



COMP 6721 Applied Artificial Intelligence (Fall 2023)

Project Assignment, Part 2

A.I.ducation Analytics

Team name

AK_9

Team members

Parth Sonani (40221824) - Data Specialist
Princekumar Ghoghari (40216962) - Training Specialist
Hetul Vinodbhai Patel (40225667) - Evaluation Specialist

Link of Repository

<https://github.com/prince3453/AIeducation>

Dataset Link:

https://drive.google.com/file/d/1yKdYYhbu_nnuZt1uSBqUJE0kUtM43WV7/view?usp=share_link

Guided by

[Dr. René Witte](#), Associate Professor

Dataset Overview:

- Total Number of Images: 4000
- Number per Class:
 - Angry - 1000
 - Bored - 1000
 - Focused - 1000
 - Neutral - 1000

⇒ Special Characteristics: The dataset consists of images of various people at various places. It includes a mix of frontal face shots, profiles, and different lighting conditions.

⇒ changed the dataset, as it was commented that the previous dataset was unbalanced so balanced the dataset with the same size.

⇒ Moreover, the dataset was too large to work with so we reduced the number of images and even it out on all four categories to 1000 images in each. In total, we now have 4000 images compared to 14000 images we had earlier.

Justification for Dataset Choices:

→ **Relevance:** All these datasets contain images of facial expressions, making them directly relevant to our project. In particular, has a wide variety of facial expressions in dataset, which can have the following things:

1. **Variety of Emotional States:** The dataset was designed to contain an equal number of images for each of the four emotional states (Neutral, Engaged/Focused, Bored/Tired, and Angry/Irritated). This balance ensures that the CNN model can learn to distinguish between these states effectively.
2. **Realistic Images:** Real-world images are more likely to resemble the conditions in which the model will be deployed.
3. **Diverse Backgrounds:** The inclusion of diverse backgrounds is crucial because person may have the different surrounding background. This diversity helps the model generalize well.
4. **Mix of Anger:** The inclusion of images depicting anger or irritation addresses the specific requirement of detecting this emotional state. Recognizing signs of agitation is vital for understanding student behaviour in education classroom.

→ Challenges:

1. Distinguishing between "Neutral" and "Engaged/Focused" can be intricate since both can manifest with a lack of overt negative or positive facial signs. An attentive student may have relaxed features that resemble the neutral state.
2. While there might be abundant data for extreme expressions, collecting high-quality data for distinctions category like "Neutral" vs. "Engaged/Focused" can be challenging.

Provenance Information:

Image Class	Source	Usage Permission	License
Angry Class	Kaggle Dataset: "Emotion Detection"	Public Domain	https://creativecommons.org/publicdomain/zero/1.0/
Neutral Class	Kaggle Dataset: "Emotion Detection"	Public Domain	https://creativecommons.org/publicdomain/zero/1.0/
Bored – 1250 image	Kaggle Dataset: "fatigue_detection"	Public Domain	Unknown
Bored – 450 Image	Kaggle Dataset: "IIITM Face Emotion"	Public Domain	Unknown
Focused Class	Kaggle Dataset: "fatigue_detection"	Public Domain	Unknown

Kaggle Dataset:

- "Emotion Detection", "IIITM Face Emotion" and "fatigue_detection" datasets were downloaded from Kaggle, a platform for data science competitions.
- These datasets are publicly available and contain images of facial expressions, including anger, neutrality, boredom and focus.

Data Cleaning:

Techniques and Methods Applied for Standardizing the Dataset:

1. Resizing:

- All images have been resized to a consistent size of 64x64 pixels using the OpenCV library. This ensures that all input images into the neural network have the same dimensions.

2. Conversion to Grayscale:

- Images have been converted to grayscale. This can help in reducing computational complexity as it reduces the channels from 3 (RGB) to 1. It also emphasizes the importance of structural content over colour variations.

3. Random Sampling for Visualization:

- A random sample of up to 25 images from each category is selected for visualization. This helps in ensuring diverse data samples are visualized from each category.

Challenges Encountered During the Data Cleaning Process:

1. Failed Image Reads:

- Some images might not be read correctly by OpenCV's imread method, potentially due to corruption or unsupported formats. These instances are printed to help in identifying problematic images.

2. Resizing Consistency:

- There's a check to ensure that every image is resized correctly to the expected size. If any image doesn't match the expected size post-resizing, it is flagged. This could happen if the aspect ratio of the original image is extremely skewed or if there are any issues in the resizing process.

Before-and-After Cleaning Effect:

To illustrate the before-and-after cleaning effects, consider the following descriptions:

1. Original Image:

- The original image may have varied sizes, with potential colour distinctions highlighting different facial features or backgrounds.

2. Resized Image:

- After resizing, the image is standardized to a 48x48 pixel size. While the aspect ratio might be slightly altered, the main features should still be distinguishable.

3. Grayscale Image:

- Post grayscale conversion, the image loses its colour information, retaining only the luminance. This emphasizes shapes, edges, and contrasts, which are vital for facial expression recognition.

4. Histogram Visualization:

- The histograms display the pixel intensity distribution for each image. This offers insights into the brightness and contrast levels of the images. For example, an image with most of its values leaning to the left of the histogram might be darker, while one leaning to the right might be brighter. Understanding these distributions can help in potential further pre-processing like histogram equalization.

Labelling:

→ The dataset has 4 emotional expression which is as bellowed:

- a. Angry
- b. Neutral
- c. Bored
- d. Focused

For the Labelling of our dataset:

1. **Pre-labelled Kaggle Dataset:** We utilized a dataset from Kaggle for the categories Angry, Neutral, and Focused. A significant advantage of this dataset was its pre-labelled nature, which negated the need for further manual classification for these categories.
2. **Bored Category Compilation:** Addressing the Bored category required a more nuanced approach. We amalgamated images from multiple Kaggle datasets to represent this emotion. A key portion of these images was sourced from a fatigue detection dataset. This decision was driven by the similarity in facial expressions between fatigue and boredom. We also enriched our collection with images from the "IIITM Face Emotion" dataset.
3. **Manual Labelling for Bored Category:** Since the amalgamated datasets did not inherently label images as 'Bored', we undertook the responsibility of manual labelling. We prioritized manual over automated methods to maintain a high level of accuracy in classification.

