

Machine Learning Engineer Nanodegree

Capstone Project

Keshri Nandan
August 27th, 2017

Project Overview

Anyone who has watched the American TV show “Lie to me”, would share my fascination of reading people’s emotional state from their faces alone. While most people are naturally able at this, machines haven’t yet advanced enough to match a human at this skill. As AI makes great advances in performing repeatable tasks with perfection, strategizing in complex games etc., emotional intelligence still remains an unexplored domain.

As we build smarter systems, contextual awareness becomes more and more important. The emotional state of the user becomes an important factor in defining the context. Reading emotions directly of the faces can be huge leap and can be beneficial in multiple applications:

1. Intelligent homes: house lighting, music etc. can be customized to the mood
2. Personalized marketing: The ad campaigns can be better targeted with the right emotional state in mind
3. Surveillance: surveillance systems could become smarter in predicting events as in “Person of Interest”

Emotions (the seven basic ones) surface on faces in universally characteristic manner. These characteristic features can be leveraged to recognize emotions on any face.

Anger



Disgust



Fear



Happy



Sad



Surprise



Neutral



This topic has garnered great amount of interest in the recent past. It became all too popular when contests were hosted on Kaggle and EmotiW2015 to devise efficient algorithms to identify facial emotions. Multiple approaches like CNN, SVM, SIFT & MKL have shown great results. CNNs have proved to be the most promising in terms of improving the accuracy further.

Problem Statement

The most common emotions have similar facial patterns. Given the image of a face, the objective is to identify the emotion expressed by the face. The task is to categorize the faces based on the emotion shown in the facial expression in to one of seven categories (0=Angry, 1=Disgust, 2=Fear, 3=Happy, 4=Sad, 5=Surprise, 6=Neutral)

Solution Approach

With characteristic features evident from faces across the globe, it suits well to extract them to identify the emotions.

The pixels in close proximity are more related with each other than with pixels that are a greater distance away. This **spatial structure** of the image needs to be leveraged.

While these patterns are universal, they surface on the faces at different locations, thus **shift invariance** is required of the solution we adopt. From the discussion it is evident that we are solving a **classification** problem here.

All of these conditions lead us to use **Convolutional neural networks**. We shall be starting with a base structure and build on it to improve its performance by tuning its hyper-parameters.



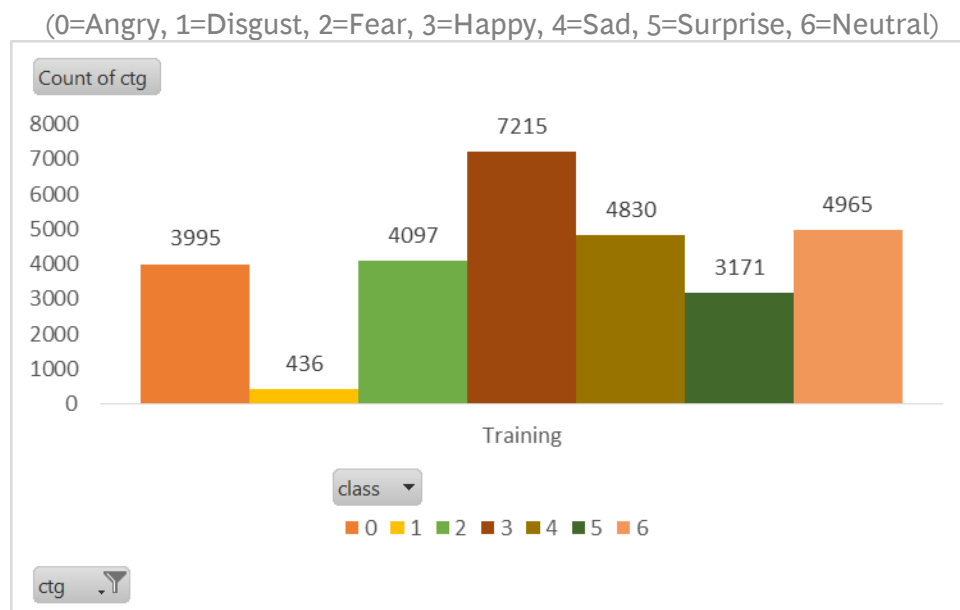
Analysis

Datasets and Inputs

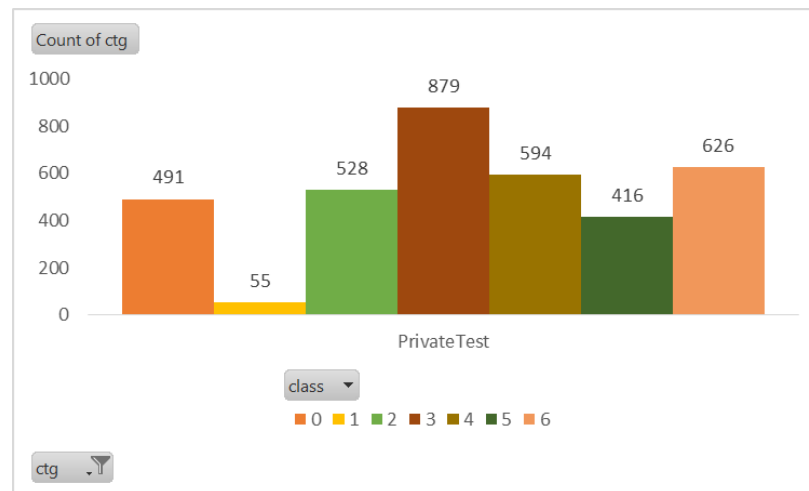
The data to be used is the **Fer2013** dataset shared in the Kaggle contest: “Challenges in Representation Learning: Facial Expression Recognition Challenge”. It has around 35000 facial images in 48x48 greyscale pixel format and the associated emotion as label. There are 7 standard emotions identified in the dataset. The dataset contains two columns, “emotion” and “pixels”. The “emotion” column contains a numeric code ranging from 0 to 6, inclusive, for the emotion that is present in the image. The “pixels” column contains a string surrounded in quotes for each image. The contents of this string are space-separated pixel values in row major order.

