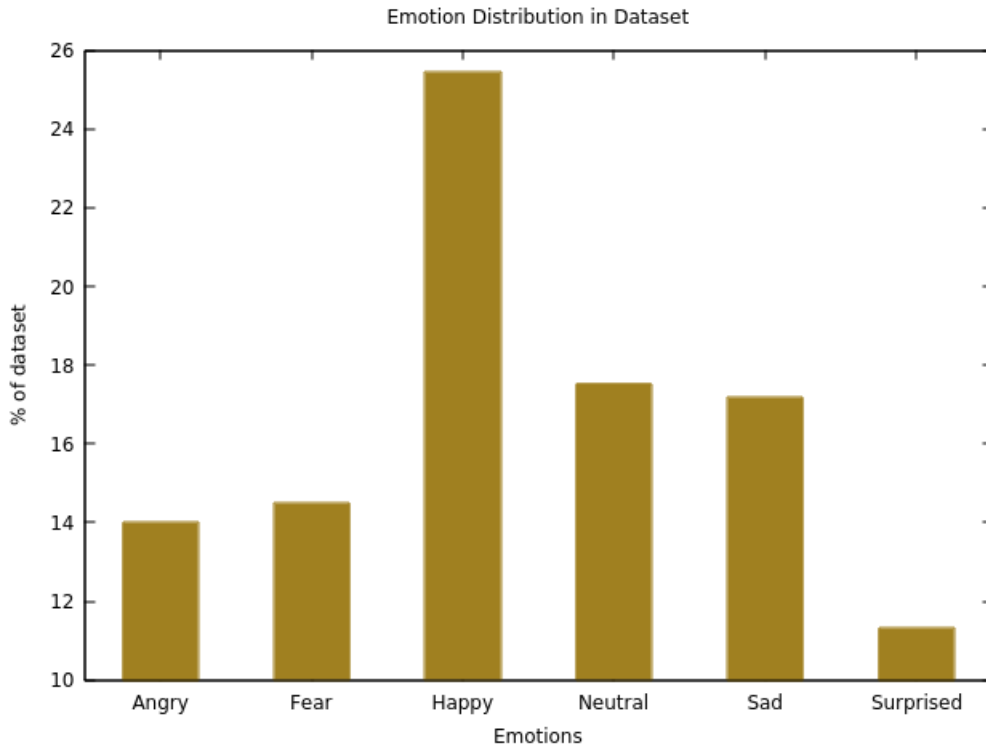


Figure 1: Emotion Distribution in Dataset

87



Figure 2: Sample Images from FERC Dataset

88 **4.2 FACS vs EMFACS**

89 FACS or Facial Action Coding System- It refers to the facial muscle movements which cause an
90 emotion. The system arranges sentiments by assessing the movements of muscles at the face, this
91 enables us to code any emotion/ expression. FACS is the standard for automated system which
92 analyses faces in videos. FACS labels each component of observed facial movement in form of
93 Action Units(AU). This is the only technique using which we can read emotions from an image/
94 video in real time. FACS eliminate the requirement of an individual to label a video or image as per
95 the sentiments.[10]

96 It was planned to use EMFACS, because it seemed to be highly appropriate for sentiment based
97 analysis, also it would save time and allow us to use a more complex architecture, however the
98 problems have been discussed below. EMFACS or emotion FACS- Under EMFACS the FACS is
99 applied selectively, only those images which are likely to have any emotional significance are coded.
100 Prototypes which are of emotional significance, the application of FACS is decide by these prototypes.
101 It is obvious that EMFACS saves on times as compared to normal FACS. However since no definite
102 standard has been defined it is difficult to obtain consistency on EMFACS coding. Moreover, the
103 set of instructions EMFACS delivers to selectively use FACS which are available only to those who
104 obtain certifications on FACS(by Paul Ekman). [11]

105 **4.3 Convolutional Neural Networks**

106 In machine learning, a convolutional neural network(CNN) is a class of deep, feed-forward neural
107 networks. CNNs use multiple hidden layers of nodes which, in an abstract sense transfer information
108 on recieving a stimulus. As images have strong spacial structures, these types of neural networks are
apt for the task at hand.