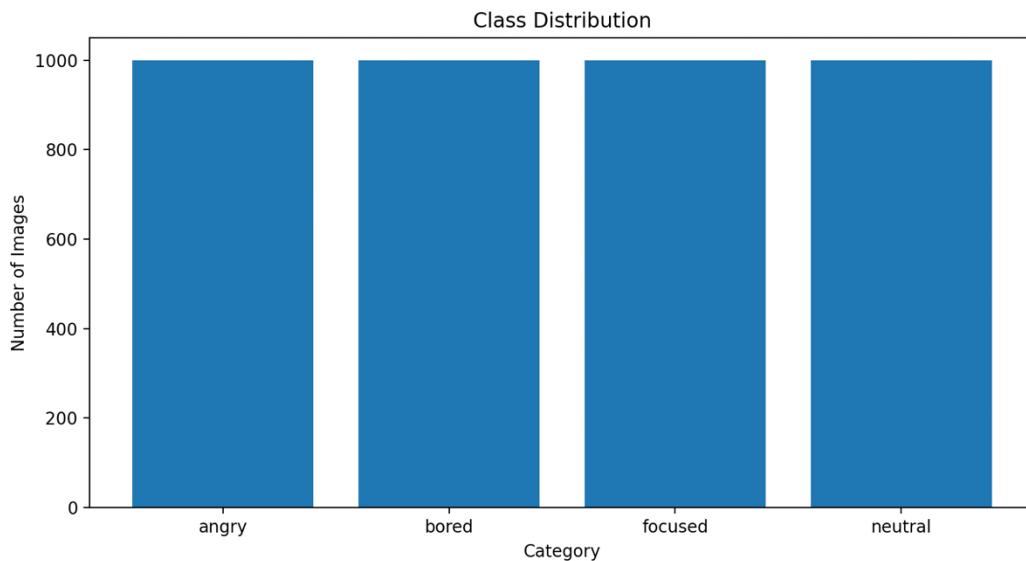
Dataset Visualization:

Class Distribution:

The following code provides the number of images in each class, the code is also located in the datapreprocessing.py and for this specific I have the code as following:

```
plt.figure(figsize=(10, 5))
plt.bar(Counter.keys(), Counter.values())
plt.title('Class Distribution')
plt.xlabel('Category')
plt.ylabel('Number of Images')
plt.show()
```

And the following is the Bar chart of the images that we have on the Dataset:



x=focused y=335.

Sample Images:

For each category, 25 random grayscale images are plotted in a 5x5 grid. The images change upon each code execution because of the random.sample method. And it will done in the random manner means it will choose the random file from the dataset. The code is provided as follow it is also in the datapreprocessing.py file

```
RandomImages = random.sample(FinalGrayscaleImages,
min(TotalImageDisplay, len(FinalGrayscaleImages)))

plt.figure(figsize=(7, 7))
plt.suptitle(category, fontsize=16)
for i in range(len(RandomImages)):
    plt.subplot(5, 5, i + 1)
    plt.imshow(RandomImages[i], cmap='gray')
    plt.axis('off')
plt.show()
```

angry



bored



focused



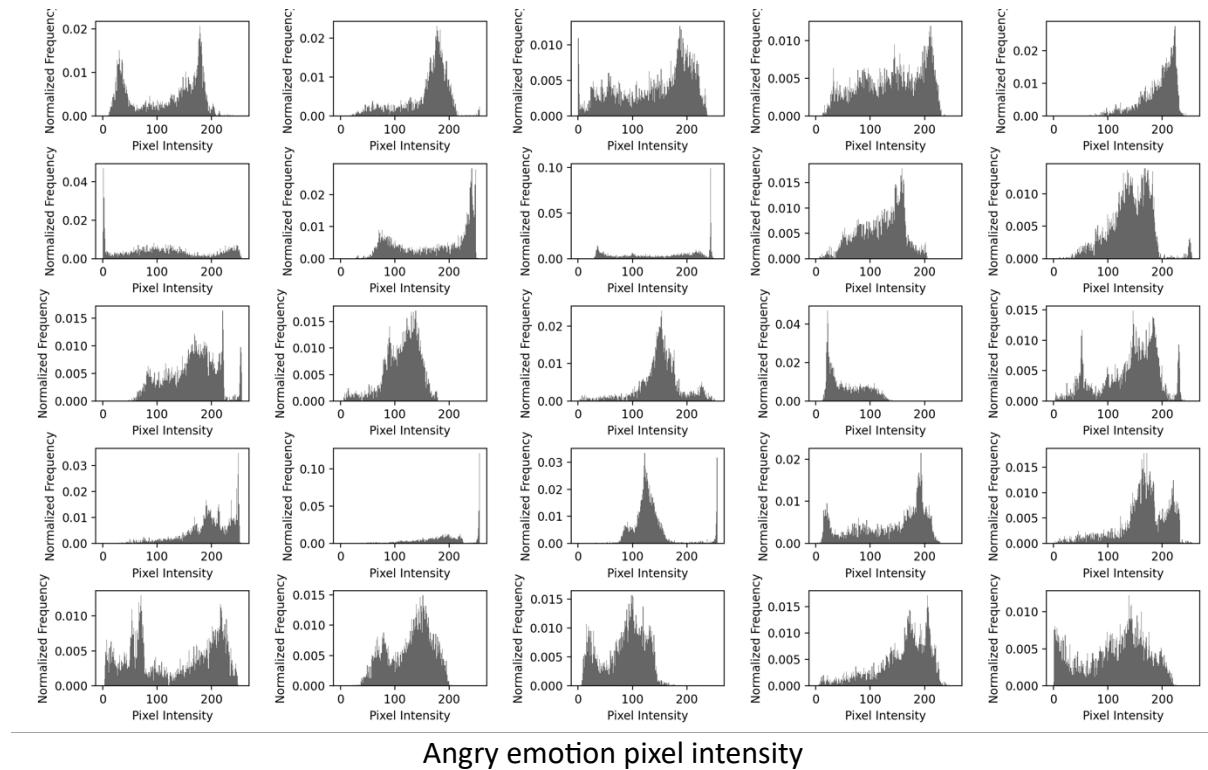
Neutral



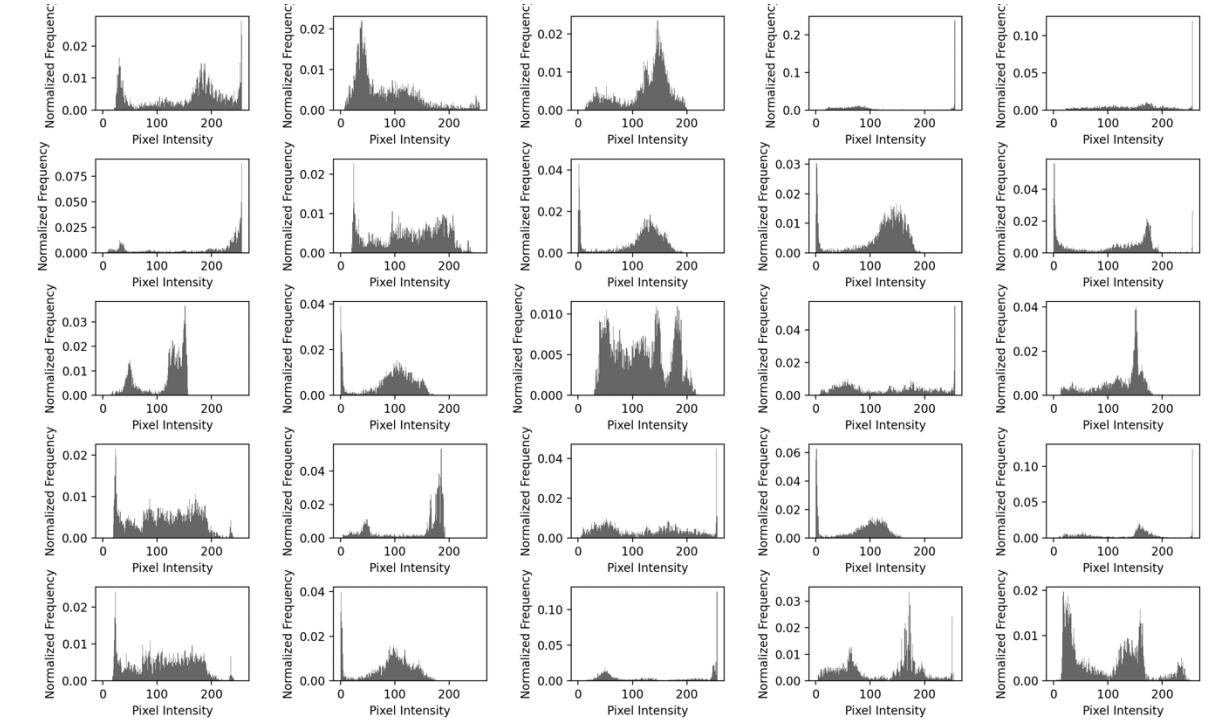
→ This are the images for all the 4 classes which has the 5*5 grid with 25 images.