The images are divided under 3 categories:

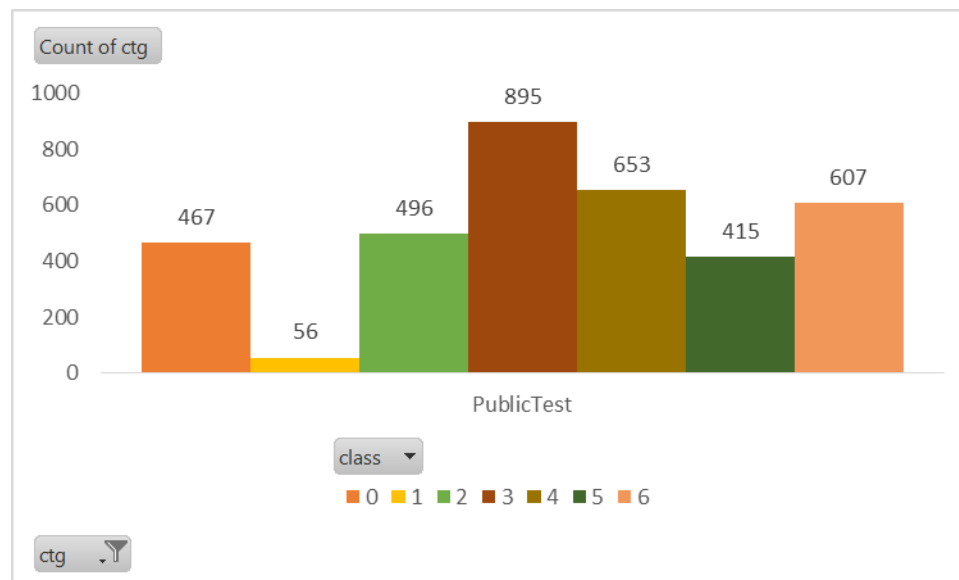
1. Training: It has 28,709 images spread across 7 emotions as below:



2. Validation: It has 3589 images distributed as:



3. Test: It has 3589 images distributed as:



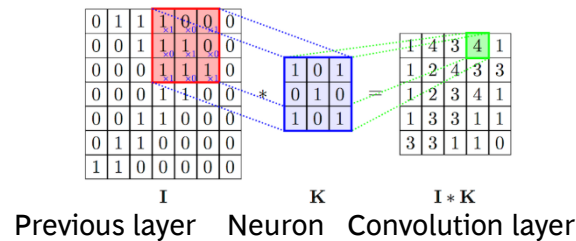
Data pre-processing

The dataset is hosted on the Kaggle competition's ("Challenges in Representation Learning: Facial Expression Recognition Challenge") page. It has three columns:

1. "Emotion" – indicating emotion of the face
2. "pixels" – it is a 2304X1 array of pixels of the 48X48 greyscale image
3. "category" – it segregates the dataset into:
 - a. Training
 - b. PrivateTest: used for validation
 - c. PublicTest: used as test set

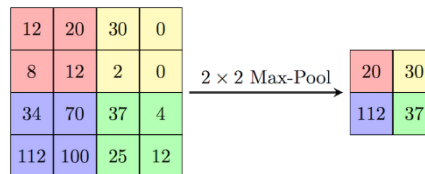
The 48X48 grid for a face looks like this:

