

Assignment #3

Speech Emotions Recognition

COLLABORATORS

i	Name	ID
1	Zeyad Zidan	19015709
2	Abdelrahman Gad	19015894
3	Abdelaziz Mohammed	19015941

You can view the ipynb file on [google colab](#) [here](#).

Problem Statement

Speech is the most natural way of expressing ourselves as humans. It is only natural then to extend this communication medium to computer applications. We define speech emotion recognition (SER) systems as methodologies that process and classify speech signals to detect embedded emotions. Below we are showing how we achieved the goal of the assignment.

Understanding the Dataset

The dataset is composed of 7442 wav-formatted audio files. Dealing with the dataset is not easy as audio files are not discrete signals. Relying on Python and sampling, we can represent the audio files as samples with a sufficient sampling rate. This method opens the way to formulate the audio files into an understandable dataset. In the following section, we are showing how we used [librosa](#) to formulate a dataset from the audio files.

Librosa

We used librosa to load the audio files and get the samples of each audio file. Adding each array of samples to an array, we successfully constructed a dataset of samples corresponding to each audio file.

Getting the labels is done by indexing the letters [9, 12] of the file name and then encoding each label to a numeric value.

Librosa also can display the waveform of each audio file.

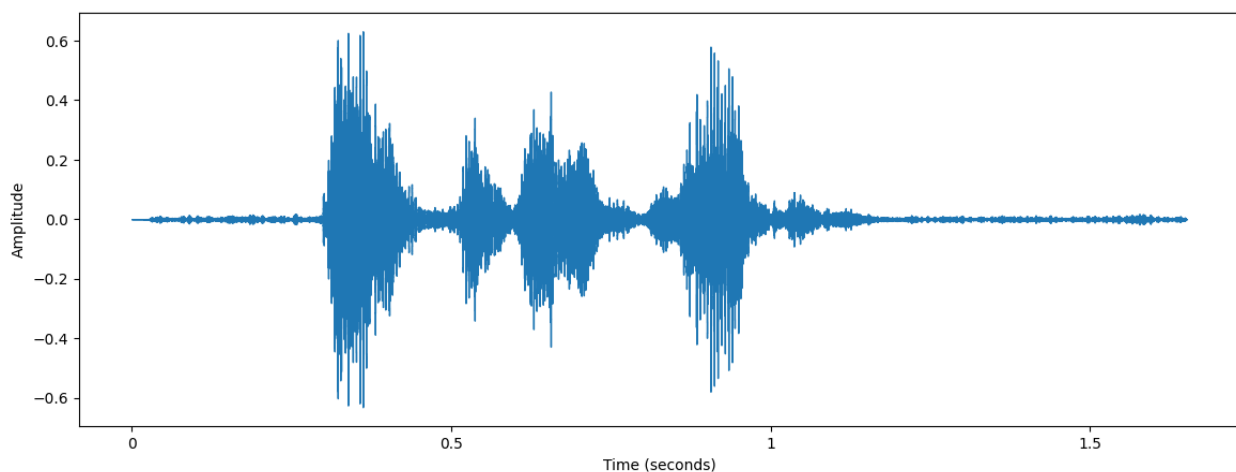


Figure 1. Sample waveform of an audio

Creating the Feature Space

Zero Crossing Rate Feature Space

The zero crossing rate is the rate of sign changes of the signal during the duration of a particular frame. To construct such feature space, librosa offers a built-in method in librosa.feature called zero_crossing_rate. Simply for each audio file, call this line of code:

```
librosa.feature.zero_crossing_rate(y=entry)[0]
```

Finally, append the previous result to the feature space list. Return the list after finishing each entry of the dataset.

Note that each entry in the dataset represents the samples array of each audio file.

Energy Feature Space

The energy of an audio file is the sum of squares of the signal values, normalized by the respective frame length. Librosa offers such an option in the same class as the zero crossing rate. To construct the energy feature space, Simply for each audio file, call this line of code:

```
frames = librosa.util.frame(entry,
frame_length=frame_size, hop_length=hop_length)

# Compute energy of each frame

energies = np.sum(frames**2, axis=0) / frame_size

feature_space.append(energies)
```

Finally, return the feature space array.

Mel Spectrogram

To construct the mel spectrogram feature space, we run the following code

```
1 def construct_melspectrogram(dataset, rate):
2     spectrogram = []
3     for i in range(len(dataset)):
4         spectrogram.append(librosa.feature.melspectrogram(y=dataset[i], sr=rate[i]))
5
6     return spectrogram
```

Building the Model

Data Augmentation

In the preparation step of the feature space, we take each entry and apply an augmentation algorithm to it. The augmentation algorithm we follow applies a number of **A** augmentations, e.g. **4** augmentations. We discuss the algorithm in the following points:

- For each augmentation, apply some random noise to a copy of the original feature space, i.e. we have a new unique feature space.
- The number of feature spaces created from the previous step will be a total of **4** newly identified random feature spaces based on the random noise addition by the end of all augmentations.
- In the same epoch (augmentation), we apply random shifting operations to a copy of the original feature space.
- The number of feature spaces created from the previous step also will be a total of **4** newly identified random feature spaces based on the random shifting operations by the end of all augmentations.
- In the same epoch, we apply random scaling operations to a copy of the original feature space. By the end of all epochs, we will have another **4** newly identified random feature spaces based on the random scaling operations.
- Finally, we concatenate all of the newly created feature spaces (**4 + 4 + 4 = 12**) to the original feature space to have a total of **((3A + 1) x N) x M**, where 3 stands for the number of applied operations (noise addition, shifting, and scaling), A stands for the number of epochs (augmentations), the 1 stands for the original feature space, N is the total number of entries in the original feature space, and M is the total number of features in the original feature space.