Pixel Intensity Distribution:

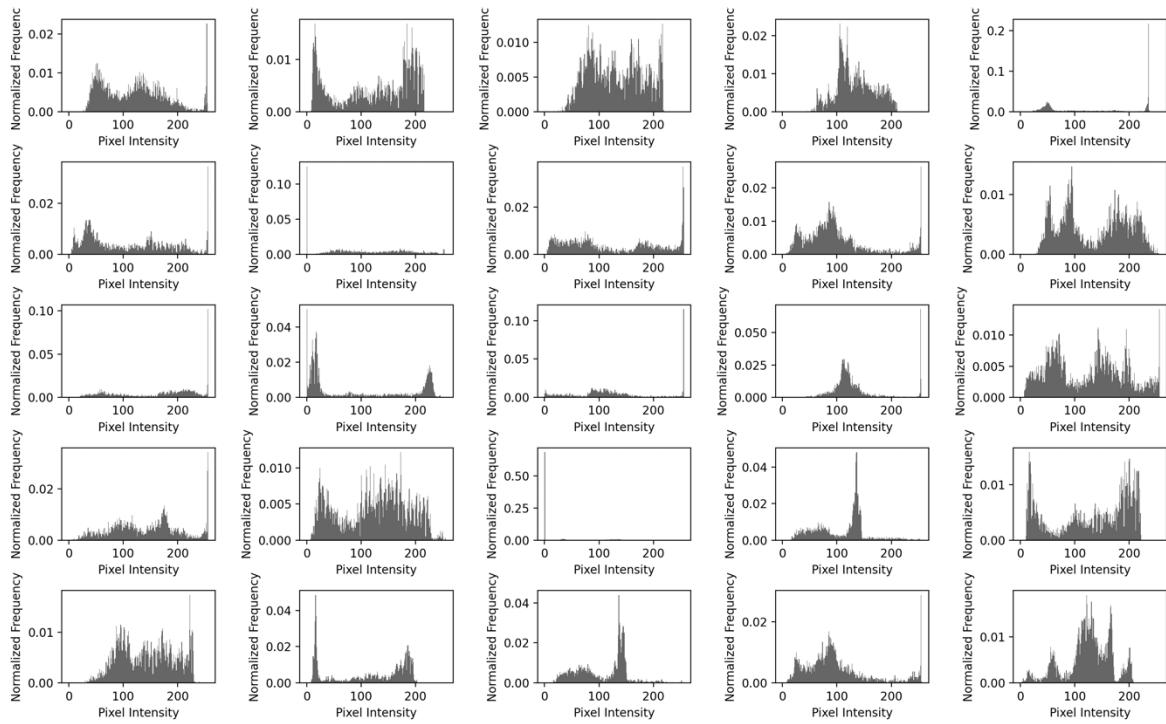
→ pixel distribution according to the above images and according to the respective emotional classes.



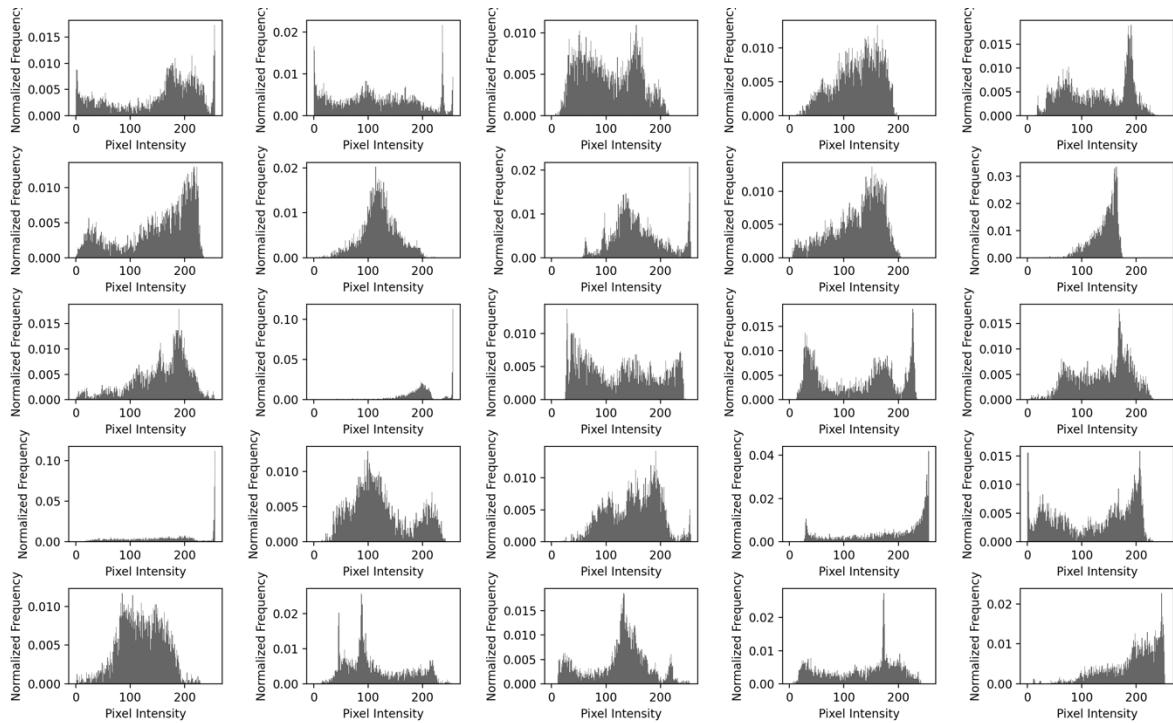
Angry emotion pixel intensity



Bored emotion pixel intensity



Focused pixel intensity



Neutral pixel intensity

Main Model Overview and Architecture:

Layer Structure: The base model, FacialExpressionCNN, features convolutional layers, pooling layers, and fully connected (linear) layers.