70	80	82	72	58	58	60	63	54	58	60	48	89	115	121	119	115	110	98	91	84	84	90	99	110	126	143	153	158	171	169	172	169	165	129	110	113	107	95	79	66	62	56	57	61	52	43	41
65	61	58	57	56	69	75	70	65	56	54	105	146	154	151	151	155	155	150	147	147	148	152	158	164	172	177	182	186	189	188	190	188	180	167	116	95	103	97	77	72	62	55	58	54	56	52	44
50	43	54	64	63	71	68	64	52	66	119	156	161	164	163	164	167	168	170	174	175	176	178	179	183	187	190	195	197	198	197	198	195	191	190	145	86	100	90	65	57	60	54	51	41	49	56	47
38	44	63	55	46	52	54	55	83	138	157	158	165	168	172	171	173	176	179	179	180	182	185	187	189	189	192	197	200	199	196	198	200	198	197	197	91	87	96	58	58	59	51	42	37	41	47	45
37	35	36	30	41	47	59	94	141	159	161	164	170	171	172	176	178	179	182	183	183	187	189	192	192	194	195	200	200	199	199	200	201	197	193	111	71	108	69	55	61	51	42	43	56	54	44	
24	29	31	45	61	72	100	136	150	159	163	162	163	170	172	171	174	177	177	180	187	186	187	189	192	192	194	195	196	197	199	200	201	200	197	201	137	58	98	92	57	62	53	47	41	40	51	43
24	35	52	63	75	104	129	143	149	158	162	164	166	171	173	172	174	178	179	187	188	188	191	193	194	195	198	199	199	197	198	197	197	197	201	164	52	78	87	69	58	56	50	54	39	44	42	
26	31	49	65	91	119	134	145	147	152	159	163	167	171	170	169	174	178	179	187	187	185	187	190	188	187	191	197	201	199	199	200	197	196	197	182	58	62	77	61	60	55	49	59	52	54	44	
22	30	47	68	102	123	136	144	148	150	153	157	167	172	173	170	171	177	179	178	186	190	186	189	196	193	191	194	190	190	192	197	201	201	199	194	189	69	48	74	56	60	57	50	59	51	41	
20	34	47	79	111	132	139	143	145	147	150	151	160	169	172	171	167	171	177	174	180	182	181	192	186	190	192	198	195	194	196	198	201	202	195	189	70	48	50	59	61	57	54	54	61	52	36	
26	40	61	93	124	135	138	142	144	146	151	152	158	165	168	168	165	161	164	173	172	167	172	167	180	198	198	193	199	195	194	198	200	198	197	195	190	65	35	68	59	62	57	60	59	50	44	
32	54	90	115	132	137	138	140	144	146	146	156	165	168	174	176	176	175	168	168	169	171	175	171	172	192	194	184	198	205	201	194	195	193	195	192	186	57	38	72	65	57	62	58	57	60	54	49
47	79	116	130	138	141	141	139	141	143	145	157	164	164	166	173	174	176	179	179	176	181	189	188	173	180	175	160	182	189	198	192	189	190	190	188	172	46	44	64	66	59	62	57	56	62	53	50
66	103	133	137	141	143	141	136	132	131	136	127	118	111	107	108	123	131	143	154	158	166	177	181	175	170	159	148	171	161	176	185	192	194	188	190	162	53	49	58	63	61	61	55	56	61	51	50
81	116	139	142	142	146	144	136	128	119	112	97	85	90	91	88	92	90	80	81	84	106	122	132	144	145	144	147	163	147	163	173	181	190	187	191	167	61	48	53	61	61	58	54	56	61	51	53
89	123	140	144	145	146	147	136	122	107	99	95	92	90	87	83	76	67	52	46	52	63	69	83	96	119	132	148	159	136	137	143	138	143	152	156	156	70	48	50	59	61	57	54	54	61	52	36
93	124	135	140	144	148	150	140	125	114	101	80	54	56	54	41	41	33	40	39	35	49	60	63	74	107	129	147	147	116	111	100	77	76	86	108	111	73	49	50	60	62	60	57	55	63	59	56
89	121	134	139	146	151	152	150	141	127	111	96	77	85	70	32	31	37	91	65	50	48	59	73	83	112	136	155	130	60	46	38	40	43	81	116	91	72	52	48	58	62	62	59	53	61	59	52
85	114	134	140	147	154	159	158	153	145	143	150	126	121	125	68	45	89	137	95	70	78	75	95	109	131	153	171	94	23	16	32	82	82	65	113	77	71	54	48	56	62	62	60	53	60	56	52
75	108	133	140	148	156	169	167	163	156	155	146	112	119	134	127	142	140	121	117	129	114	120	129	146	174	191	98	46	33	33	109	147	98	109	67	73	55	50	56	64	64	61	58	61	53	54	
64	106	129	140	148	159	169	175	176	174	165	159	156	145	120	115	124	127	131	133	141	147	142	141	147	161	182	202	154	114	96	100	158	158	153	123	61	76	57	48	56	64	64	63	62	61	54	55
44	97	131	137	147	158	168	177	181	183	179	170	168	169	165	155	152	151	152	154	162	165	158	151	158	168	187	206	186	147	135	144	145	152	178	115	57	74	58	48	58	64	63	59	63	55	53	
66	104	130	132	144	153	162	170	180	185	187	181	178	182	180	177	173	171	177	176	172	164	161	167	164	185	207	197	173	152	141	141	161	193	104	54	69	60	48	57	65	62	60	57	64	55	50	
94	111	124	130	135	150	159	163	172	179	184	184	178	178	177	173	171	174	177	178	176	169	165	161	163	161	180	205	201	183	171	177	178	180	194	101	55	65	60	47	55	65	63	59	58	63	57	52
90	105	117	122	130	133	153	157	163	171	174	182	183	182	178	174	175	175	177	175	172	163	161	159	157	162	178	200	201	188	181	172	177	187	198	98	57	63	61	48	52	61	64	63	60	65	57	51
95	104	113	117	127	136	145	152	156	162	165	173	177	182	183	183	180	181	177	165	153	154	155	160	174	193	200	200	188	185	180	182	192	196	101	60	60	56	49	50	60	66	64	62	64	59	53	
99	104	111	112	118	132	142	147	155	158	160	159	162	171	176	184	186	183	180	169	154	141	135	145	155	164	180	196	205	188	189	188	189	193	192	98	61	64	55	49	49	60	66	63	64	63	60	57
99	105	108	112	113	125	139	143	150	155	158	164	169	174	176	182	183	182	177	163	141	133	147	151	164	170	185	200	210	194	188	192	186	185	180	88	64	67	60	46	50	59	65	64	64	64	59	56
101	103	108	109	109	118	134	143	143	147	155	159	166	171	174	177	179	178	172	153	129	143	161	159	166	173	186	197	207	203	185	191	183	179	164	73	67	67	66	48	50	57	65	65	63	64	61	57
103	108	114	112	110	115	128	138	144	145	152	156	159	164	168	172	172	169	161	139	125	147	156	161	162	164	180	188	188	197	185	187	181	180	137	65	70	68	70	52	47	53	62	63	63	65	61	58
105	109	112	120	113	112	122	134	141	149	150	153	155	159	164	167	167	162	152	134	115	126	119	106	99	109	141	158	150	155	175	184	176	175	106	63	70	68	68	50	46	50	57	63	63	64	61	59
107	110	117	117	114	117	128	137	147	148	150	153	156	161	162	163	156	150	148	105	70	45	26	25	47	73	74	79	128	177	180	173	157	77	66	68	67	68	52	49	51	56	62	62	62	62	60	
101	107	108	114	115	114	117	125	134	143	148	149	152	154	158	160	158	155	160	158	132	88	73	73	64	52	66	91	138	160	174	173	171	125	64	67	63	64	68	54	50	49	54	60	60	60	62	60
98	105	108	112	113	125	139	143	150	155	158	164	169	174	176	182	183	182	177	163	141	133	147	1																								



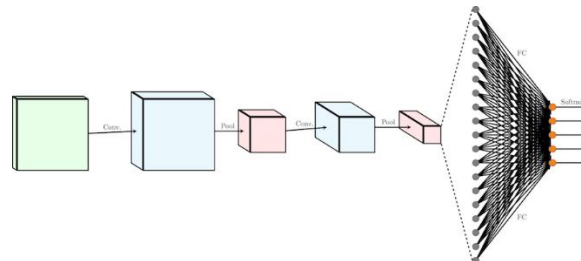
2. Pooling layers:

After each convolutional layer, there may be a pooling layer. The pooling layer takes small rectangular blocks from the convolutional layer and subsamples it to produce a single output from that block. There are several ways to do this pooling, such as taking the average or the maximum, etc. Our pooling layers will always be max-pooling layers; that is, they take the maximum of the block they are pooling. A diagrammatical illustration of $2 \times 2 \times 2$ max-pooling is given below