2D Convolution Operation

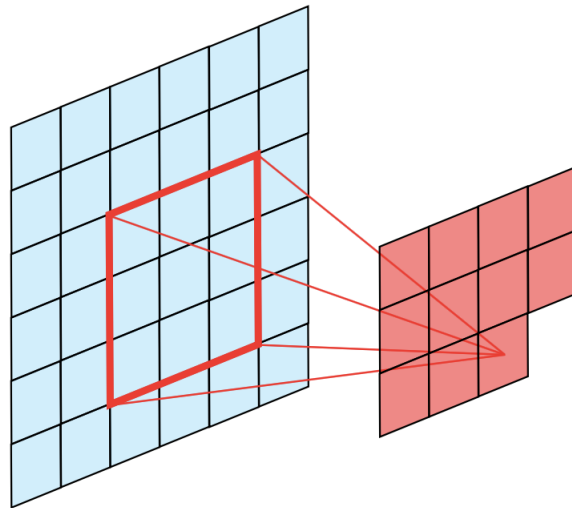


Figure 3: 2D Convolution¹

109
110

111 As the images have many pixels, it is not possible to use fully connected layers and hence it is useful
112 to use some modifications to the network architecture so that we are able to efficiently train the CNN
113 on many images from the dataset.

114 In a convolutional neural network, layers are sparsely connected with parameter sharing which
115 drastically reduces the number of parameters to be learnt and thus, reducing the training and testing
116 time. These are the type of local operations which are desired to capture objects like images (finding
117 patterns locally) and are called convolution operations. These operations determine features like

¹Image courtesy: Purushottam Kar , Course CS771 IITK

118 edges,etc using kernels.
 119 The second type of operation used is the pooling operations which aims to make the network
 120 insensitive to small/minor changes. Some nodes at fixed intervals (strides) are selected and a function
 121 is applied to them in order to create only one node representing all of them. The common type of
 122 functions used here are the max or average function.

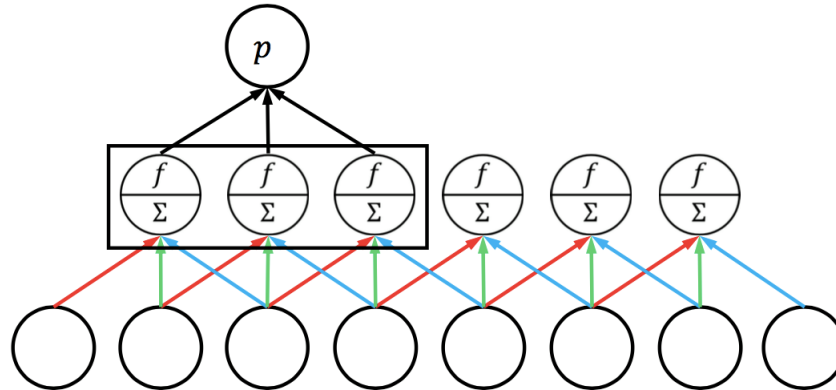


Figure 4: Pooling Operation²

123 4.3.1 Variants of CNNs tried

124 1. LeNet

125 These are very basic multilayer neural networks which constitute the backpropagation
 126 algorithm. These Gradient based classifiers can be used to extract high level complex
 127 features depending on the network architecture. As discussed in [12] Lenets were initially
 128 used for the purpose of document extraction. Lenets are seen as learning relevant features
 129 in initial layers and training on the given set of features for the coming layers. Although
 130 artificial intelligence is useful but it is not possible to avoid the bias, which manifests in
 131 the bias on the architecture of the network, therefore the architecture should be specific to
 132 the problem at hand. The output of each layer is expressible as a function of inputs from
 133 previous layers and the edge weights, using which the gradients are found and optimum
 134 values of edge weights are found, The lenet used is basically a 5 layer neural network for
 135 the task, learning parameters using gradient descent. It consists of repeated convolutional
 136 followed by pooling(max pooling), followed by a final hidden layer and an output layer. The
 137 activation function is ReLU in hidden layers and finally softmax at the output. However the
 138 accuracy obtained through normal lenet is not satisfactory, hence we tried to modify it to
 139 improve accuracy.

