

Emotion Detection from Images using Machine Learning

Abstract

Convolutional Neural Networks (CNNs) have revolutionized image classification by effectively learning and extracting features from images. In this experiment, proposing a novel CNN architecture with multiple convolutional layers, pooling layers, and fully connected layers. The approach incorporates ReLU activation functions and dropout regularization to enhance performance and mitigate overfitting. Extensive experiments on CIFAR-10 and ImageNet datasets demonstrate the superiority of our CNN in terms of classification accuracy. The learned hierarchical representations capture discriminative features, contributing to improved performance.

Task Description

Building a machine learning model capable of detecting the emotions expressed by a person in an image using Python. Emotion detection from images is an exciting and challenging problem in the field of computer vision and artificial intelligence. The goal is to leverage the power of machine learning algorithms to accurately predict the emotional state of individuals.

Dataset

The data consists of 48x48 pixel grayscale images of faces. The faces have been automatically registered so that the face is centered and occupies about the same amount of space in each image.

The task is to categorize each face based on the emotion shown in the facial expression into one of seven categories (0=Angry, 1=Disgust, 2=Fear, 3=Happy, 4=Sad, 5=Surprise, 6=Neutral). The training set consists of 28,709 examples and the public test set consists of 3,589 examples.

Table of Contents

Abstract	1
Task Description.....	1
Dataset	1
Imports.....	3
EDA	3
Methods & Terminology.....	5
1. Handling Class Imbalance.....	5
2. Data Augmentation.....	6
3. Sequential	6
4. Optimizers	7
5. Imshow()	8
6. Confusion Matrix	8
Implementation	9
TRIAL 1: Baseline Model	9
TRIAL 2: Increased Complexity of Model	11
Transfer Learning.....	14
Conclusion	17
Future Work.....	18
References	18

Imports

```
import pandas as pd
import math
import numpy as np
import time
import matplotlib.pyplot as plt
from skimage.io import imread, imshow
import tensorflow as tf
from tensorflow.keras.models import Model
from sklearn.model_selection import train_test_split
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from sklearn.metrics import f1_score, precision_score, recall_score, confusion_matrix, classification_report
from tensorflow.keras.utils import to_categorical
from tensorflow import keras
from tensorflow.keras import layers
from tensorflow.keras.layers import Flatten, Dense, GlobalAveragePooling2D
from tensorflow.keras.models import Sequential
```

EDA

Started with Exploratory Data Analysis to check what kind of data we will be dealing with and how to prepare the data for the model.

Dataset Information:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 35887 entries, 0 to 35886
Data columns (total 3 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   emotion     35887 non-null  int64
 1   pixels      35887 non-null  object
 2   Usage       35887 non-null  object
dtypes: int64(1), object(2)
memory usage: 841.2+ KB
```

The dataset consists of 35887 instances which represent grayscale images of faces with the emotions associated with those images.

At any given instance in the dataset, the size of image is **48 x 48**.

The dataset consists of 7 emotions (0=Angry, 1=Disgust, 2=Fear, 3=Happy, 4=Sad, 5=Surprise, 6=Neutral)

```
1 df['emotion'].unique()

array([0, 2, 4, 6, 3, 5, 1])
```

Checking for Data Distribution :

From this analysis we can notice that the samples for happy emotion are in excess and for disgust emotion the samples are very limited which indicates a class imbalance.