When to use this approach?

We use this approach after splitting the data into training and testing data to not mess with the original testing volume of data, that is we apply data augmentation on the **training data only** after split and exclude the testing data out of operation.

Best Model Results

We have tried multiple models, and so far this was the best model results

Layer (type)	Output Shape	Param #
conv1d_23 (Conv1D)	(None, 192, 1024)	6144
max_pooling1d_23 (MaxPooling1D)	(None, 94, 1024)	0
conv1d_24 (Conv1D)	(None, 90, 1024)	5243904
max_pooling1d_24 (MaxPooling1D)	(None, 43, 1024)	0
conv1d_25 (Conv1D)	(None, 39, 512)	2621952
max_pooling1d_25 (MaxPooling1D)	(None, 18, 512)	0
conv1d_26 (Conv1D)	(None, 14, 256)	655616
max_pooling1d_26 (MaxPooling1D)	(None, 5, 256)	0
flatten_8 (Flatten)	(None, 1280)	0
dense_17 (Dense)	(None, 512)	655872
dropout_2 (Dropout)	(None, 512)	0
dense_18 (Dense)	(None, 256)	131328
dropout_3 (Dropout)	(None, 256)	0
dense_19 (Dense)	(None, 6)	1542
Total params: 9,316,358 Trainable params: 9,316,358 Non-trainable params: 0		

Big Picture

In this section we show a comparison between a variety of models we have tried for each class of models. We discuss the 1D model first, breaking it down into 5 different models with different accuracies and confusing classes, then we discuss the 2D models with a variety of 3 different models.

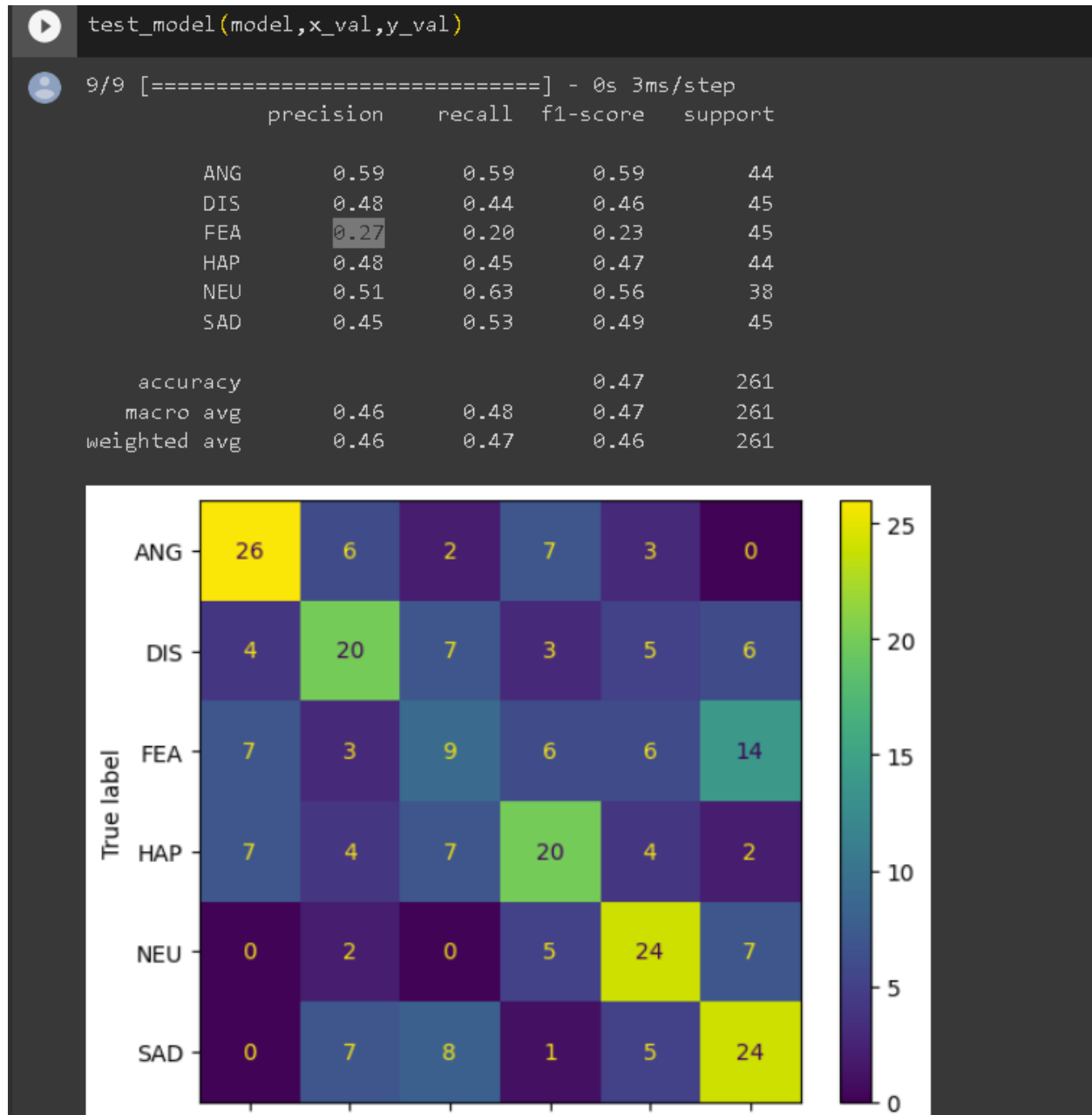
1D Models

Model #	Most Confusing Class	Most Confusing Class Precision	Accuracy
1	FEAR	0.27	0.47
2	FEAR	0.40	0.49
3	NEUTRAL	0.40	0.49
4	NEUTRAL	0.34	0.47
5	FEAR	0.33	0.44

Best accuracy → 0.49 in model #2 and model #3.

