# Hand Gesture Classification with Convolution Neural Networks (CNN)

NATHANIEL SOLON PEDERSEN, University of Stavanger, Norway

Hand gestures plays an important role in non-verbal communication between humans. With machine learning we can predict and classify which hand gesture someone makes. In this study we will be experimenting with different CNNs by tuning different parameters. The dataset we are using consists of labeled Grey-scale images [1]. Source code can be found on Github[3]

## 1 Data preprocessing

The raw dataset consists of Grey-Scale images captured under pretty good conditions, but to improve the model performance, we will need to process the images further.



Fig. 1. Image[1] of an 'l' gesture before processing

### 1.1 Shrinking

Resizing all images into 32x32. This will ensure consistency between all images input into the model. This is also a great way to reduce the amount of data and speed up training time.

### 1.2 Grey-Scale

We will reduce the dimensionality of the input data from 3 channels (RGB) to a single channel. The colors are not that relevant for hand gesturing and will only add a big computational overhead.

Author's Contact Information: Nathaniel Solon Pedersen, University of Stavanger, Stavanger, Rogaland, Norway.

### 1.3 Blur

A Gaussian blur will be used to remove out noise from the images to improve model accuracy.

### 1.4 Edge Detection

We are using Canny edge detection to highlight edges. This will just make the gesture generally easier to identify.

### 1.5 Normalize

Normalization is redundant when using the canny edge detection, but the job of normalization would be to mitigate gradients exploding or vanishing.
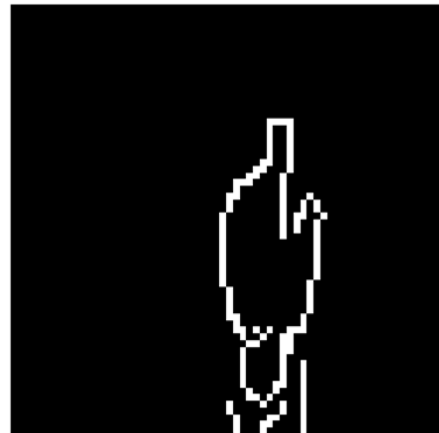


Fig. 2. Image[1] of an 'l' gesture after processing

## 2 Feature Extraction and Splitting

We will convert all images into numpy arrays, so we can use that as input for tensorflow keras. We will also one hot encode all labels. All data is splitted into training data, validation data, and testing data. The split is 80/10/10. Each data-split will have an equal amount of variety between each label.

## 3 Models

We will be using 2 different CNN models. Below is the pseudocode for the first model.

```
Input(shape=(64, 64, 1))


Conv2D(32, (3, 3), activation="relu")
MaxPooling2D(pool_size=(2, 2))


Conv2D(64, (3, 3), activation="relu")
MaxPooling2D(pool_size=(2, 2))
```

```
Conv2D(64, (3, 3), activation="relu")
MaxPooling2D((2, 2))

Flatten()
Dense(128, activation="relu")
Dense(10, activation="softmax")

compile(loss="MAE")
```

Below is the pseudocode for the second model

```
Input(shape=(64, 64, 1))

Conv2D(32, (3, 3), activation="relu")
MaxPooling2D(pool_size=(2, 2))

Flatten()
Dense(128, activation="relu")
Dense(10, activation="softmax")

compile(loss="categorical_crossentropy")
```
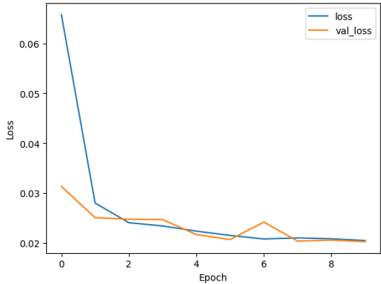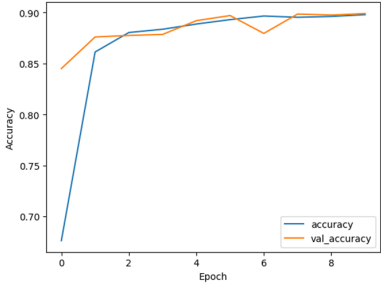
## 4 Results

When using the confusion matrix, refer to this table.

| Number | Label |
|--------|-------|
| 1 | palm |
| 2 | l |
| 3 | fist |
| 4 | fist moved |
| 5 | thumb |
| 6 | index |
| 7 | ok |
| 8 | palm moved |
| 9 | c |
| 10 | down |