```
1 df['emotion'].value_counts().max()
```

8989

```
1 df['emotion'].value_counts().min()
```

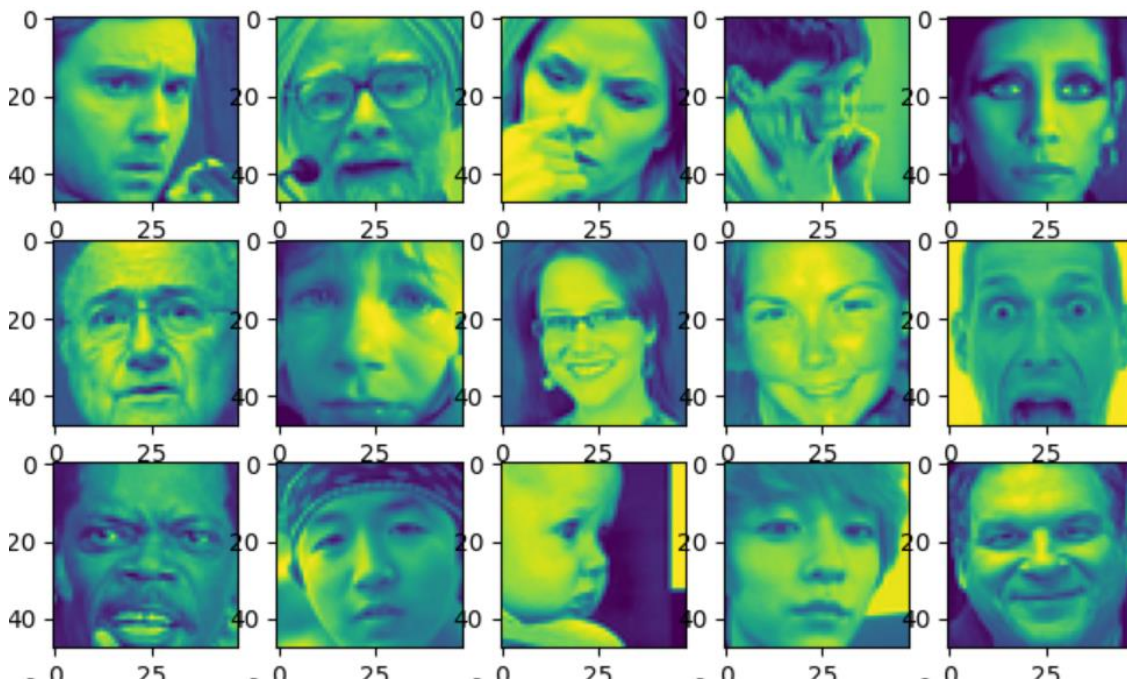
547

Checking for any null values to ensure proper learning for model:

```
1 df['emotion'].isnull().sum()
```

0

Sample Data:



Methods & Terminology

1. Handling Class Imbalance

compute_class_weight() : This function from scikit-learn is used to compute class weights for imbalanced classification problems. It calculates the weights inversely proportional to the class frequencies, allowing the model to give more importance to the minority classes during training.

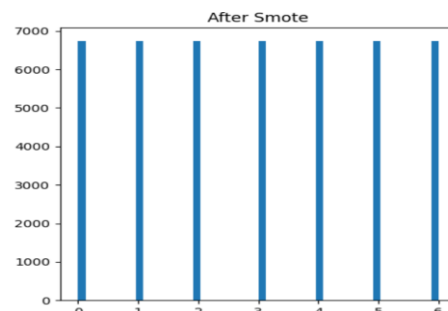
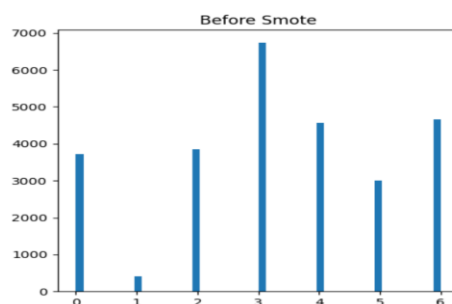
The 'balanced' strategy adjusts the weights such that they are inversely proportional to the class frequencies in the input data.

SMOTE[Not Implemented due to poor performance after trials] (Synthetic Minority Over-sampling Technique) is a popular algorithm used for oversampling the minority class in imbalanced datasets. It generates synthetic samples by interpolating between existing minority class samples.

SMOTE works by creating synthetic samples in the feature space of the minority class.

It selects a minority class instance and finds its k nearest neighbors.

It then synthesizes new instances by interpolating between the selected instance and its nearest neighbors.



2. Data Augmentation

tf.keras.preprocessing.image.ImageDataGenerator()

Efficiently loading, augmenting and preprocessing image data during the training of deep learning models.

shear_range: Shear transformation tilts the image along a particular axis, deforming it in a linear fashion.

zoom_range : Magnifies the Image.

width_shift_range : The width shift range parameter controls the range of horizontal (left or right) translations applied to the images.

height_shift_range: Like the width shift range, the height shift range parameter controls the range of vertical (up or down) translations applied to the images,

rotation_range : A rotation range of 30 means the image can be randomly rotated between -30 and +30 degrees.

3. Sequential

Sequential model is a linear stack of layers used for building deep learning models. It allows you to create a neural network by stacking multiple layers sequentially, where the output of one layer serves as the input to the next layer.

BatchNormalization()

Normalizes the activations of a neural network layer by subtracting the batch mean and dividing by the batch standard deviation. By normalizing the activations, Batch Normalization reduces the dependence of the model on the initialization and learning rate, leading to faster convergence.

Conv2D()