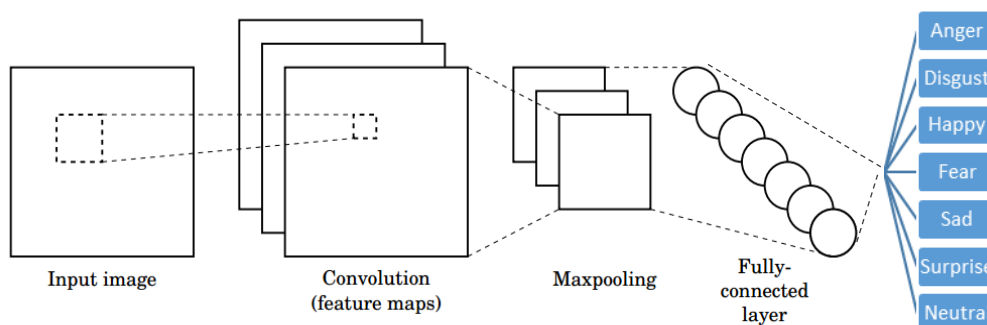
3. Fully connected layers:

After several convolutional and max pooling layers, the high-level reasoning in the neural network is done via fully connected layers. A fully connected layer takes all neurons in the previous layer (be it fully connected, pooling, or convolutional) and connects it to every single neuron it has. Fully connected layers are not spatially located anymore.



Input – Conv – MP – Conv – MP – Fully connected layer

The architecture of a CNN is designed to take advantage of the 2D structure of an input image. This is achieved with local connections and tied weights followed by some form of pooling which results in translation invariant features. Another benefit of CNNs is that they are easier to train and have many fewer parameters than fully connected networks with the same number of hidden units.



Evaluation Metrics

1. Accuracy: The accuracy of predicting which of the seven classes should be assigned to each image

The classification accuracy A_i of an individual program i depends on the number of samples correctly classified (*true positives plus true negatives*) and is evaluated by the formula:

$$A_i = \left(\frac{t}{n}\right) * 100$$

Where t = number of sample cases correctly classified
 n = total number of sample cases

Accuracy is the most readily comprehensible metric even to a layman. The dataset has fair representation for all classes except the (1 = Disgust) class. Hence we can use this as the performance metric.

2. Optimizer – ‘rmsprop’: The learning rates for the weights are configured to be tuned by dividing the learning rate for a weight by a running average of the magnitudes of recent gradients for that weight

While training the CNN, with increasing epochs the learning rate needs to be scaled with the gradients as too small steps shall slow down the training progress.

3. Loss function – ‘categorical crossentropy’: For classification problems the standard way of calculating losses (to be minimized) is by evaluating the cross-entropy between the predicted classes and true classes

Mathematically, it is defined as:

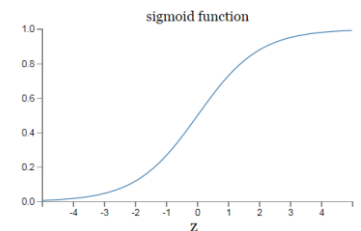
$$H(p, q) = - \sum_x p(x) \log q(x).$$

Where, $p(x)$ are the true classes and $q(x)$ probabilities of predicted classes

For neural nets, the learning rates are dependent on $(\partial C / \partial w)$ and $(\partial C / \partial b)$ where:

C = cost function, w = weights and b = biases.

If the cost function is the cross-entropy function, these terms are directly proportional to the error $(\sigma(z) - y)$ i.e. (predicted – actual). This ensures that learning at initial stages is not exceedingly slow as in case of the quadratic cost function where these terms are proportional to $\sigma'(z)$ (sigmoid)



Benchmark Models

1. Yu & Zhang reported accuracy of 61% in the EmotiW2015 challenge. They used CNN with **randomized perturbation**
2. Mollahosseini and Mahoor used **inception layers with convolution layers** and reported varying accuracy (60% - 94%) across multiple datasets (FER2013 – 66%)
3. Kim, Roh, Dong & Lee trained multiple deep **convolutional neural networks (deep CNNs) as committee members and combined their decisions** for solving a seven-class problem of static FER in the wild for the EmotiW2015, and achieved a test accuracy of 61.6 %.

Implementation

The code structure is explained below:

The python notebook “emoticon” is structured in the steps we take in the process of developing the model.

1. Step 0: Prepare image repository
This section of the code reads in the dataset (fer2013.csv) row-wise and generates the images and saves them in the associated folder structure.
This needs to be run only once while setting up the project.
2. Step 1: Import image data
The image paths are then read into model according to their category (Train, test, validate)
Their associated emotion classes are read as categorical variables across the three categories
The next cell defines the functions that read in the image from the image path and return their corresponding 4D tensors
3. Step 2: Rescale image data
The pixel data is rescaled from 0-255 to 0-1 by dividing it with 255
4. Step 3: Define the model architecture
The initial model and further tuning to it are employed here. After defining the model it is compiled and trained in the subsequent cells. While training, the model is configured to save the learning history and the best set of weights.
After we have trained the model, the best set of models are loaded.
For every tuning iteration, steps 3 and 4 are repeated
5. Step 4: Test the model
The model is now tested with the “test” dataset and its accuracy is recorded.
6. Step 5: Model prediction
The final step involves generating sample predictions from the model and comparing with actual classes

Please note:

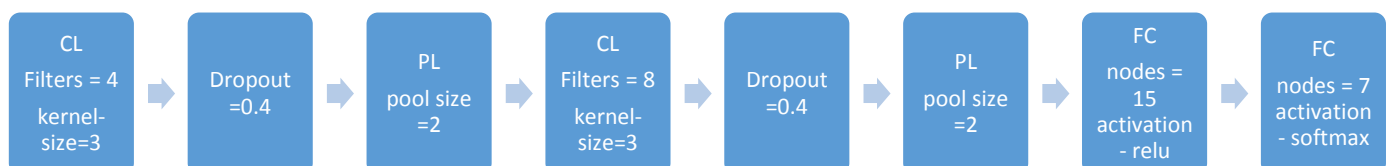
- While training the model in step 3, it is advisable to have the “verbose” option set to 1 if one is running it for 10-15 epochs as it helps keeping track of the training and validation errors.
Beyond 20 epochs, one must set it to 0 as it slows down the system (it did for me). To keep track of errors and accuracies, one should record these in a log using the “**CSVlogger**” callback.
- The model also uses the “**ModelCheckpoint**” callback. This option saves the weights of the model for every epoch. It has been configured to save the best set of weights i.e. only when the validation loss improves.