140 2. Modified LeNet

141 As the observed accuracy on lenet was very poor, we decided to tweak the network parameters
 142 and other hyper parameters with the aim to obtain better accuracy. Deep learning networks
 143 are very hard to train (in general NP hard) and thus we have no other option than trying
 144 new parameters and functions . The activation function used was the same , ReLU . We
 145 now used the Adam Optimiser which is can used instead of the classical stochastic gradient
 146 descent procedure to update network weights iterative based in training data. We now have 7
 147 layers in our architecture instead of 5. The training time increased as we are now using more
 148 layers than before. The perfomance increased a little than before but was still unsatisfactory.
 149 The performance for lenet and modified lenet are tabulated below.

²Image courtesy: Purushottam Kar , Course CS771 IITK

CNN Model Used	Top-1 Accuracy (%)
LeNet	29.6
Modified LeNet	36.3

Table 1: Top-1 Accuracy for Lenet and Modified Lenet

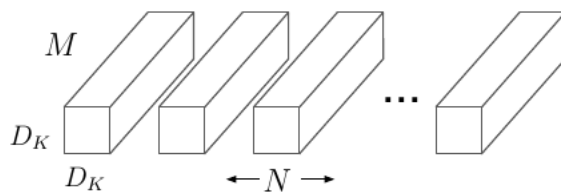
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166

3. MobileNets[13]

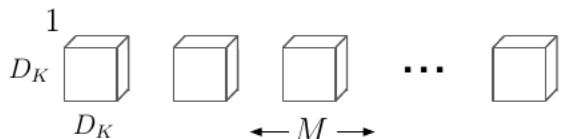
MobileNets are basically streamlined light weight network architectures for mobile and embedded applications. They are more efficient with respect to size and speed, while maintaining almost similar accuracies. They are able to achieve this by inculcating depthwise separable convolutions in their architecture, which basically breaks the interaction between the number of output channels and the size of the kernel. This convolution technique is explained below.

It is made up of 2 layers, depthwise convolutions and pointwise convolutions. The depthwise convolutions apply a single filter to each of the input channels, but does not combine them further to create new features for the output layer. This is where the 1×1 pointwise convolutional layer comes in. It computes a linear combination of the output of the depthwise convolution layer, via 1×1 convolution, for each of the channels of the output to produce the final output's feature map.

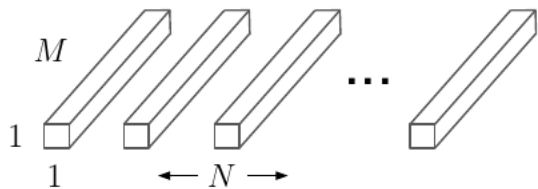
By using this, we get an overall complexity reduction of $\frac{1}{N} + \frac{1}{D_K^2}$, compared to the standard, fully connected CNN; where N is the number of output channels and $D_K \times D_K$ is the spatial dimension of the kernel.



(a) Standard Convolution Filters



(b) Depthwise Convolutional Filters



(c) 1×1 Convolutional Filters called Pointwise Convolution in the context of Depthwise Separable Convolution

Figure 5: Depthwise Separable Convolutions[13]

167 In Figure 5 above, M is the no. of channels in the input layer, N is the number of channels
 168 in the output layer and $D_K \times D_K$ is the spatial dimension of the kernel.