The convolution operation involves sliding a small filter (also known as a kernel or window) over the input data. At each position of the sliding window, an element-wise multiplication is performed between the filter and the corresponding input data patch. The results of these multiplications are summed up, producing a single value that represents the output of the convolution operation at that position.

kernel_size defines the size of the filters.

Padding adds extra border pixels to the input data, which can help preserve spatial dimensions and reduce information loss at the borders.

Stride defines the step size of the sliding window during the convolution operation. A stride of 1 moves the window one pixel at a time, while a larger stride value skips pixels.

Activation function to introduce non-linearity and increase the model's expressive power. This can be specified using the **activation** parameter of the **Conv2D** layer, such as 'relu' for the rectified linear unit (ReLU) activation.

MaxPooling2D() a downsampling operation commonly used in convolutional neural networks (CNNs) to reduce the spatial dimensions of feature maps while retaining the most prominent features. For each pooling window, MaxPooling2D outputs the maximum value within that region while discarding other values.

4. Optimizers

The **Adam optimizer** adapts the learning rate for each parameter during training. It incorporates the concept of momentum, which helps accelerate convergence and navigate narrow ravines or flat regions of the loss function. The Adam optimizer maintains two moving average estimates: the first moment (mean) and the second moment (variance) of the gradients. The first moment captures the average direction of the gradients, while the second moment captures the scaling of the gradients.

The **SGD** is an iterative algorithm that gradually adjusts the model's parameters to minimize the loss function. By repeatedly updating the parameters using mini-batches, SGD approximates the gradient of the entire dataset, making it computationally efficient.

5. Imshow()

For grayscale images, the **imshow** function displays the image with a colormap that represents different intensity values as different colors. By default, it uses a colormap that ranges from black (minimum intensity) to white (maximum intensity).

6. Confusion Matrix

	Predicted Negative	Predicted Positive
Actual Negative	TN	FP
Actual Positive	FN	TP

Visualizes the performance of a classification model by summarizing the predictions and actual class labels.