Keras offers multiple callbacks one could leverage to diagnose one’s model.

Refinement

1. **Model 1:** We begin with a simple CNN with 2 convolution layers (CL), 2 pooling layers (PL), 2 fully connected layers (FCL) : Epochs = 50

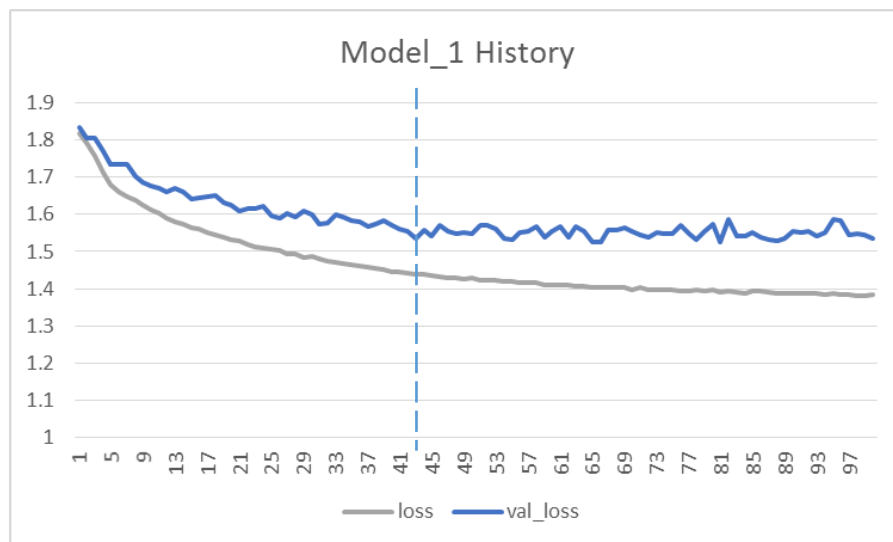
Weights (for CL) initialized with std-dev. = 0.01



Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 48, 48, 4)	112
dropout_1 (Dropout)	(None, 48, 48, 4)	0
max_pooling2d_1 (MaxPooling2)	(None, 24, 24, 4)	0
conv2d_2 (Conv2D)	(None, 24, 24, 8)	296
dropout_2 (Dropout)	(None, 24, 24, 8)	0
max_pooling2d_2 (MaxPooling2)	(None, 12, 12, 8)	0
flatten_1 (Flatten)	(None, 1152)	0
dense_1 (Dense)	(None, 15)	17295
dense_2 (Dense)	(None, 7)	112

If we follow the loss (training vs validation) history over training epochs, we observe that the validation error flattens out at around epoch = 55.

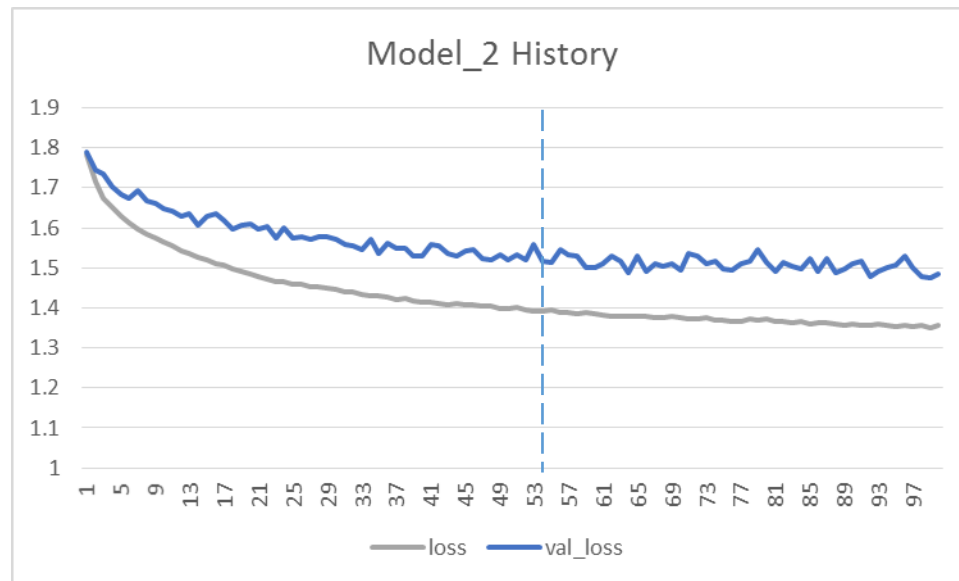
The best set of weights are saved in the folder 'saved_models'. We load these weights to check the accuracy of the model.



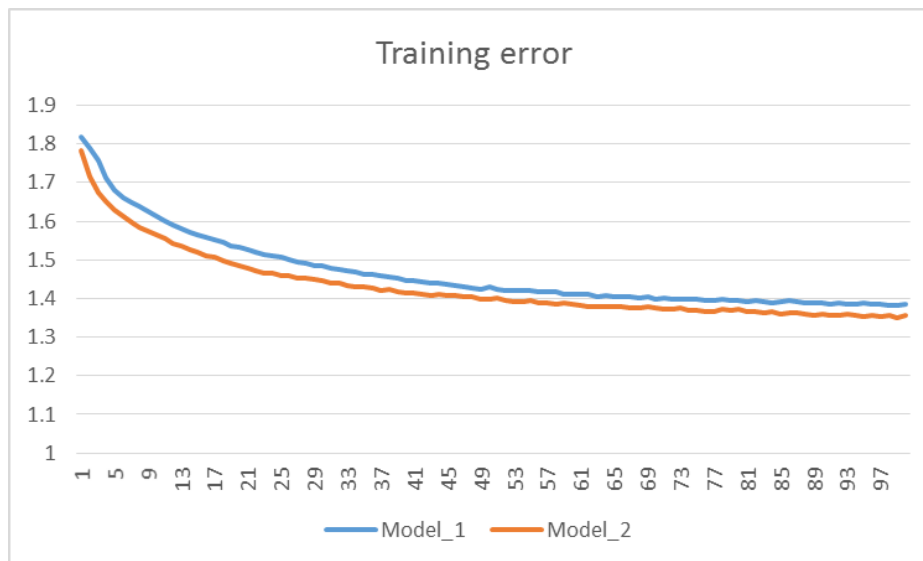
Test accuracy - 38.9245%

This lends us a good starting point to begin with. We shall tweak other parameters to further improve the accuracy of the model