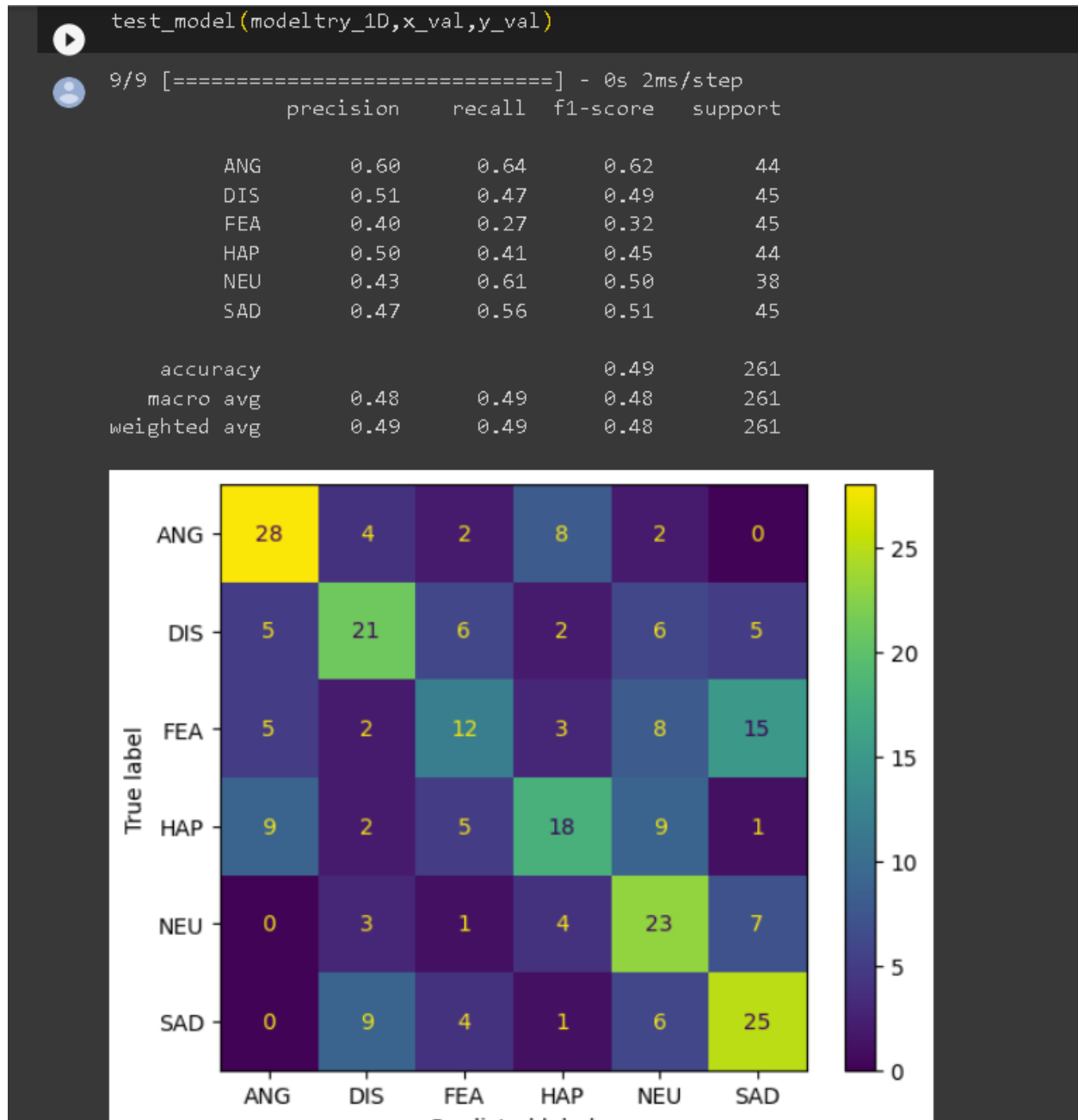
First Model

In this model, most confusing class is **FEAR** as its precision equals **0.27**. The accuracy of this model is **0.47**.