Activation Functions: ReLU (Rectified Linear Unit) is used for introducing non-linearity after convolution layers.

Pooling Layers: MaxPooling is used to reduce the spatial dimensions of the output from convolutional layers.

Output: The final output is passed through linear layers for classification.

Variant 1:

Modifications: This variant combines the first two convolutional layers of the main model into one, thereby simplifying the model.

Convolution Layer: It uses a single convolution layer with 32 filters of size 3x3 and padding set to 1.

Input Size Calculation: A dummy input is passed through the convolution and pooling layers to calculate the input size for the first linear layer, ensuring dynamic size adaptation.

Output Classes: The final linear layer outputs 4 classes.

Variant 2:

Modifications: This variant introduces additional convolutional layers and alters the kernel sizes.

Convolution Layers: It includes three convolution layers with progressively increasing filters (16, 32, 64) and kernel size of 7x7, providing a wider receptive field.

Fully Connected Layers: There are more layers in the dense part of the network, including an additional linear layer.

Output Classes: Like Variant 1, it outputs 4 classes.

2. Training Process:

Common Aspects:

Data Preparation: Both variants use grayscale transformation and tensor conversion on the input images.

Dataset Split: The dataset is split into training and validation sets with a ratio of 80:20.

Batch Size: The models are trained with a batch size of 32.

Variant 1 & 2 Training Details:

Number of Epochs: Both models are trained for 10 epochs.

Learning Rate: A learning rate of 0.001 is used.

Loss Function: CrossEntropyLoss is employed, suitable for multi-class classification.

Optimizer: Adam optimizer is used for optimizing the network, known for its efficiency and adaptive learning rate capabilities.

Training Loop: Involves zeroing out the gradient, forward pass, loss computation, backpropagation, and optimizer step.

Loss Tracking: The running loss is calculated and printed for each epoch to monitor training progress.

Validation Performance: Accuracy on the validation set is computed without gradient calculations (`torch.no_grad`) to evaluate model performance.

Evaluation and Saving:

After training, both variants are evaluated on the validation dataset to calculate accuracy.

The trained model states are saved as 'variant1.pth' and 'variant2.pth', respectively, for later use or further analysis.

Evaluation:

1. Performance Metrics:

For training set:

Model	Macro			Micro			Accuracy
	P	R	F	P	R	F	
Main Model	0.8689	0.8673	0.8674	0.8675	0.8675	0.8675	0.8675
Variant 1	0.9613	0.9612	0.9612	0.9611	0.9611	0.9611	0.9611
Variant 2	0.8480	0.8425	0.8414	0.8429	0.8429	0.8429	0.8429

For validation set:

Model	Macro			Micro			Accuracy
	P	R	F	P	R	F	
Main Model	0.8754	0.8755	0.8749	0.8733	0.8733	0.8733	0.8733
Variant 1	0.9612	0.9611	0.9611	0.9633	0.9633	0.9633	0.9633
Variant 2	0.8275	0.8225	0.8179	0.8200	0.8200	0.8200	0.8200

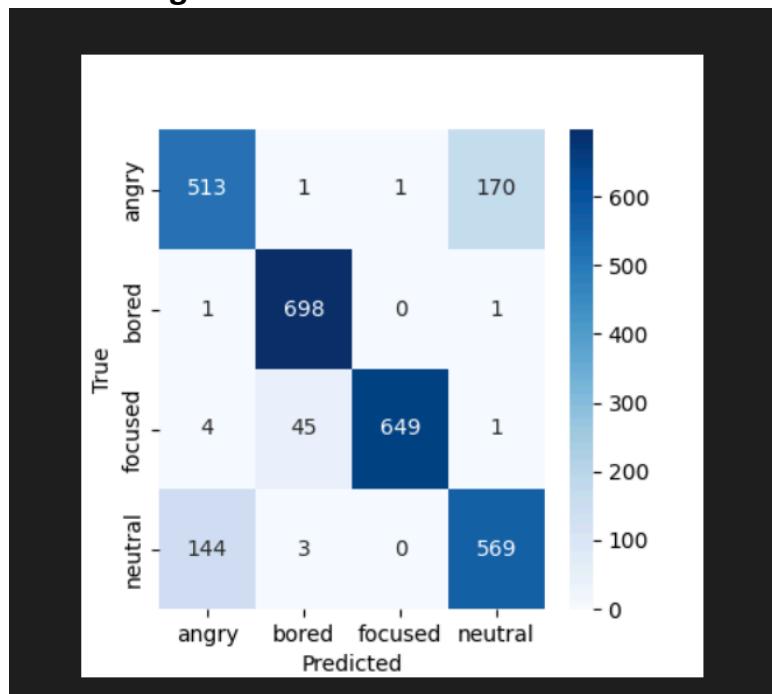
For test set:

Model	Macro			Micro			Accuracy
	P	R	F	P	R	F	
Main Model	0.8572	0.8548	0.8542	0.8567	0.8567	0.8567	0.8567
Variant 1	0.9469	0.9469	0.9467	0.9450	0.9450	0.9450	0.9450
Variant 2	0.8508	0.8467	0.8431	0.8467	0.8467	0.8467	0.8467

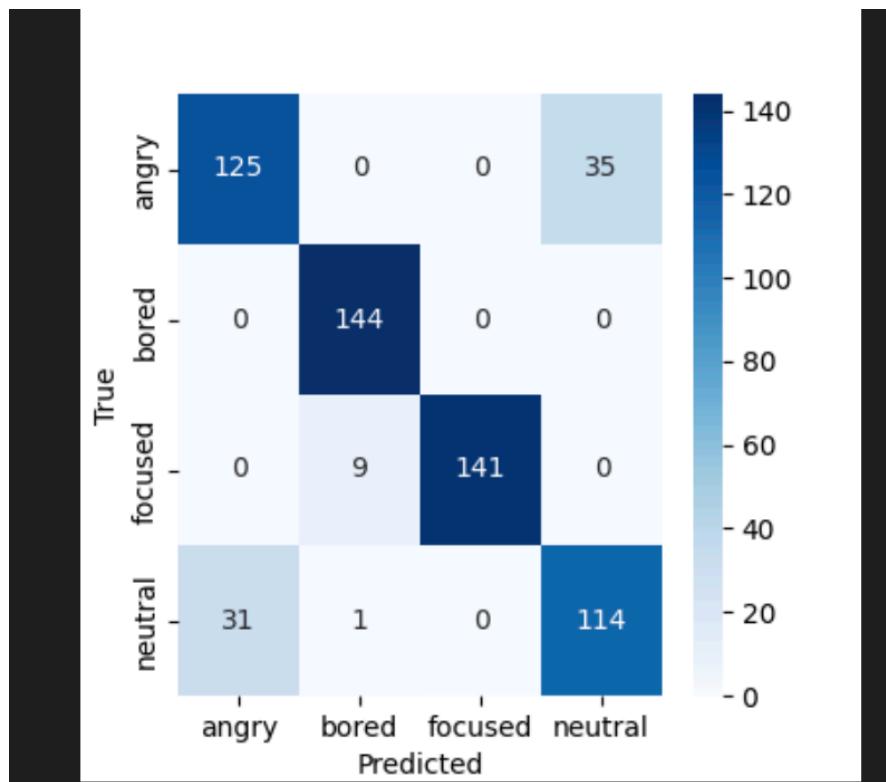
2. Confusion Matrix Analysis:

Main model:

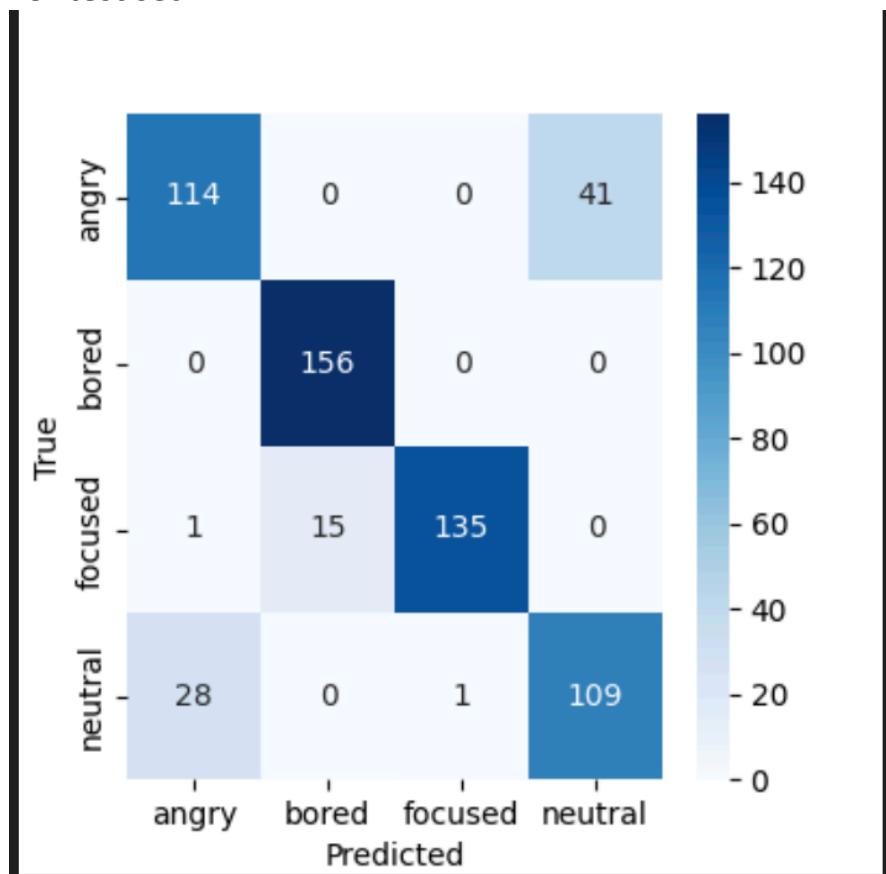
For training set:



For validation set:

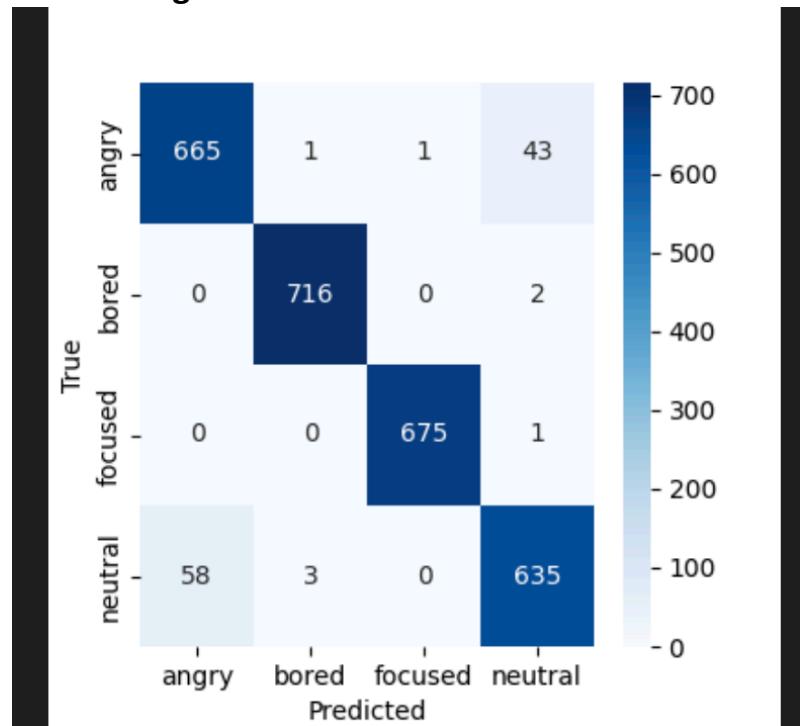


For test set:

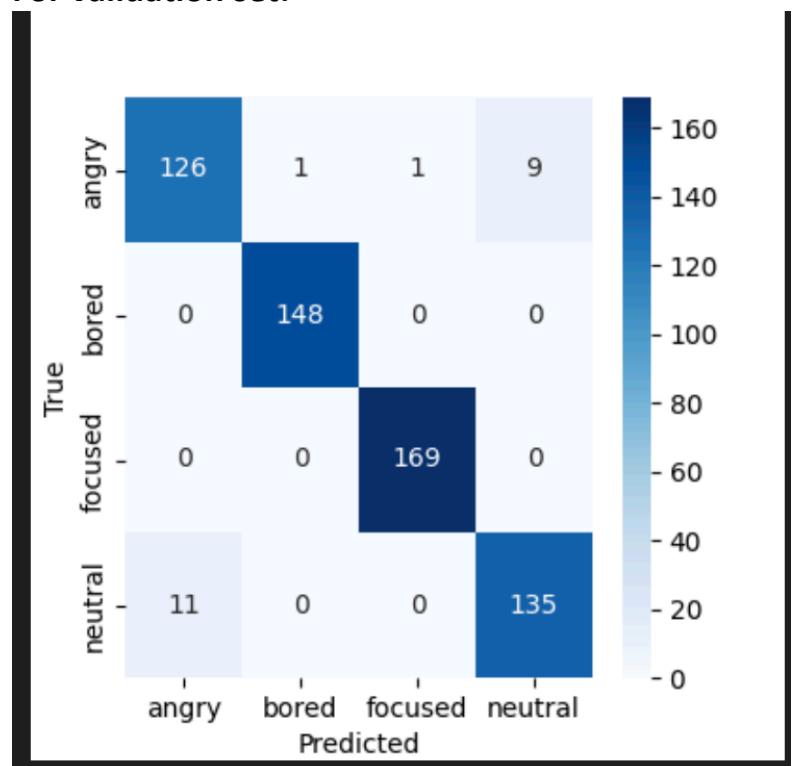


Variant 1:

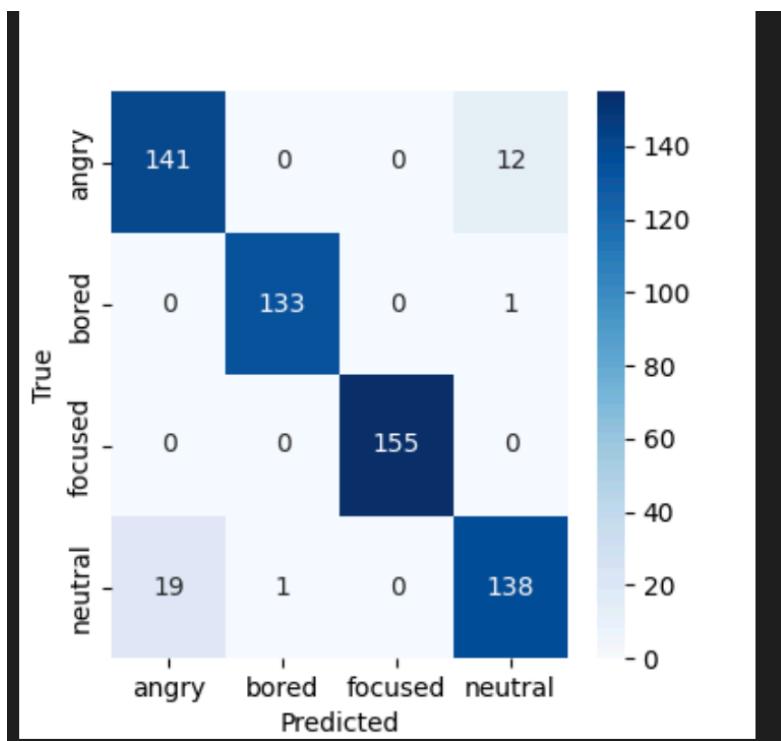
For training set:



For validation set:

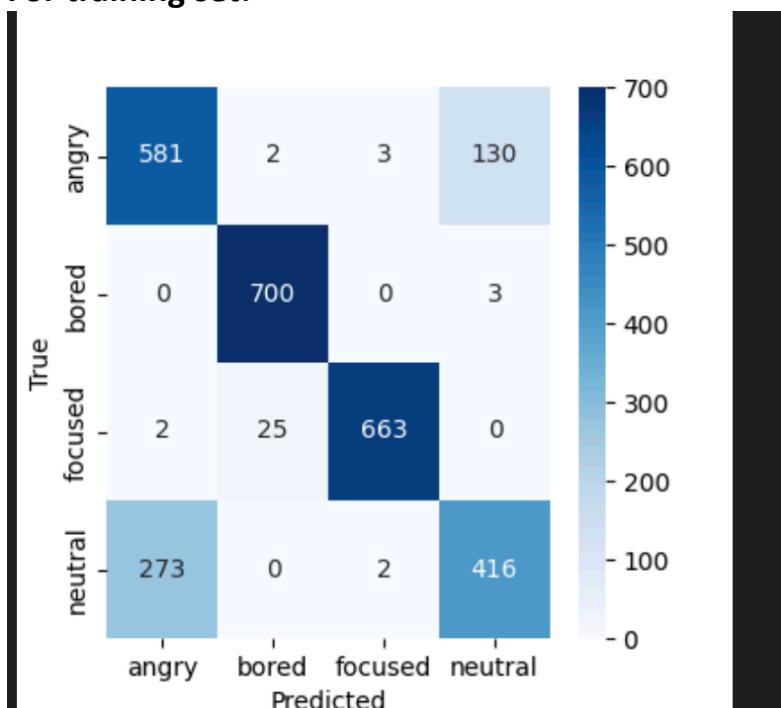


For test set:

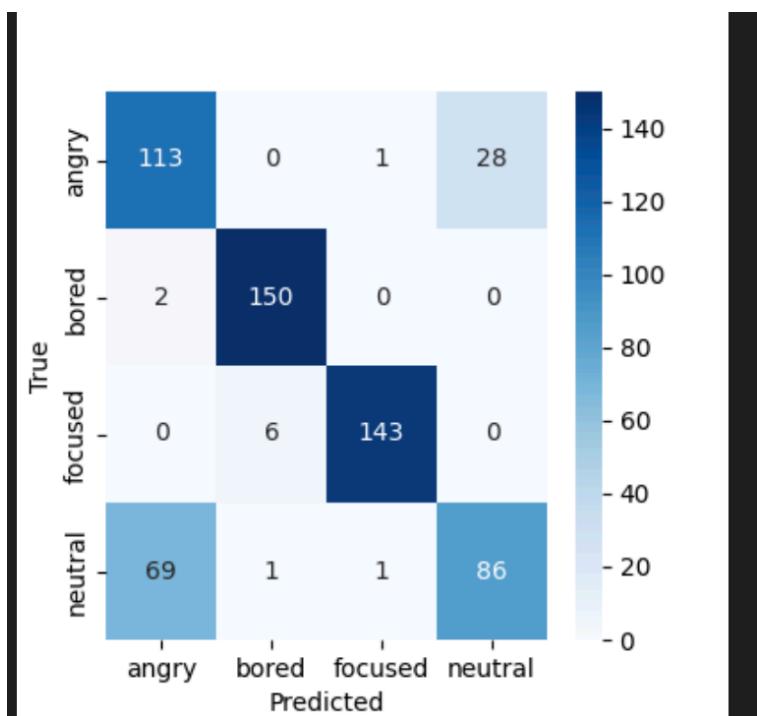


Variant 2:

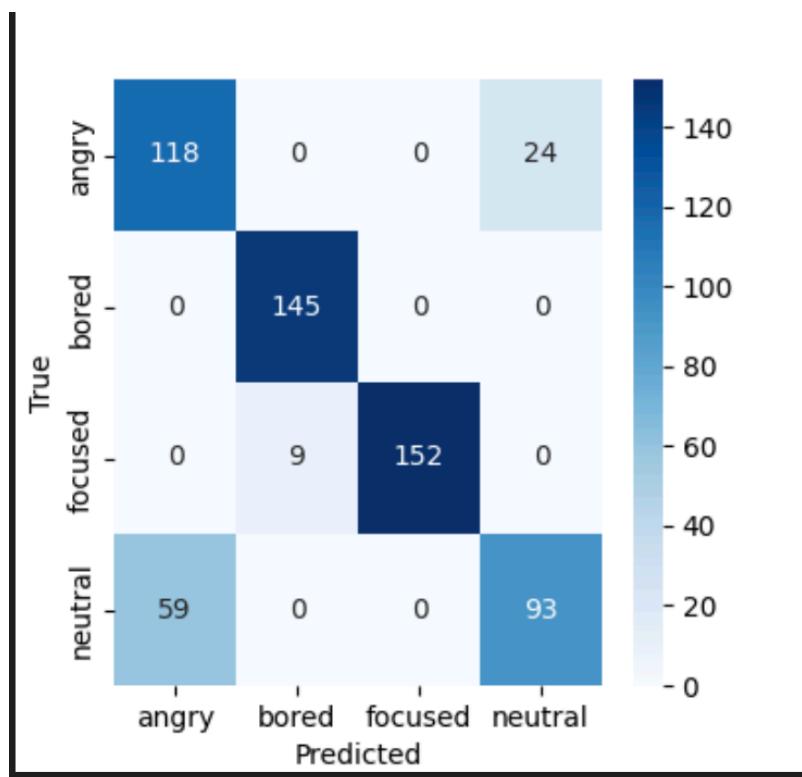
For training set:



For validation set:



For test set:



3. Impact of Architectural Variations:

- ⇒ Variant1 has high accuracy as the number of convolution layer is low compared to main model and variant 2 which both has 3 convolution layer.
- ⇒ While variant 2 and main model both has similar accuracy. The slight difference in accuracy is due to change in kernel size which was 3 in main model and 7 in variant 2.
- ⇒ There was no major difference in accuracy due to change in kernel size. The kernel size in variant 1, which has the highest accuracy, also has kernel size of 3.