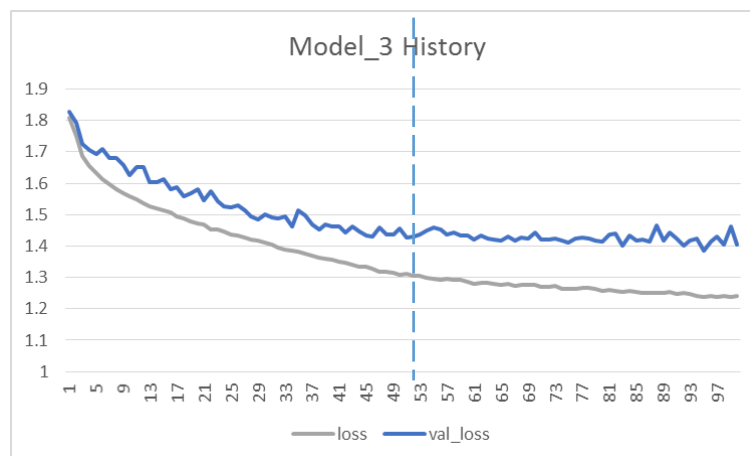
2. **Model 2:** In the next iteration we use the same architecture but initialize the CL weights by Xavier initialization (glorot_normal)
In this case the validation loss flattens out around epoch = 65.
The test accuracy jumps to **43.1039%**



If we compare the training error for model_1 & model_2, we see improvement that reflects in the test accuracy.

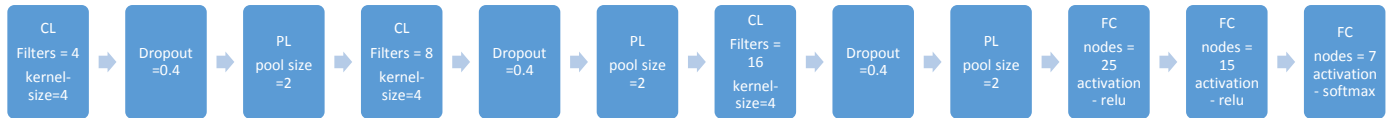


3. **Model 3:** Adding another fully connected layer (FC-25) before (FC-15), the model improves and test accuracy rises to **46.8376%**