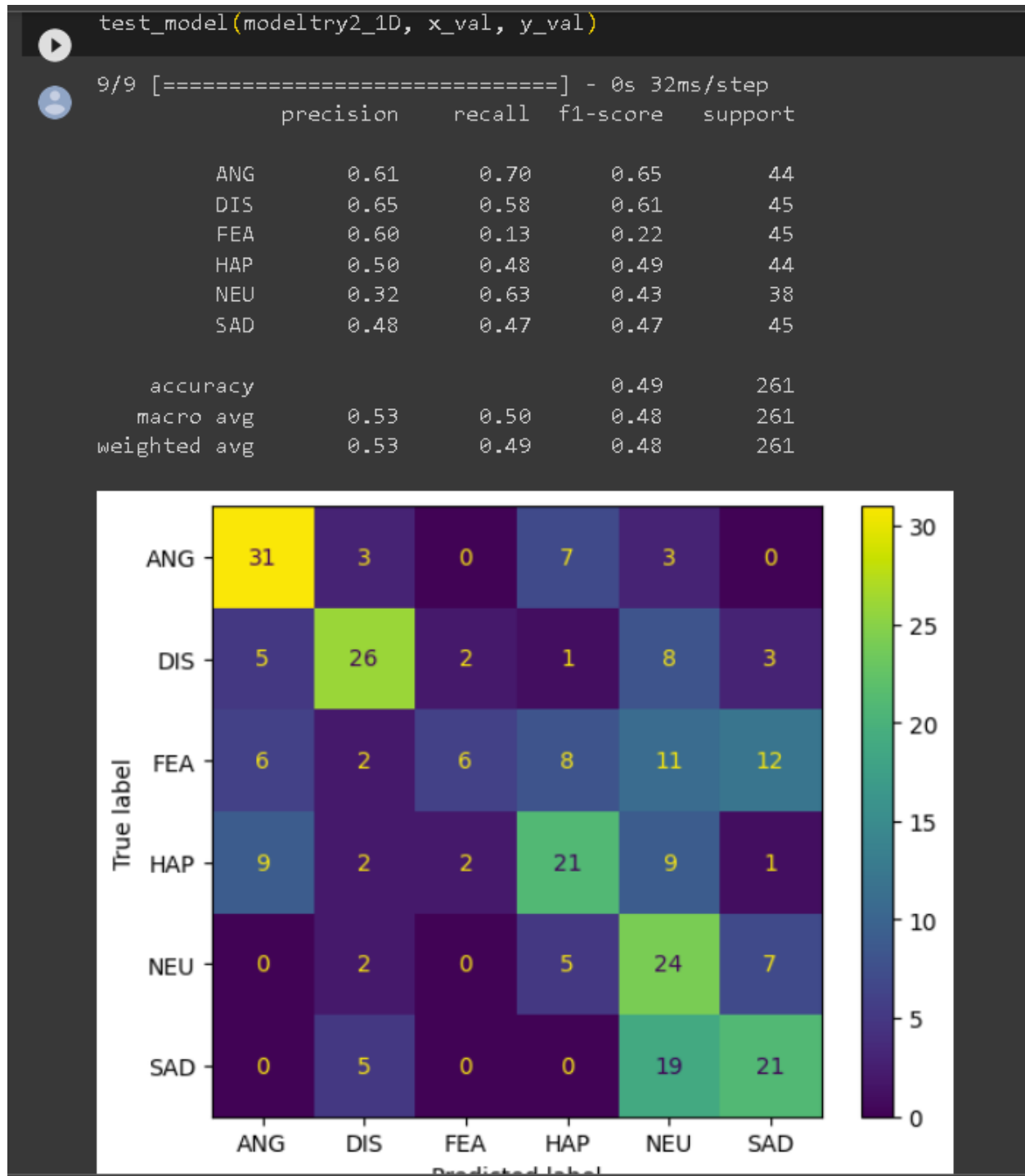
Second Model

In this model, most confusing class is **FEAR** as its precision equals **0.40**. The accuracy of this model is **0.49**.