Accuracy measures the overall correctness of the model's predictions and is calculated as $(TP + TN) / (TP + TN + FP + FN)$

Precision quantifies the proportion of correctly predicted positive instances among all instances predicted as positive. It is calculated as $TP / (TP + FP)$

Recall measures the proportion of correctly predicted positive instances among all actual positive instances. It is calculated as $TP / (TP + FN)$

Specificity measures the proportion of correctly predicted negative instances among all actual negative instances. It is calculated as $TN / (TN + FP)$

F1 score is the harmonic mean of precision and recall, providing a balanced measure of the model's performance. It is calculated as $2 * (precision * recall) / (precision + recall)$

Implementation

TRIAL 1: Baseline Model
Goal: To run an initial CNN learning on FER Dataset

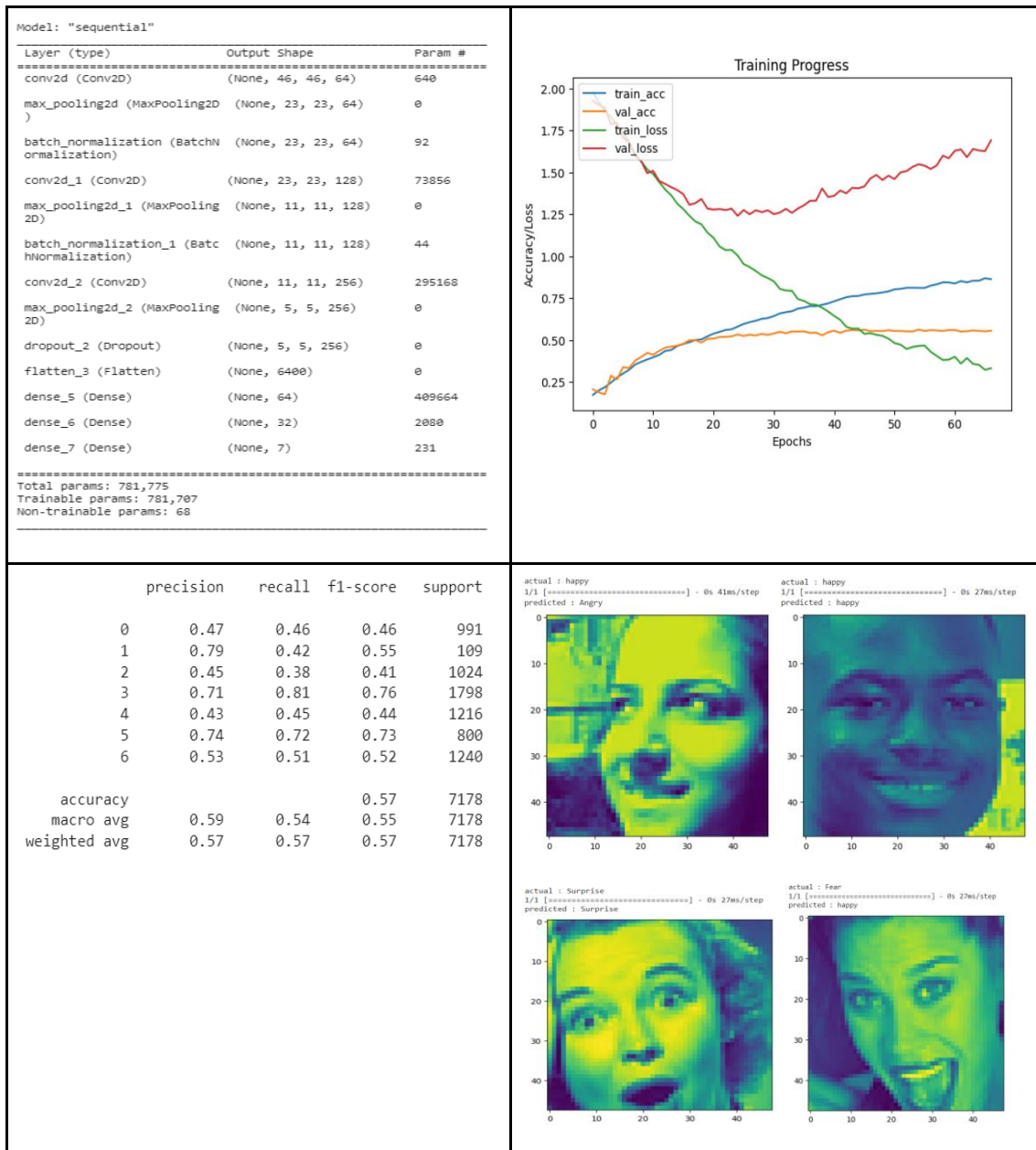
Dataset	Facial Expression Recognition(FER)Challenge Kaggle
Image size	48 x 48 x 1

Parameters	<ol style="list-style-type: none">1. optimizer = 'adam'2. Learning rate = 0.0013. Loss = 'sparse_categorical_crossentropy'4. Batch size = 64 , Epochs = 1005. Imbalanced data handled by weight balancer
------------	--

Insights:

1. Trained a baseline model with 3 convolutional layers with kernel size(3,3) followed by max-pooling layer with activation function set to relu with and callbacks set for early stopping if no improvement after 12 epoch training will be stopped.
2. Handled class imbalance with weights according to the proportion of the samples in the dataset, following are the weights per class used while training :
{0: 1.034993270524899, 1: 9.378048780487806, 2: 1.0010413954699298, 3: 0.5703055473153367, 4: 0.8435717419921018, 5: 1.2812395868043986, 6: 0.8272375215146299}
3. The callback was triggered after 65 + epochs .
4. After evaluation with **F1-score** metric the overall accuracy gained was **0.57** on test set and according to Loss and Accuracy curves , it shows although the training loss was converging, and the validation loss started increasing after **25-30** epochs and this eventually leads the model to be **Overfit** and model stops learning and the same is concluded from accuracy curves as well.

- The reason for this **Overfit** is exclusion of regularization from middle layers of model and **dropout** was only included in last layer.
- The following visuals are the representation of model parameters, Training graphs, F1 score and some predictions.