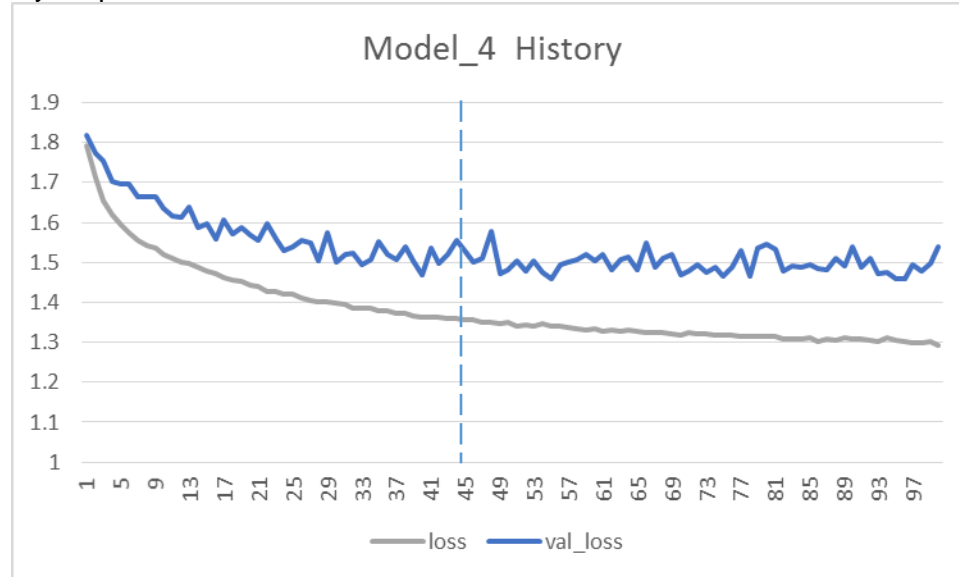


4. **Model 4:** We continue with model 3 with modifications

- Add a convolution layer + pooling layer
- Kernel sizes: 4 – 4 – 4

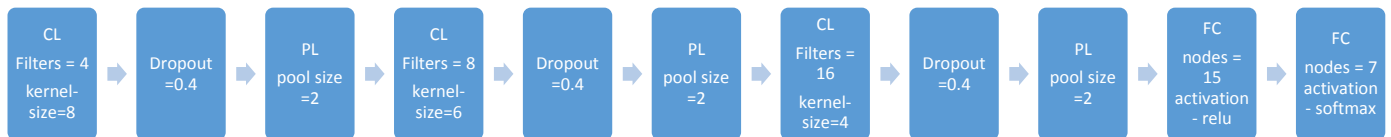


The test accuracy drops to **44.1627%**

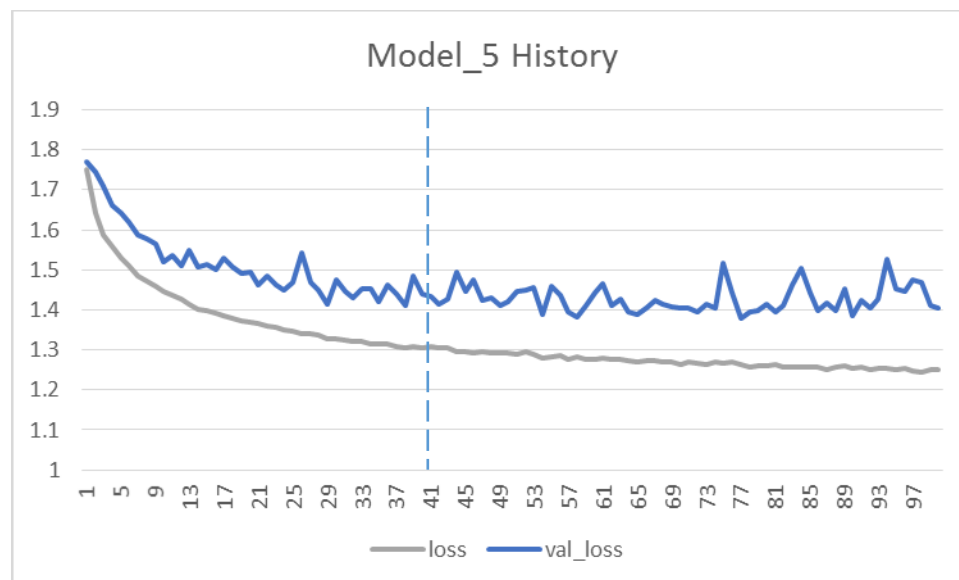


5. **Model 5:** We continue with model 4 with modifications

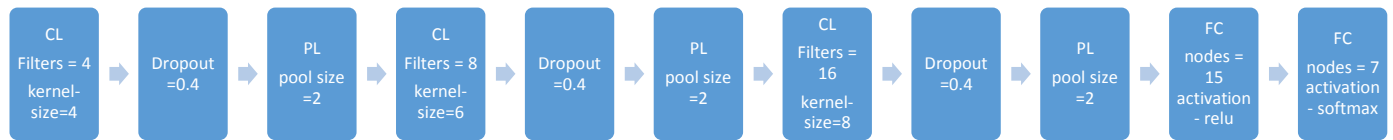
- Modify kernel-size progressively (8 – 6 – 4)



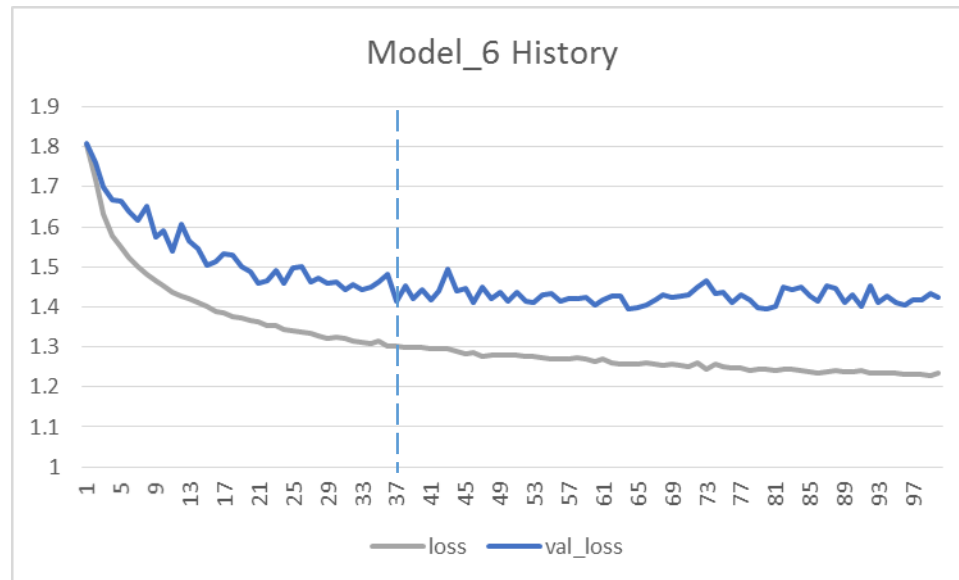
The test accuracy is reported to be **46.9211%**



6. Model 6: We continue with model 4 but reverse the kernel sizing scheme (4 – 6 – 8)



This pushes the test accuracy to **47.5899%**



Results

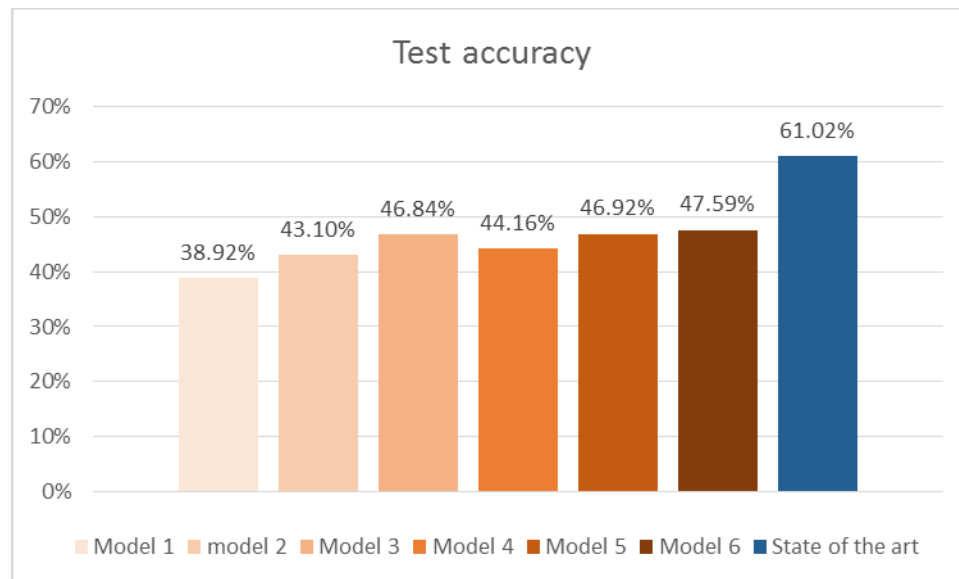
From the training history it also comes out that, as the model performance improves, the validation loss plateaus at an earlier epoch:

Model	Epoch	Accuracy
Model – 4	45	44.1627%
Model – 5	41	46.9211%
Model – 6	37	47.5899%

Model – 6 has the maximum accuracy achieved in this exercise which falls short of the state of the art accuracy (**61%**) by **13 pts**.