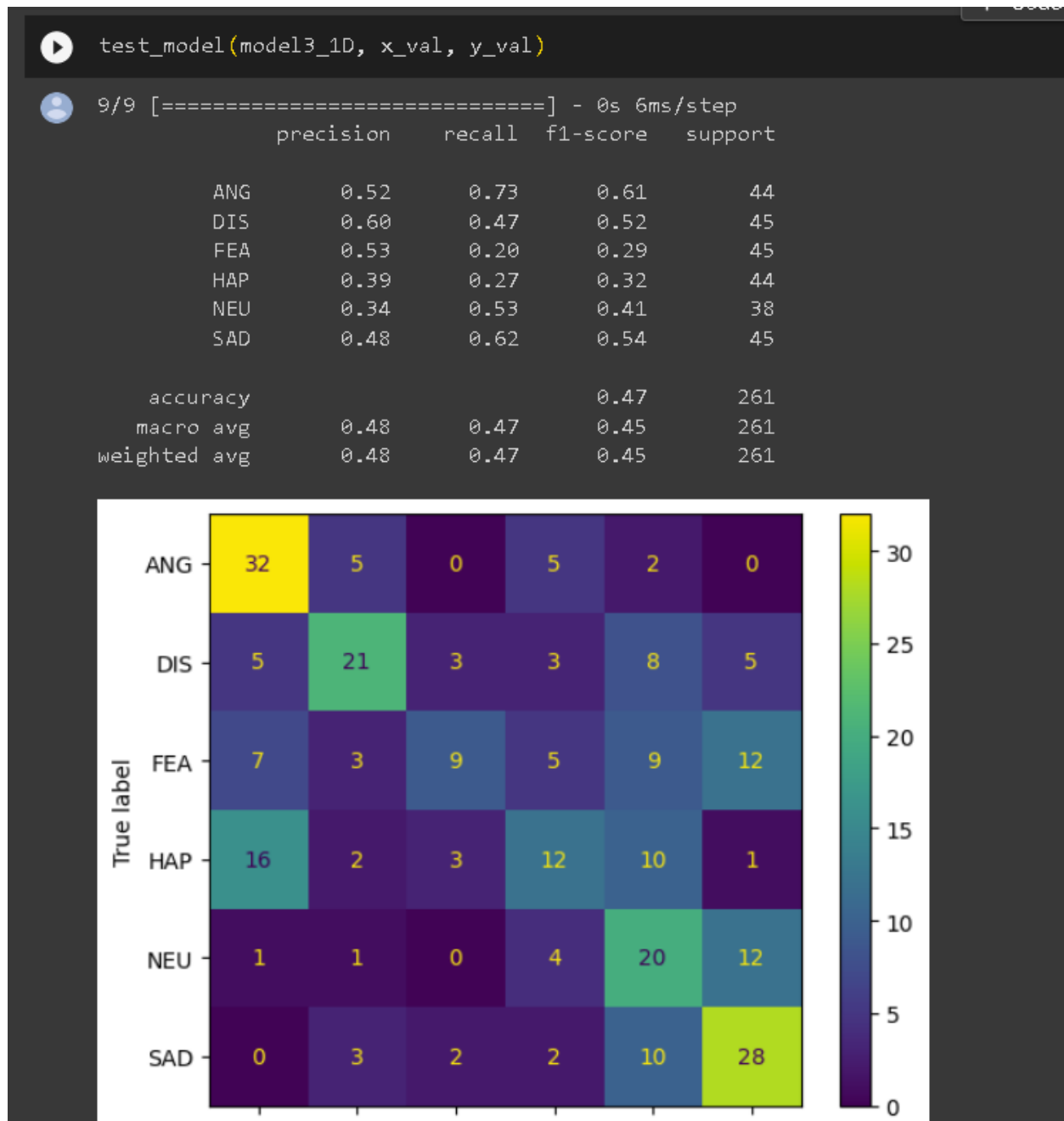
Third Model

In this model, most confusing class is **NEUTRAL** as its precision equals **0.40**. The accuracy of this model is **0.49**.



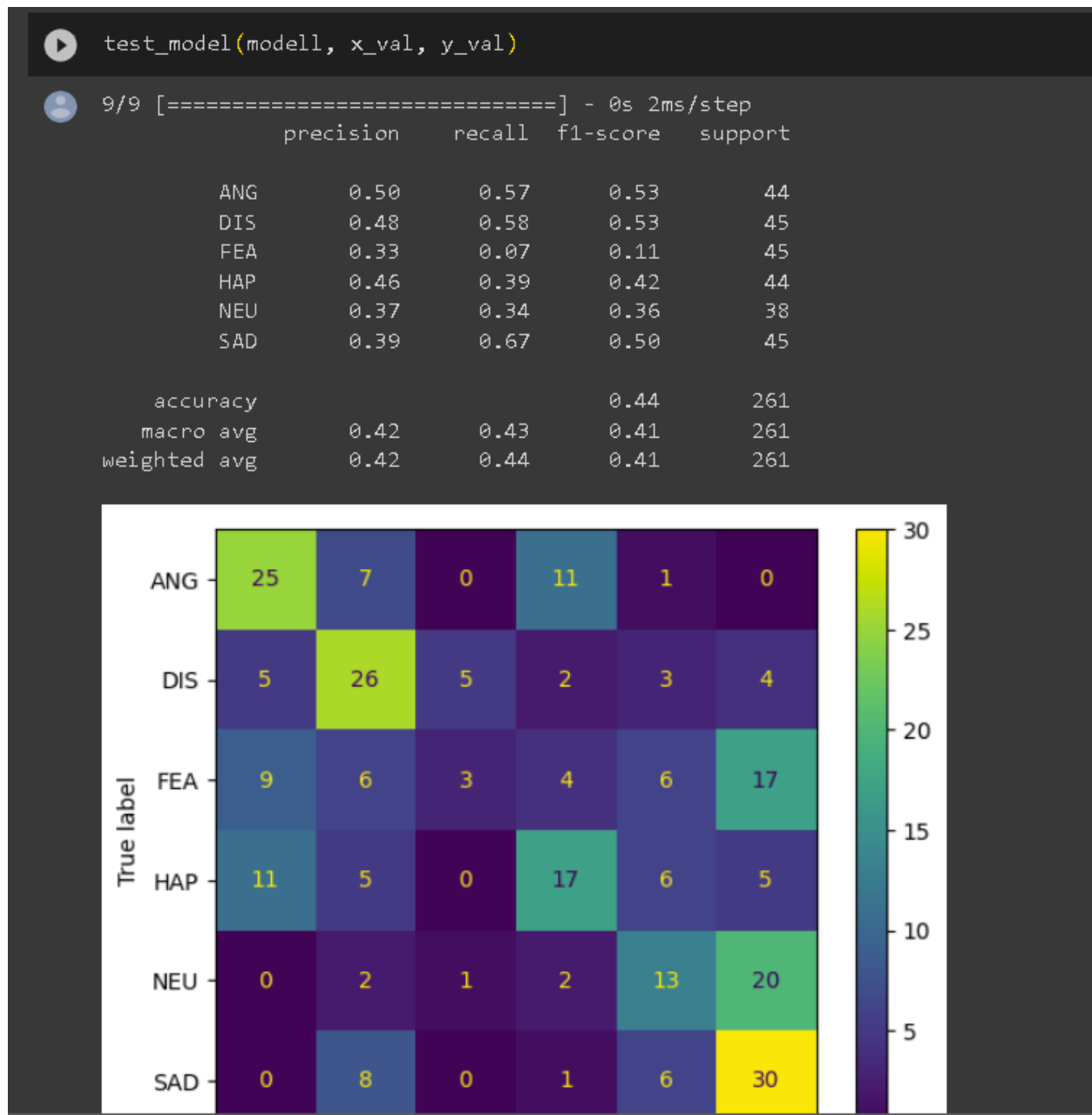
Fourth Model

In this model, the most confusing class is **NEUTRAL** as its precision equals **0.34**. The accuracy of this model is **0.47**.