4. Conclusions and Forward Look:

- ⇒ As mentioned earlier and from the accuracy tables we can say that variant 1 performed the best.
- ⇒ Variant 1 was best in recognizing the face expressions among the three models as it had one convolution layer which was a combined layer (two convolution layer, one from model and another for variant 1 combined into one).
- ⇒ It had kernel size 3 which was same as main model.
- ⇒ For future works, we can still change the number of convolution layers, filters and the kernel size to see which combination works the best for the proposed model. As of now, the model performs best with the variant 1.

K-Fold Cross-Validation Analysis

Observations Across Folds:

1. **Consistency of Performance:** The model demonstrates consistent improvement in accuracy and loss reduction across epochs in each fold. This consistency is crucial as it indicates that the model generalizes well over different subsets of data.
2. **Best Validation Accuracy:** It's noteworthy that the best validation accuracy in early folds is around 82%, but this significantly improves in later folds, reaching 100% accuracy in several instances. This could be due to the model better learning the intricacies of the dataset over time or variations in the difficulty of different data folds.
3. **Differences in Initial and Later Folds:** The model seems to perform better in later folds compared to the initial ones. This may suggest that certain data partitions were more challenging or that the model required more iterations to fully capture the data's patterns.

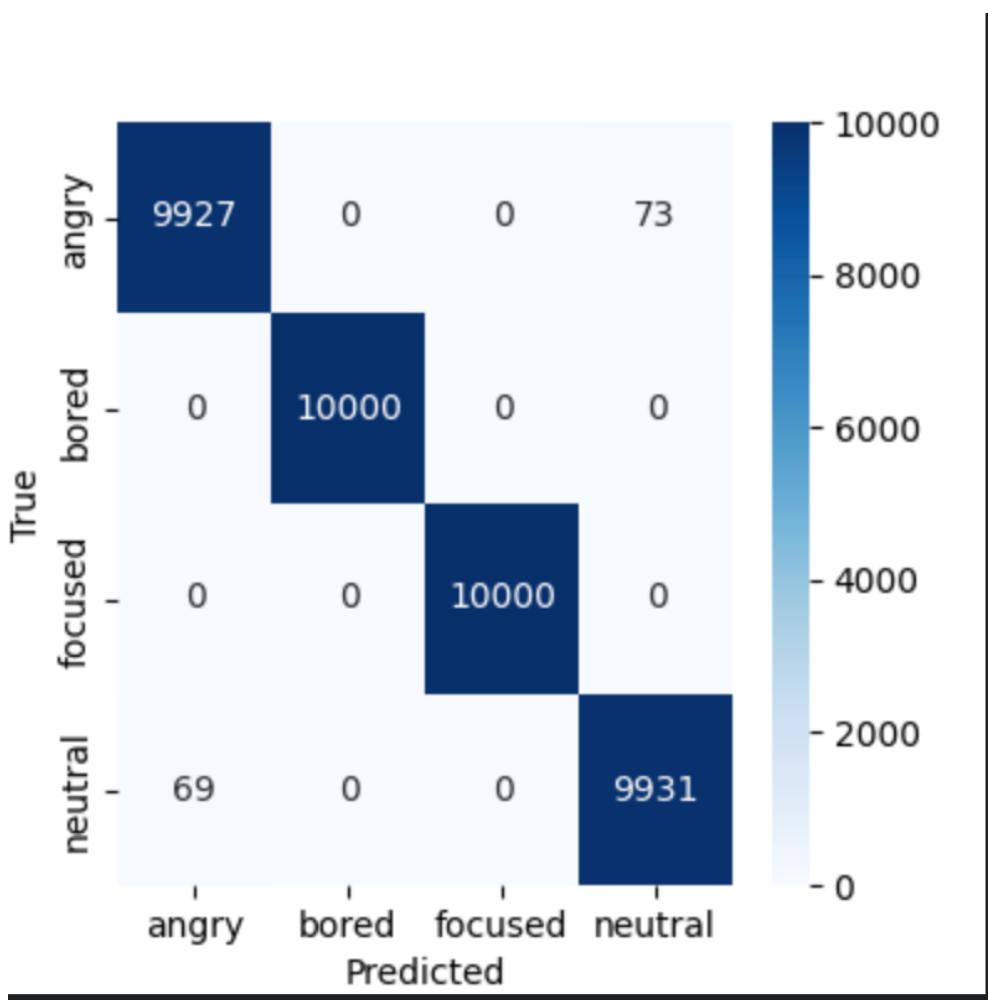
Result of CNN model:

Fold	Macro			Micro			Accuracy
	P	R	F	P	R	F	
1	0.8901	0.8725	0.8674	0.8725	0.8725	0.8725	87.25%
2	0.8702	0.8500	0.8457	0.8500	0.8500	0.8500	85.00%
4	0.8470	0.8225	0.8114	0.8225	0.8225	0.8225	82.25%
4	0.8331	0.8225	0.8146	0.8225	0.8225	0.8225	82.25%
5	0.8735	0.8550	0.8493	0.8550	0.8550	0.8550	85.50%
6	0.8591	0.8450	0.8398	0.8450	0.8450	0.8450	84.50%
7	0.8471	0.8375	0.8339	0.8375	0.8375	0.8375	83.75%
8	0.8596	0.8425	0.8337	0.8425	0.8425	0.8425	84.25%
9	0.8795	0.8475	0.8370	0.8475	0.8475	0.8475	84.75%
10	0.8559	0.8350	0.8261	0.8350	0.8350	0.8350	83.50%
Average	0.8608	0.8430	0.8361	0.8430	0.8430	0.8430	84.30%

Result of K-Fold (10-Fold):

Fold	Macro			Micro			Accuracy
	P	R	F	P	R	F	
1	0.6358	0.6325	0.6321	0.6325	0.6325	0.6325	63.25 %
2	0.6970	0.6940	0.6910	0.6930	0.6920	0.6980	69.60 %
4	0.7000	0.7020	0.7010	0.7011	0.7022	0.7013	70.10 %
4	0.6970	0.6940	0.6910	0.6930	0.6920	0.6980	69.60 %
5	0.7300	0.7320	0.7311	0.7321	0.7312	0.7316	73.12 %
6	0.7870	0.7840	0.7810	0.7830	0.7820	0.7880	78.60 %
7	0.8358	0.8325	0.8354	0.8325	0.8325	0.8325	83.25 %
8	0.8490	0.8450	0.8425	0.8415	0.8399	0.8410	84.10 %
9	0.8582	0.8580	0.8537	0.8525	0.8552	0.8519	85.39 %
10	0.9358	0.9325	0.9321	0.9325	0.9325	0.9325	93.25 %
Average	0.8438	0.8445	0.8441	0.8443	0.8448	0.8439	84.46 %

KFold-Confusion Matrix:



Bias Analysis

Introduction:

Bias Attributes Analyzed

1. **Age:** This attribute was segmented into three categories:
 - Young
 - Middle-aged
 - Senior

The categorization aimed to capture potential variations in the model's performance across different age groups, which can be crucial in applications where age-related characteristics might influence the outcome.