The best models are fairly deep as compared to the Model_6 and require greater computational resources than a laptop.

The objective of this exercise was to test out novel architecture schemes and gain significant performance gains with these successive tweaks. Further experimentation may result in better schemes that can improve existing models.



Model-6 description

The model employs varying size of kernels (4 – 6 – 8). This might help us extract multiple levels of abstraction from the images to better capture features specific to emotions. Further tuning of the kernel-sizes and their counts can help us achieve the optimum configuration.

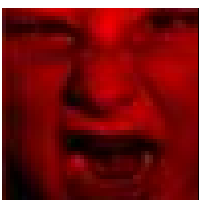
For this model the best weights are saved for epoch = 54, where validation accuracy = 45.0822%, and train accuracy = 48.9359%. Since the accuracies are within the same range (45% to 49%), and training accuracy slightly greater than validation accuracy, we can rule out overfitting.

The current value of dropout (0.4) is doing a good job of regularizing the model as we see multiple troughs and ebbs in the validation error curve unlike the monotonous decline of training error. We can explore model's performance with multiple values of dropouts.

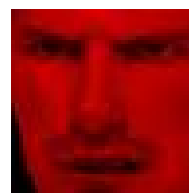
The learning curves for each iteration have been depicted with the epochs highlighted where validation errors flatten out. We have configured the model to save the best weights, this ignores the set of weights that may result due to overfitting.

Conclusion

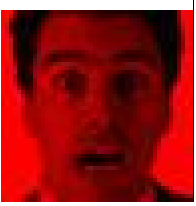
Free-form visualization



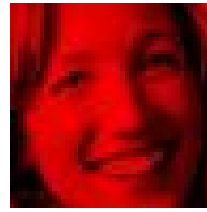
Emotion detected: Happy
Actual emotion: Angry



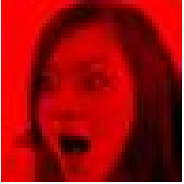
Emotion detected: Neutral
Actual emotion: Disgust



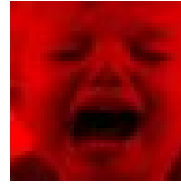
Emotion detected: Fear
Actual emotion: Fear



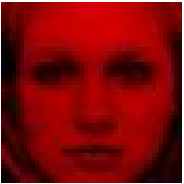
Emotion detected: Happy
Actual emotion: Happy



Emotion detected: Surprise
Actual emotion: Surprise



Emotion detected: Angry
Actual emotion: Sad



Emotion detected: Neutral
Actual emotion: Neutral

From above results we can summarize that,

- Characterizing emotions on children's faces can be tricky. ☺
- "Disgust" if expressed subtly can hard to distinguish from "Neutral" faces
- The model fares fairly well for adult faces irrespective of gender

Reflections

The steps taken in the project can be summarized in the following points:

1. An interesting problem with readily available dataset was found.
2. The dataset was used by multiple teams in past, so the benchmark was well established.
3. The methodology was settled on to with evident reasoning
4. A basic model was further refined by applying best practices and tuning the parameters
5. The best model achieved was then chosen after it consistently resulted in high accuracies

The difficult part of the process was to manually tune the architecture of the model and run multiple iterations of these versions. While parameter optimization techniques exist, they proved to be computationally expensive and don't warrant novel schemes.

The interesting part of the project was that something as trivial/natural for humans was not evident for machines and that it is part of active research across the AI community. Seeing something out of sci-fi becoming real was fascinating.

The exercise has been a great learning experience and has sparked my curiosity to employ ML techniques in various aspects of my work.

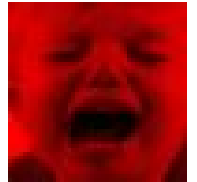
Scope of improvement

There are multiple avenues to explore to improve the accuracy of the model.

1. Augmentation of the training dataset
 - a. The human face being symmetric across the vertical axis, adding laterally inverted images in the training set might improve accuracy
2. Hyper-parameter optimization
 - a. Currently the manual approach has been followed as I have limited computational resources (no GPUs)
 - b. One could employ grid-search cross-validation, Gaussian process etc.
 - c. These are resource intensive processes and, hence manual tuning was employed
 - d. It is recommended to use SKlearn with Keras to deploy **GridSearchCV**
3. Explore deploying the deep forest model
 - a. The model has shown promising results in deep learning problems
 - b. This could not be attempted owing to paucity of time
4. Fractional max-pooling and batch-normalization
 - a. These concepts have been tested by Raghuvanshi and Choksi [10] and achieved similar accuracies
 - b. These can be tested with the current architecture to see if performance improves

While we could improve our model in various ways, we need to reassess our initial assumption that all emotions surface in the same fashion universally.

The definition of emotion appears to be malleable with age-groups. For example, the image of the kid classified as “Sad” could easily be construed as “Angry” for an adult.



Emotion detected: Angry
Actual emotion: Sad

References

- 1 <https://www.humintell.com/2010/06/the-seven-basic-emotions-do-you-know-them/>
- 2 <http://ufldl.stanford.edu/tutorial/supervised/ConvolutionalNeuralNetwork/>
- 3 Image based Static Facial Expression Recognition with Multiple Deep Network Learning
<http://www.contrib.andrew.cmu.edu/~yzhiding/publications/ICMI15.pdf>
- 4 Going Deeper in Facial Expression Recognition using Deep Neural Networks
<https://arxiv.org/pdf/1511.04110.pdf>
- 5 Hierarchical committee of deep convolutional neural networks for robust facial expression recognition
<https://link.springer.com/article/10.1007/s12193-015-0209-0>
- 6 Challenges in Representation Learning: Facial Expression Recognition Challenge
<https://www.kaggle.com/c/challenges-in-representation-learning-facial-expression-recognition-challenge>
- 7 Deep Forest: Towards An Alternative to Deep Neural Networks
<https://arxiv.org/abs/1702.08835>
- 8 <http://andrew.gibiansky.com/blog/machine-learning/convolutional-neural-networks/>
- 9 <https://cambridgespark.com/content/tutorials/convolutional-neural-networks-with-keras/index.html>
- 10 http://cs231n.stanford.edu/reports/2016/pdfs/023_Report.pdf