Fifth Model

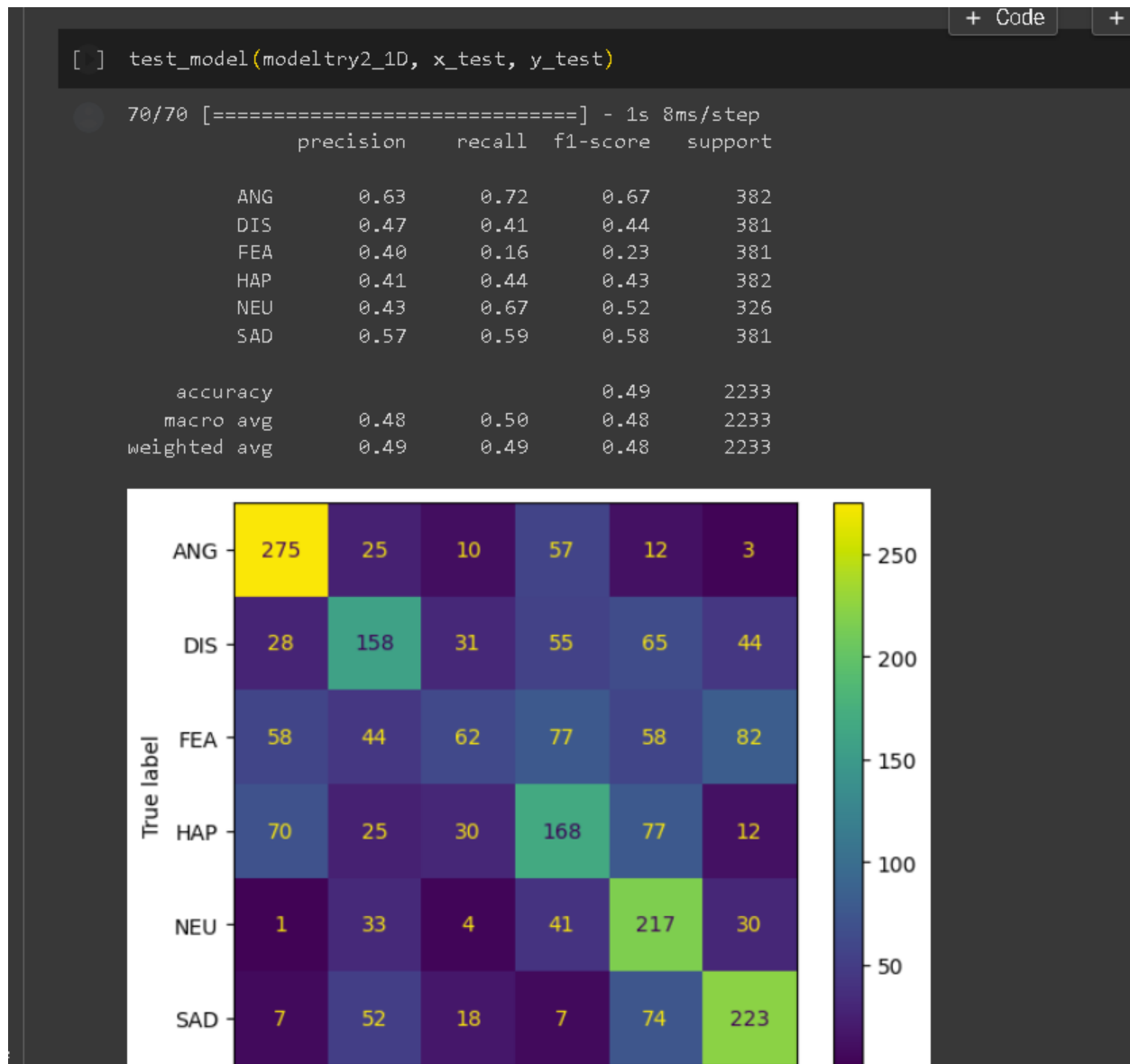
In this model, the most confusing class is **FEAR** as its precision is **0.33** with an accuracy of **0.44**.



Conclusion

So it appears that the model (**modeltry2_1D**) has the best accuracy in the validation so we apply this model on testing.

Accuracy: **0.49**.