169 Due to low levels of accuracy using the above methods, research material was referred which
 170 could help to modify the CNN architecture in a task-specific manner, as discussed in [14],
 171 the problem involves observing high level abstractions, hence a number of fully connected
 172 nodes were also employed. Many changes were applied to the existing CNN, however the
 173 improvements in accuracy were not significant, therefore we shifted to a larger CNN model

174 4. AlexNet[15]

175 Some features of the AlexNet are:

- 176 • 7 hidden weight layers
- 177 • 650K neurons
- 178 • 60 million parameters
- 179 • 630 million connections

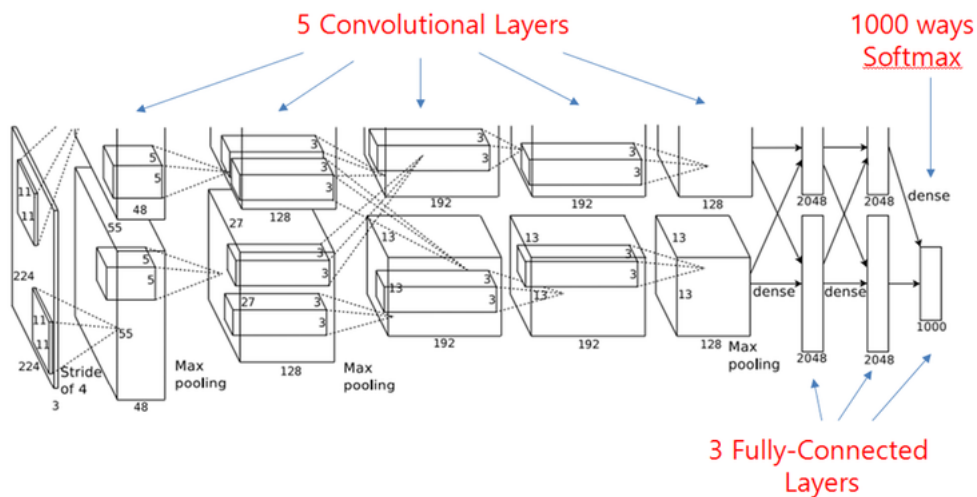


Figure 6: AlexNet Architecture³

180 AlexNet is a larger neural network with 7 hidden layers, the first 5 layers being convolutional
 181 layers and last 2 fully connected layers. But it contains around 60 million parameters, and
 182 is meant to classify among 1000 classes. Also, it requires a lot of computational resources
 183 to train an AlexNet from scratch. Since we only need to classify among the six emotion
 184 classes, we tried a smaller modified version, with only 3 convolutional layers and 2 fully
 185 connected layers along with reduced number of parameters for these layers.

186 Meanwhile, we researched as why the above mentioned neural network models were
 187 performing so poor. One fact that came out was we were training all these models for atmost
 188 2000 steps, which is around 3 epochs for a batch size of 32 images and a very less quantity.
 189 So for this model, we gradually increased the number of epochs and found that the accuracy
 190 of our model improved with increase in the number of epochs.

191 4.4 Pre-processing of Data

192 Open CV Haar Cascade- It is a classifier for object detection, used to detect faces in our case. The
 193 number of neighbors in this classifier were tuned and finally set at 5. This helps to eliminate false
 194 positive detections from the image, it means that any bounding box needs to be surrounded by atleast
 195 5 other bounding boxes for it to be classified as a face. This also ensured that images when being sent
 196 to the classifier during training and test time (consisting of 48 x 48 gray scale images) are centred at
 197 the face. [16]

³Image Source: <https://world4jason.gitbooks.io/research-log/content/deepLearning/CNN/Model%20%20ImgNet/alexnet/img/alexnet2.png>

198 **4.5 Results**

199 The accuracies for various CNNs employed are tabulated below, different values for MobileNet
 200 correspond to different values of hyper-parameters α (the width parameter) and ρ (the resolution parameter).