Confusion Matrix:

This shows the predictability of each class, and we can conclude that the class with most samples (happy) has a better performance than other classes and the class with the lowest samples (Disgust) has the lowest performance.

Confusion Matrix

	0	1	2	3	4	5	6
0	451	4	94	118	185	29	110
1	13	46	7	13	19	4	7
2	143	3	394	88	198	101	97
3	72	1	55	1460	78	36	96
4	140	3	149	156	543	9	216
5	38	1	77	47	29	579	29
6	98	0	106	169	210	22	635
Actuals	Predictions						

TRIAL 2: Increased Complexity of Model

Goal: To run an CNN learning for better performance than baseline model

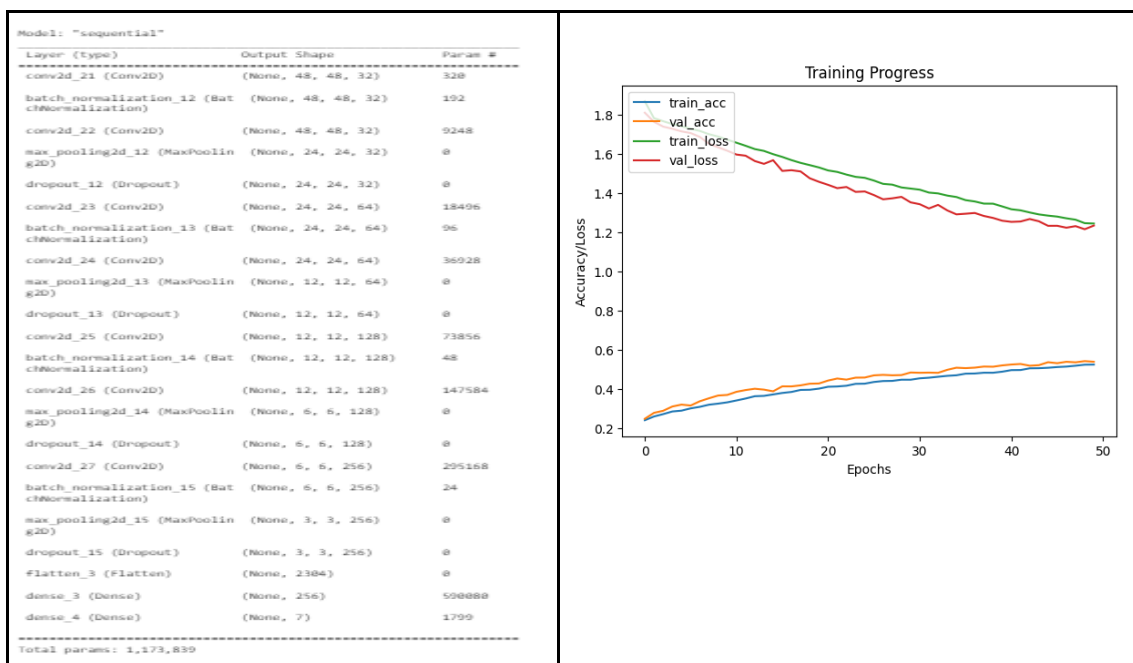
Dataset	Facial Expression Recognition(FER)Challenge Kaggle
Image size	48 x 48 x 1

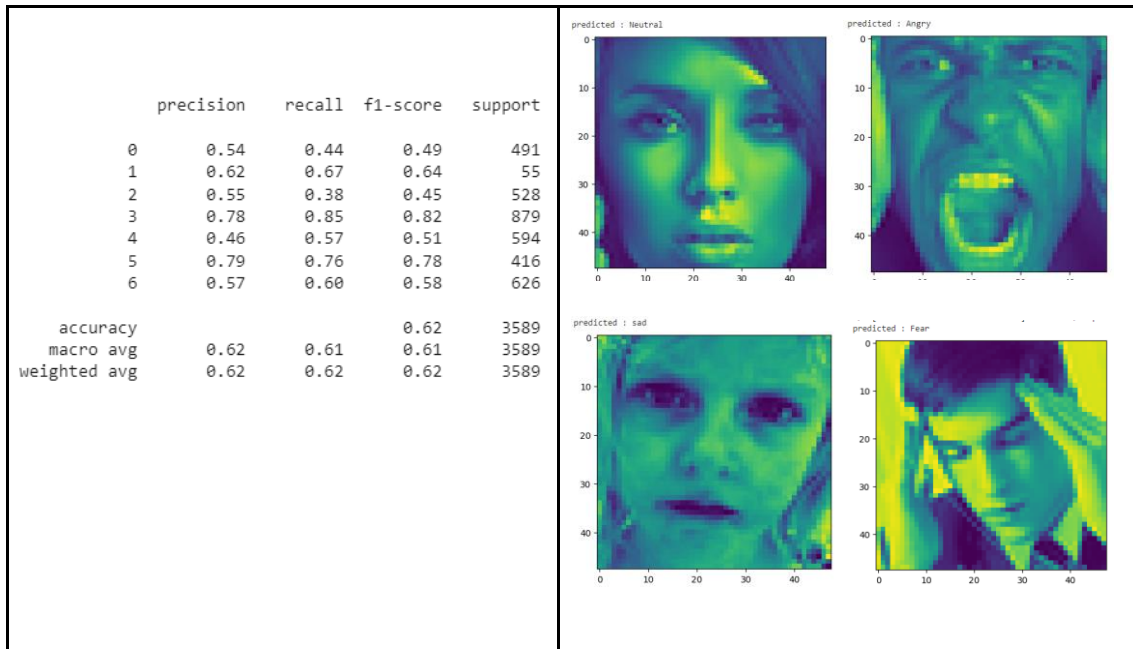
Parameters	<ol style="list-style-type: none"> optimizer = 'adam' Learning rate = 0.001 Loss = 'categorical_crossentropy'
------------	--

4. Batch size = 128 , Epochs = 50
5. Data Augmentation

Insights:

1. Trained a complex model with 7 convolutional layers with kernel size(3,3) followed by max-pooling layer with activation function set to relu and callbacks set for early stopping if no improvement after 12 epochs the training will be stopped.
2. Performed **Data Augmentation** before passing the data to model as the model complexity is now increased and there may be information loss while training and also it helpful to increase **predictability** if model comes across some unseen and real time data.
3. After evaluation with **F1-score** metric the overall accuracy gained was **0.62** which compared to baseline model is **5 %** more and according to Loss and Accuracy curves , it shows the training loss are moving towards convergence and is not **overfitting** as much compared to the baseline model and the accuracy curves depicts the same. This is because of the inclusion of dropouts in middle layers.
4. The following visuals are the representation of model parameters, Training graphs, F1 score and some predictions.





Confusion Matrix:

This shows the predictability of each class, compared to the baseline the overall predictability is increased significantly and there is difference in samples as I used different data split for TRIAL 2 to check if there are any improvements after using original data distribution.

		Confusion Matrix						
		0	1	2	3	4	5	6
Actuals	0	218	9	40	52	103	10	59
	1	7	37	3	4	3	0	1
	2	57	5	199	41	108	41	77
	3	20	1	12	751	45	20	30
	4	44	2	52	46	339	6	105
	5	16	2	26	21	17	318	16