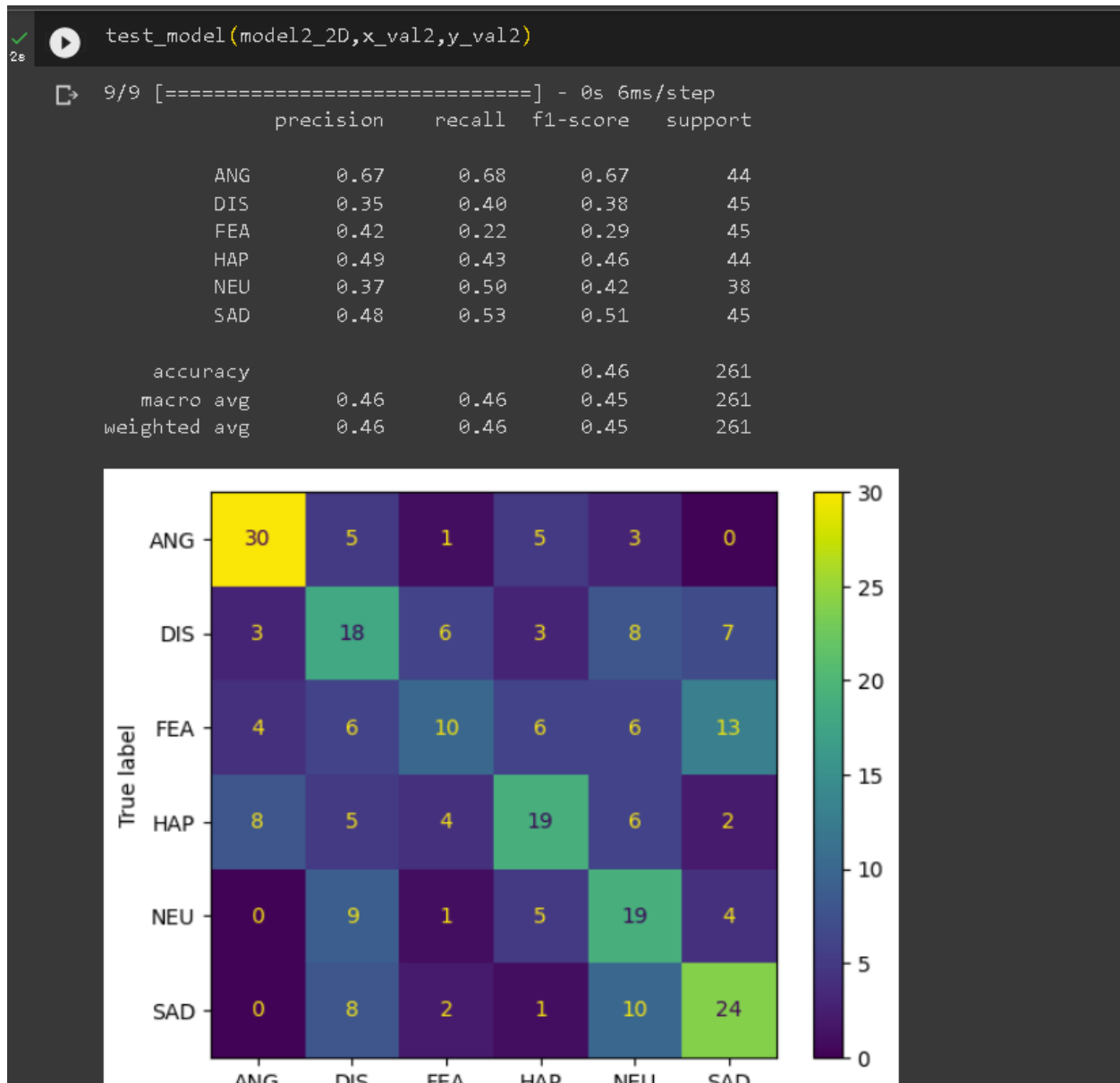
2D Models

Model #	Most Confusing Class	Most Confusing Class Precision	Accuracy
1	DISSECTED	0.35	0.46
2	DISSECTED	0.44	0.51
3	ANGRY	0.17	0.17

Best accuracy → 0.51 in model #2.