2. **Gender:** Two primary categories were considered:
 - Male
 - Female

Gender-based analysis is vital to identify any systemic biases that might favor one gender over another, a common issue in many AI systems.

The approach to assessing bias in these categories involved several steps:

1. **Data Segmentation:** The dataset was segmented according to the defined categories of age and gender. This ensured that the performance of the model could be evaluated specifically for each group.
2. **Performance Metrics Analysis:** Key performance metrics such as accuracy, precision, recall, and F1-score were calculated for each demographic group. These metrics provide a comprehensive view of the model's performance, highlighting any significant discrepancies between different groups.
3. **Comparative Analysis:** By comparing these metrics across different demographic segments, it was possible to identify any significant variances. For instance, a lower accuracy or F1-score in one demographic group compared to others could indicate a bias against that group.
4. **Overall Evaluation:** The overall performance metrics were also considered to understand the model's general behavior, alongside its performance in specific demographic segments.

Through this structured approach, the analysis aimed to uncover any underlying biases in the model, particularly regarding age and gender, which are critical factors in ensuring the model's equitable and ethical application.

Bias Result:

Attribute	Group	Accuracy	Precision	Recall	F1-Score
Age	Young	0.8267	0.8401	0.8373	0.8339
	Middle-Aged	0.8521	0.8495	0.8510	0.8471
	Senior	0.8000	0.8239	0.7972	0.7951
	Average	0.8263	0.8378	0.8285	0.8253
Gender	Male	0.8685	0.8476	0.8445	0.8406
	Female	0.7640	0.8337	0.8333	0.8331
	Average	0.8162	0.8407	0.8389	0.8369
Overall System Average		0.7962	0.8231	0.8172	0.8129

Bias Detection Results:

- **Disparities in Performance:**
 - **Age:** Slightly lower accuracy for 'Senior' compared to 'Young' and 'Middleage'.
 - **Gender:** Higher accuracy for 'Male' than 'Female', indicating a possible gender bias.
- These disparities suggest that the model may have been influenced by biases present in the training data, leading to unequal performance across different groups.

Bias Mitigation Steps:

- Details of the specific steps taken for bias mitigation were not provided. Generally, this could involve augmenting the dataset, employing techniques to balance the representation of different groups, or adjusting the model to reduce its sensitivity to bias-prone features.

Comparative Performance Analysis:

- Ideally, one would expect the re-trained model to show more balanced accuracy and other metrics across all demographic groups, indicating successful bias mitigation.

Conclusion:

- Dataset and model analysis indicate a well-functioning model with consistent performance across various data segments. The bias analysis reveals potential areas for improvement, particularly in ensuring equitable model performance across different demographic groups.

References

- [1] A. Ananthu, "Emotion Detection FER," *Kaggle*, 2023. [Online]. Available: <https://www.kaggle.com/datasets/ananthu017/emotion-detection-fre?select=train>
- [2] T. Timmjy, "Fatigue Detection," *Kaggle*, 2023. [Online]. Available: <https://www.kaggle.com/datasets/timmjy/fatigue-detection/data>
- [3] R. R. S. Diat, "IIITM Face Emotion," *Kaggle*, 2023. [Online]. Available: <https://www.kaggle.com/datasets/rrsdiat/iiitm-face-emotion>
- [4] "Python 3.10.0 documentation," *Python.org*, 2023. [Online]. Available: <https://docs.python.org/3/>
- [5] "Matplotlib: Visualization with Python," *Matplotlib.org*, 2023. [Online]. Available: <https://matplotlib.org/stable/index.html>
- [6] "OpenCV-Python Tutorials," *OpenCV.org*, 2023. [Online]. Available: https://docs.opencv.org/4.x/d6/d00/tutorial_py_root.html
- [7] "Google Scholar," *Scholar.Google.ca*, 2023. [Online]. Available: <https://scholar.google.ca>
- [8] "Pyplot tutorial," *Matplotlib.org*, 2023. [Online]. Available: <https://matplotlib.org/stable/tutorials/pyplot.html>
- [9] "Loading Images – Python OpenCV Tutorial," *PythonProgramming.net*, 2023. [Online]. Available: <https://pythonprogramming.net/loading-images-python-opencv-tutorial/>
- [10] "Image Processing in Python with PIL/Pillow," *Auth0.com*, 2023. [Online]. Available: <https://auth0.com/blog/image-processing-in-python-with-pillow/#>
- [11] "Ways to Convert Image to Grayscale in Python using Skimage, Pillow and OpenCV," *MachineLearningKnowledge.ai*, 2023. [Online]. Available: <https://machinelearningknowledge.ai/ways-to-convert-image-to-grayscale-in-python-using-skimage-pillow-and-opencv/>
- [12] "Bar Plot in Matplotlib," *GeeksforGeeks.org*, 2023. [Online]. Available: <https://www.geeksforgeeks.org/bar-plot-in-matplotlib/>
- [13] "Max Pooling," *Deepai.org*, 2023. [Online]. Available: <https://deepai.org/machine-learning-glossary-and-terms/max-pooling>
- [14] "A Gentle Introduction to Pooling Layers for Convolutional Neural Networks," *machinelearningmastery.com*, 2023. [Online]. Available: <https://machinelearningmastery.com/pooling-layers-for-convolutional-neural-networks/>

[15] " Build Your Own Convolution Neural Network in 5 mins," *towardsdatascience.com*, 2023. [Online]. Available: <https://towardsdatascience.com/build-your-own-convolution-neural-network-in-5-mins-4217c2cf964f>

[16] " Precision, Recall, and F1 Score: A Practical Guide Using Scikit-Learn," *proclusacademy.com*, 2023. [Online]. Available: <https://proclusacademy.com/blog/practical/precision-recall-f1-score-sklearn/>

[17] "K Fold Cross Validation with Pytorch and sklearn"
<https://medium.com/dataseries/k-fold-cross-validation-with-pytorch-and-sklearn-d094aa00105f>

[18] "K-fold cross validation on custom dataset"
<https://discuss.pytorch.org/t/k-fold-cross-validation-on-custom-cataset/185808>

[19] "A Gentle Introduction to k-fold Cross-Validation"
<https://machinelearningmastery.com/k-fold-cross-validation/>