	6	44	4	29	45	122	6	376
		Predictions						

Transfer Learning
Goal: To run a VGG19 Model with Imagenet dataset weights to check for better performance compared to the normal models

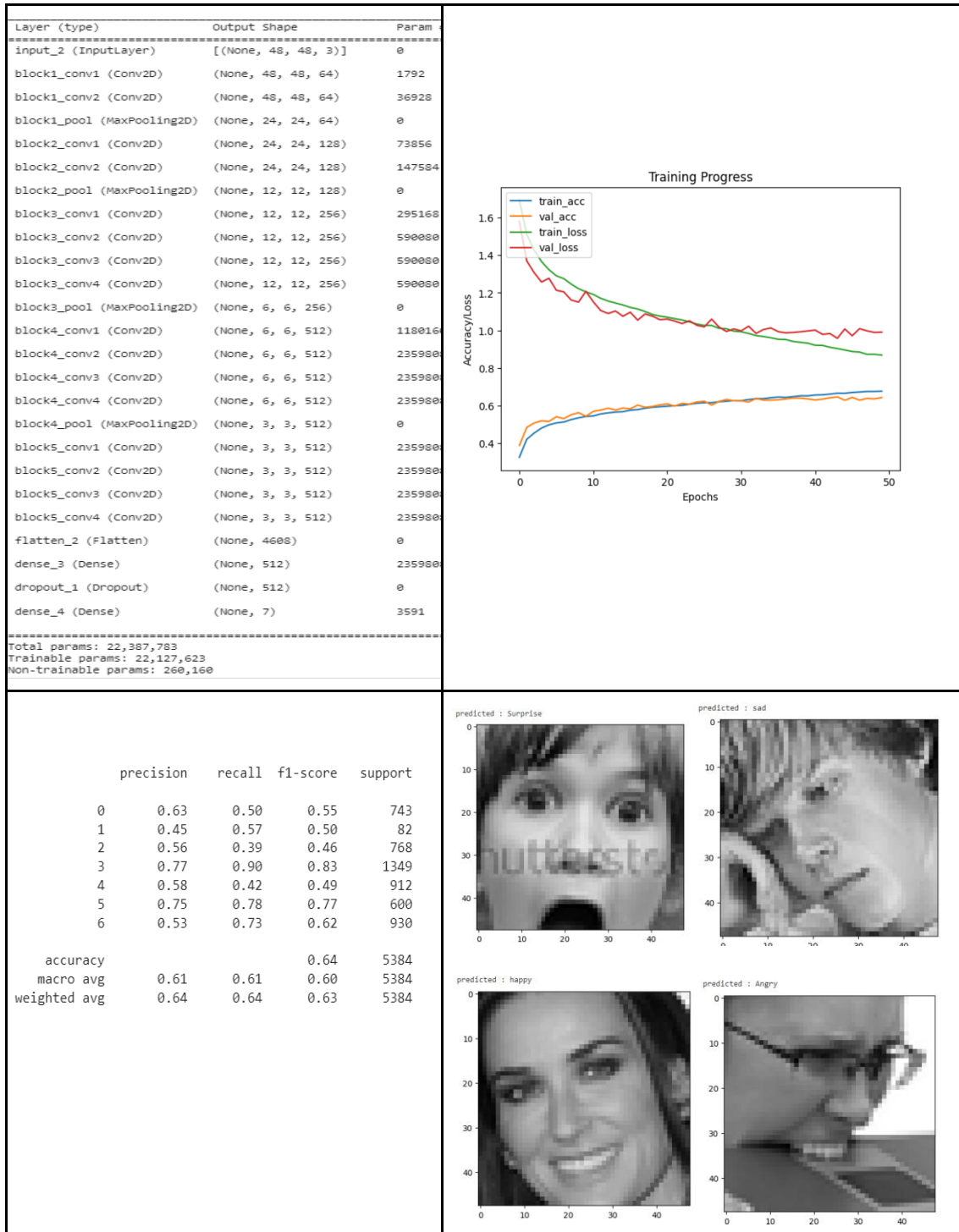
Dataset	Facial Expression Recognition(FER)Challenge Kaggle
Image size	48 x 48 x 3 (Modified for Transfer Learning models)

Parameters	<ol style="list-style-type: none">1. optimizer = 'SGD'2. Learning rate = 0.0013. Loss = 'categorical_crossentropy'4. Batch size = 32 , Epochs = 505. Data Augmentation
------------	--

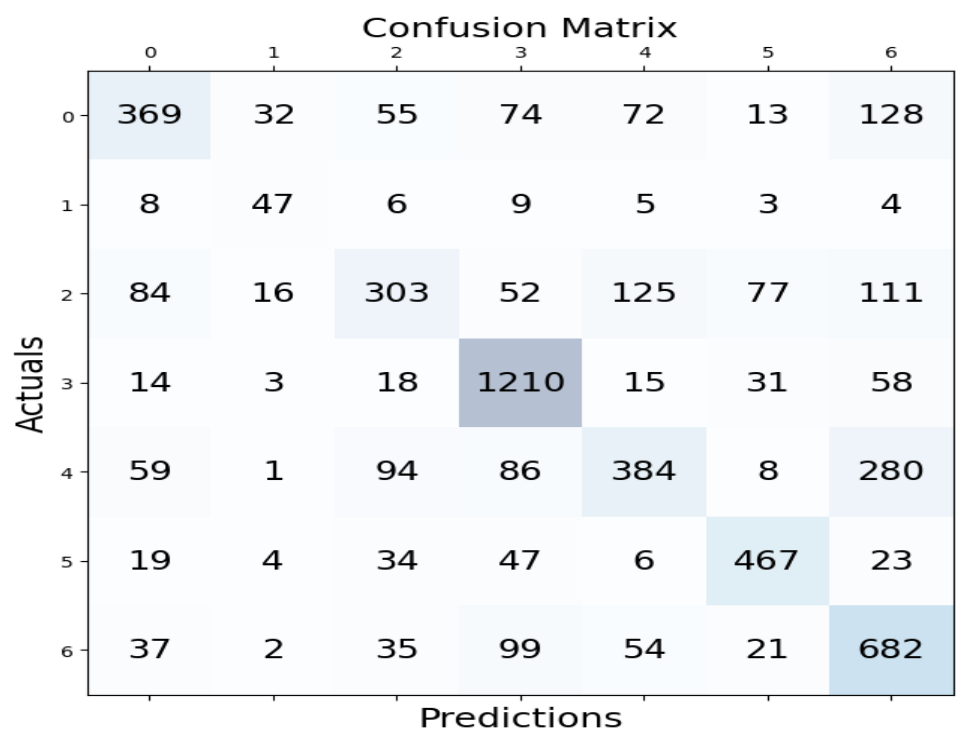
Insights:

1. Trained a model with VGG19 backbone and froze the first two CNN blocks to retain the **weights** of **Imagenet** and other layers will be trained by the FER dataset as Imagenet is more optimized for **Face** recognition than the **Face Expression** recognition and callbacks set for early stopping if no improvement after 12 epoch training will be stopped.
2. Performed **Data Augmentation** before passing the data to model as the model complexity is now increased and there may be information loss while training and also it helpful to increase **predictability** if model comes across unseen and real time data.
3. After evaluation with **F1-score** metric the overall accuracy gained was **0.64** which compared to **TRIAL 2** model is slight better and according to Loss and Accuracy curves , it shows a slight increase in **overfitting** at the end as validation loss is saturated and training loss is still moving towards convergence.
4. I have trained till 50 epochs due to limited GPU resources in Colab and avoiding runtime disconnection while training the model.
5. Model predictability also increased compared to the previous models as shown below.

6. The following visuals are the representation of model parameters, Training graphs, F1 score and some predictions which shows increase in performance.



Confusion Matrix:



Conclusion

- This report has provided a comprehensive overview of Convolutional Neural Networks (CNNs) and its application on **FER dataset**.
- Explored the basic building blocks of CNNs, such as convolutional layers, pooling layers, and fully connected layers, were examined, along with how these elements affect the network's capacity to learn and recognize complicated patterns.
- Best performing model with **64%** F1- score was achieved and good predictability on the test set as well and an average of performance of **61%** was achieved while experimenting with all the 3 models above.