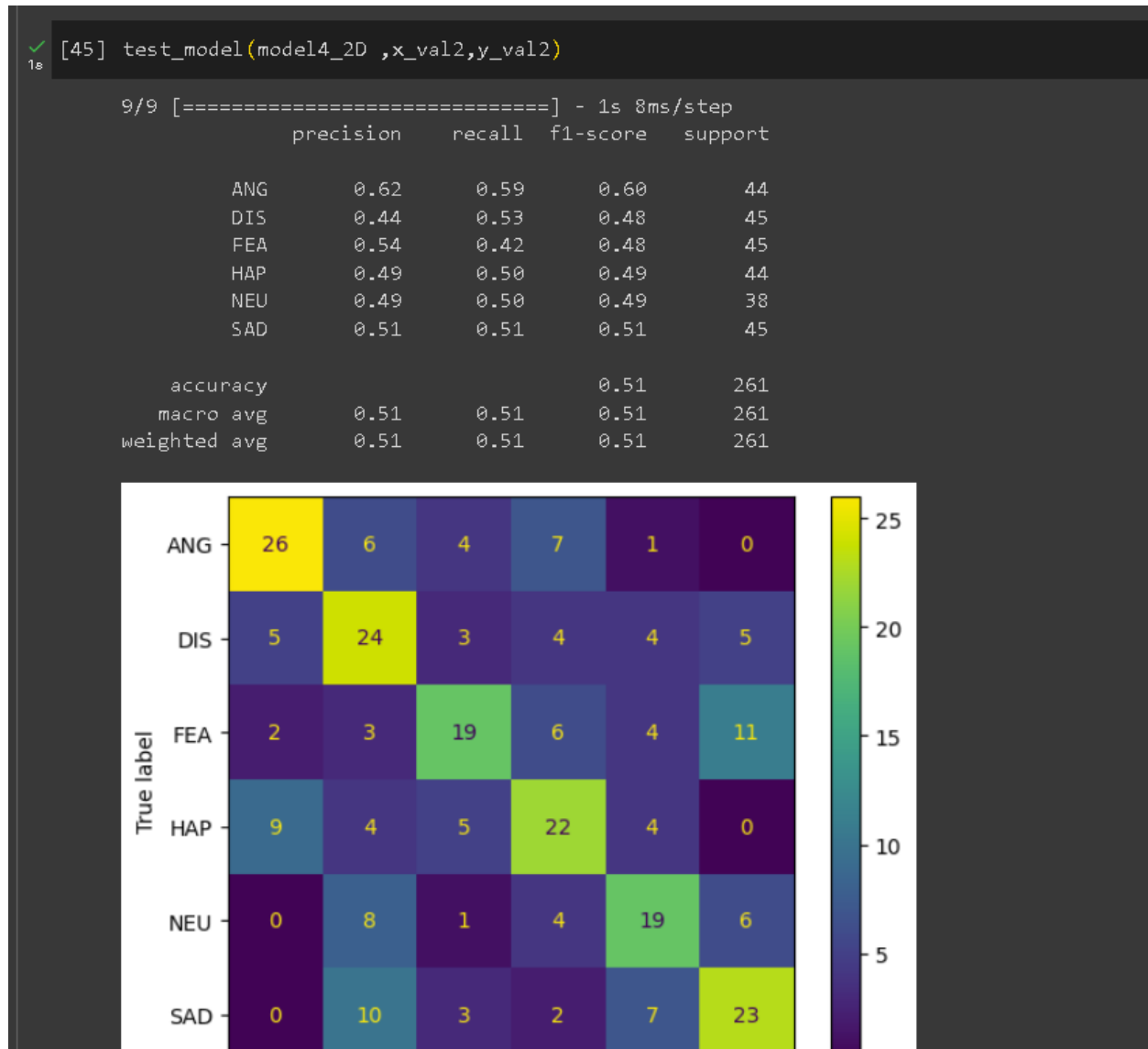
First Model

This model get accuracy **0.46** and the most confusing classes here is **DISSECTED** as it appears in the matrix.



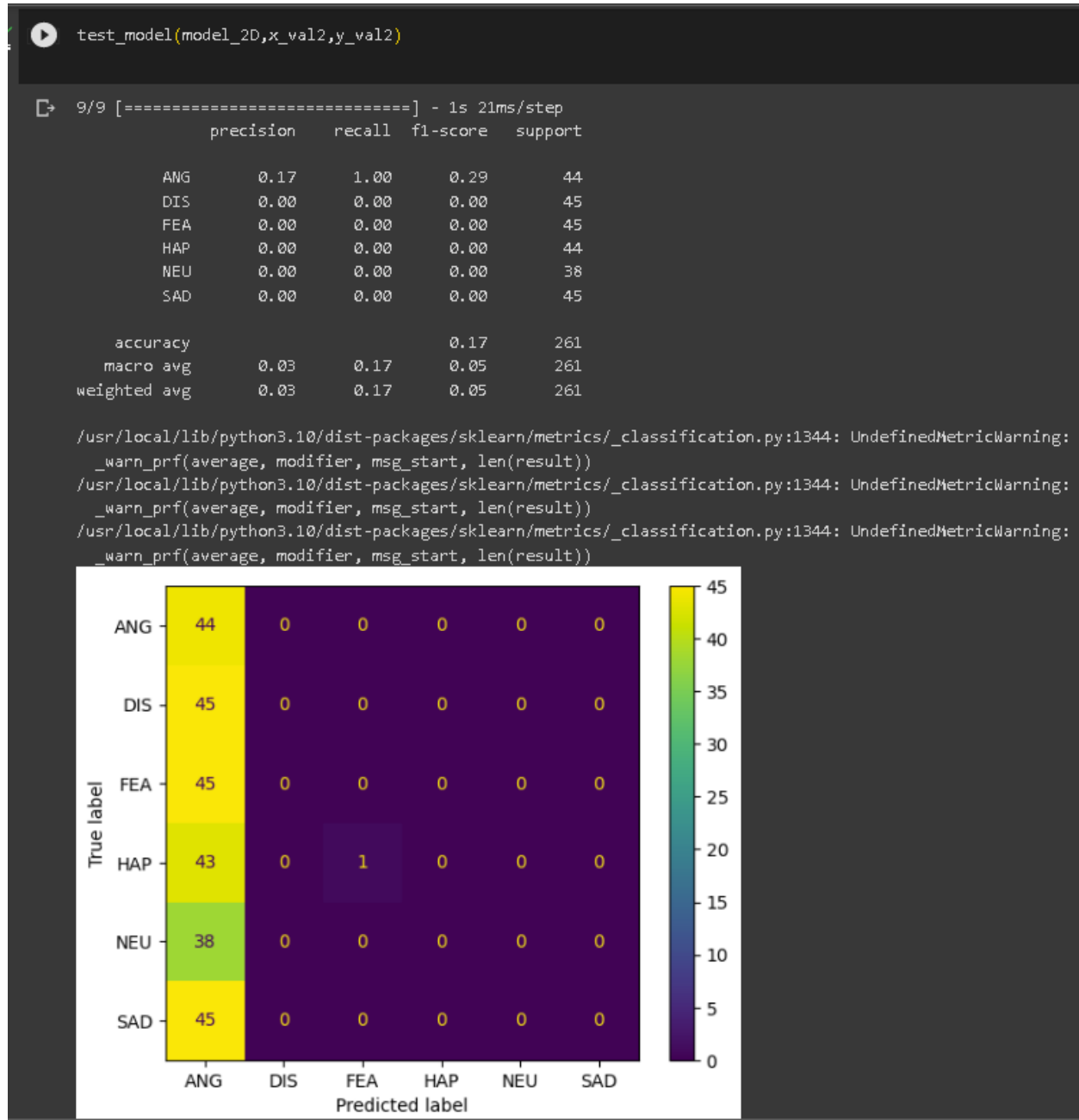
Second Model

This model4_2D accuracy is **0.51 (2D)**, and the most confusing class is **DISSECTED** as it appears in the matrix.



Third Model

This model accuracy is **0.17** in which it appear it detect all sounds as an angry voice.



Conclusion

So it appears that the best one for the 2d validation is model4_2D with an accuracy that equals **0.51**, So we apply testing to it and it appears that the accuracy is **0.49** and the most confusing matrix is **DIS**.