CNN Model Used	Top-1 Accuracy (%)
LeNet	29.6
Modified LeNet	36.3
MobileNet_0.25_128	38.2
MobileNet_0.25_160	44.4
MobileNet_0.25_192	42.4
MobileNet_0.25_224	35.0
MobileNet_0.50_128	37.6
MobileNet_0.50_160	43.9
MobileNet_0.50_192	43.2
MobileNet_0.50_224	44.1
MobileNet_0.75_128	35.4
MobileNet_0.75_160	44.2
MobileNet_0.75_192	41.3
MobileNet_0.75_224	42.5
MobileNet_1.0_128	43.6
MobileNet_1.0_160	44.0
MobileNet_1.0_192	42.9
MobileNet_1.0_224	38.9
Smaller AlexNet	74.7

Table 2: Top-1 Validation Accuracies for various CNNs employed

201 The following graph shows the detected labels as a function of time, this is what a content creator
 202 would desire, a temporal feedback of the content viewed by the user.

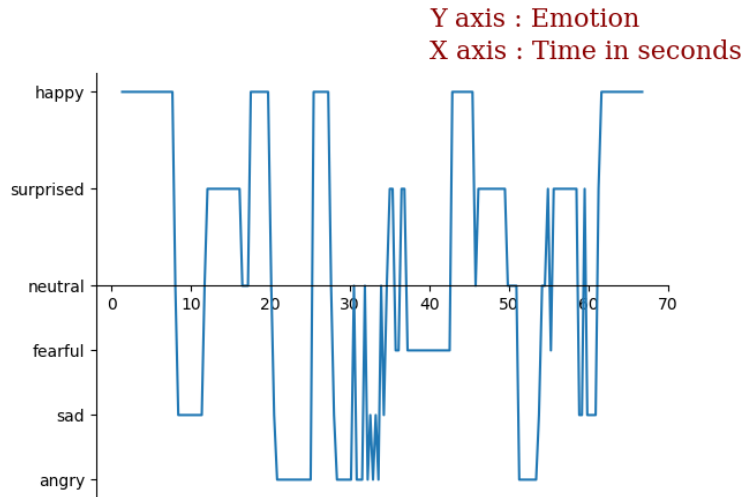


Figure 7: Peak Labelled Emotion during Testing for Final Classifier

203

204 The Recall and Precision for various labels and the overall accuracy for the final model has been
 205 tabulated below:

Emotion Label		Top-1 Accuracy (%)	Top-2 Accuracy
Angry	Recall	30.69	54.24
	Precision	37.01	65.42
Fearful	Recall	22.62	49.94
	Precision	31.78	70.18
Happy	Recall	70.52	80.72
	Precision	70.72	80.94
Sad	Recall	39.33	70.47
	Precision	28.38	50.85
Surprised	Recall	32.34	74.33
	Precision	61.73	73.60
Neutral	Recall	49.70	71.89
	Precision	44.84	64.86
Overall		53.29	71.16

Table 3: Final Model Results over Test Data

205

206 5 System limitations and Privacy Concern

207 With privacy being a major concern with our proposed system which claims to monitor viewer's
 208 activity through a webcam feed, we propose the following solutions to the posed problem.

- 209 • Ask permission of the user to allow sharing of his feed right before they start viewing your
 210 sentiment data
- 211 • Do not include this as a feature of the final application but use this system on a small set of
 212 volunteers to identify the strong and weak points

213 6 Possibilities for Future Work

214 The data generated by using this CNN is rich enough to be exploited in several applications to derive
 215 better insight about both, the users as well as the content.