- Data Distribution of classes is a very critical in improving overall predictability and the dataset had class imbalance and should be handled to achieve a well generalized model.
- Due to computation limitations the models were trained for less epochs for Trial 2 and Transfer learning because of their complexity , if trained for more epochs the models can perform better comparatively and learn all the parameters smoothly.
- It's crucial to understand that CNNs have their limitations. They make training and deploying large-scale models difficult since they need a lot of computational resources. Furthermore, CNNs are vulnerable to adversarial assaults, wherein slight changes to the input data might result in classification errors. By investigating model compression approaches, transfer learning, and enhancing resilience against adversarial attacks, researchers and practitioners are actively attempting to address these difficulties.

Future Work

- Working on data distribution of models with **SMOTETomek** which is a hybrid data sampling technique used in machine learning to address the issue of imbalanced datasets. It combines two methods: Synthetic Minority Over-sampling Technique (SMOTE) and Tomek links.
- Exploring on Generative Deep learning to generate reconstructed samples for better performance and generalization of model with reference to https://keras.io/examples/generative/vq_vae/
- **Bonus Task:** Planning to implement a model with enhancements from above techniques for emotion detection using video inputs with approach referenced to <https://github.com/JustinShenk/fer> . I have Implemented the similar in my academic project which included the processing of video to images and building a model for **Image Segmentation**.
- **GridSearchCV** for hyperparameter optimizations and to find the best parameters such as batch size , learning rates , activation functions , optimizers.
- Exploring ways to implement **Convolutional Transformers** and its benefits include parameter-efficiency and self-attention to process long-range and interactions between different regions in an image.

References

- <https://keras.io/api/applications/vgg/>
- <https://machinelearningmastery.com/smote-oversampling-for-imbalanced-classification/>
- <https://www.tensorflow.org/>
- <https://scikit-learn.org/stable/index.html>
- <https://imbalanced-learn.org/stable/references/generated/imblearn.combine.SMOTETomek.html>
- <https://keras.io/examples/vision/cct/>