- 216 • **Improvement of existing online recommendation systems**
 217 The users' emotional state while browsing the web can be monitored and this can be used to
 218 decide what content/advertisements should be recommended next to the user. This will help
 219 the recommendation system reach its predictions, using both, the past online behaviour of
 220 the user with the website/application, as well as his present emotional state.
- 221 • **Depression detection applications**
 222 Highly frequent, content independent instances of sad/disgusted/angry sentiments can be
 223 used by applications to detect possibilities of depression in a user well before the problem
 224 becomes chronic.
- 225 • **Exploiting the information along temporal dimension using RNNs**
 226 Since the data from a video, is in form of a stream or sequence of images, it would be much
 227 better if the algorithm could give feedback to the future evaluations, this is inherently the
 228 idea behind the working of an RNN. We can implement this to make the output robust, as it
 229 would take in feedback and thus it will be less prone to noise.
- 230 • **Improving Robot Interaction Experience**
 231 The responses of robots while interacting with humans can be made more emotionally
 232 informed by using this technique with a camera mounted on the head of the robot.

233 **References**

- 234 [1] “A discussion on youtube forum, are likes and dislikes really important.” [https://](https://productforums.google.com/forum/#!topic/youtube/3jbHpNIZpT8)
235 productforums.google.com/forum/#!topic/youtube/3jbHpNIZpT8.
- 236 [2] M. Z. Asghar, S. Ahmad, A. Marwat, and F. M. Kundi, “Sentiment Analysis on Youtube: A
237 Brief Survey,” *CoRR*, vol. abs/1511.09142, 2015.
- 238 [3] Z. Hussain, T. Patanam, and H. Cate, “Group Visual Sentiment Analysis,” *ArXiv e-prints*, Jan.
239 2017.
- 240 [4] V. Gajarla and A. Gupta, “Emotion detection and sentiment analysis of images,” in *Georgia*
241 *Institute of Technology*, 2015.
- 242 [5] A. Bermingham, M. Conway, L. McNerney, N. O’Hare, and A. F. Smeaton, “Combining social
243 network analysis and sentiment analysis to explore the potential for online radicalisation,” in
244 *2009 International Conference on Advances in Social Network Analysis and Mining*, pp. 231–
245 236, July 2009.
- 246 [6] A. Go, L. Huang, and R. Bhayani, “Twitter sentiment analysis,” *Entropy*, vol. 17, p. 252, 2009.
- 247 [7] H. Wang, D. Can, A. Kazemzadeh, F. Bar, and S. Narayanan, “A system for real-time twitter
248 sentiment analysis of 2012 us presidential election cycle,” in *Proceedings of the ACL 2012*
249 *System Demonstrations*, pp. 115–120, Association for Computational Linguistics, 2012.
- 250 [8] “Details of cohn kahnade dataset.” <http://www.pitt.edu/~emotion/ck-spread.htm>.
- 251 [9] “Ised database.” <https://sites.google.com/site/iseddatabase/>.
- 252 [10] “Facs.” https://en.wikipedia.org/wiki/Facial_Action_Coding_System;[https://](https://imotions.com/blog/facial-action-coding-system/)
253 imotions.com/blog/facial-action-coding-system/;[https://www.paulekman.](https://www.paulekman.com/product-category/facs/)
254 [com/product-category/facs/](https://www.paulekman.com/product-category/facs/).
- 255 [11] “Emfacs.” <https://www.paulekman.com/facs-faq/>.
- 256 [12] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document
257 recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- 258 [13] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and
259 H. Adam, “Mobilenets: Efficient convolutional neural networks for mobile vision applications,”
260 *CoRR*, vol. abs/1704.04861, 2017.
- 261 [14] Q. You, J. Luo, H. Jin, and J. Yang, “Robust image sentiment analysis using progressively
262 trained and domain transferred deep networks.,” in *AAAI*, pp. 381–388, 2015.
- 263 [15] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional
264 neural networks,” in *Advances in Neural Information Processing Systems 25* (F. Pereira, C. J. C.
265 Burges, L. Bottou, and K. Q. Weinberger, eds.), pp. 1097–1105, Curran Associates, Inc., 2012.
- 266 [16] “Haar cascade.” https://docs.opencv.org/3.3.0/d7/d8b/tutorial_py_face_
267 [detection.html](https://docs.opencv.org/3.3.0/d7/d8b/tutorial_py_face_detection.html).