Table 1. Mapping of Numbers to Hand Gesture Labels

### 4.1 Experiment 1

We are using the first model, and training it with an Epoch of 10.









```
Test Loss:  0.020285068079829216
Test Accuracy:  0.8985000252723694
Test Precision:  0.8985000252723694
Test Recall:  0.8985000252723694
Test F1 Score:  0.8985000252723694
```
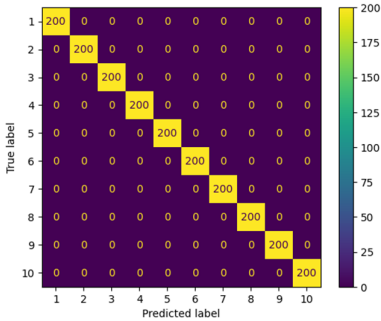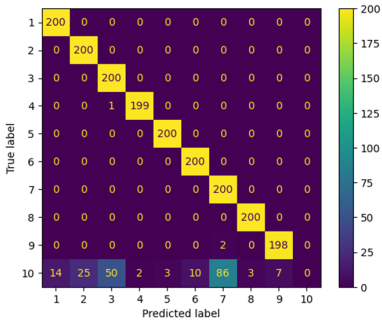
We see that the model is performing with minimal losses, though it struggles with identifying the down gesture.
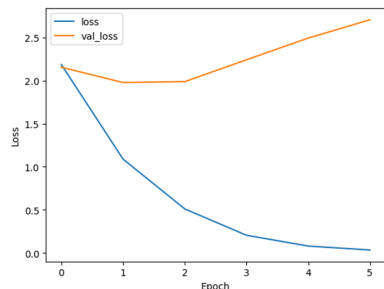
### 4.2 Experiment 2
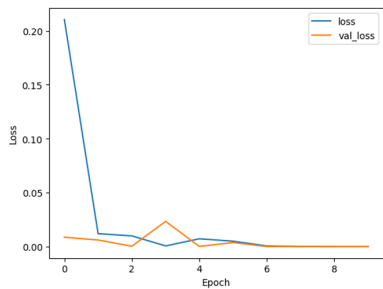
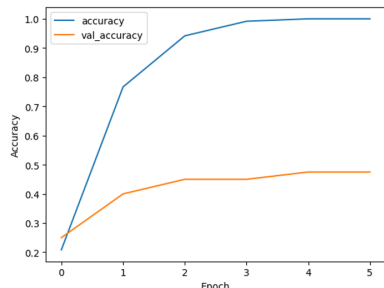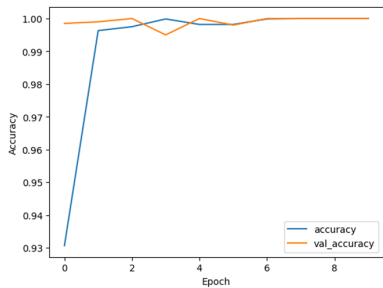We swapped out from MAE to categorical crossentropy

```
Test Loss:   4.887561317445943e-06
Test Accuracy:  1.0
Test Precision:  1.0
Test Recall:  1.0
Test F1 Score:  1.0
```

```
Test Loss:  2.219278573989868
Test Accuracy:  0.515500009059906
Test Precision:  0.611940324306488
Test Recall:  0.4715000092983246
Test F1 Score:  0.532617920251381
```
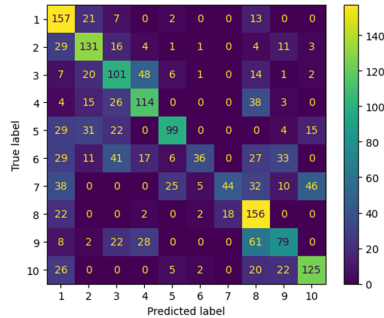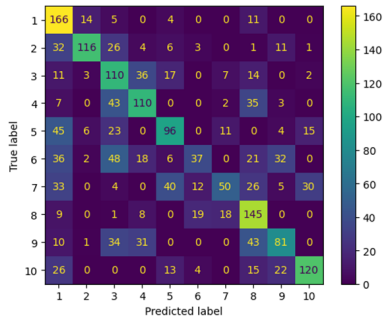
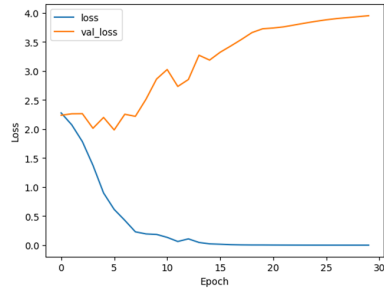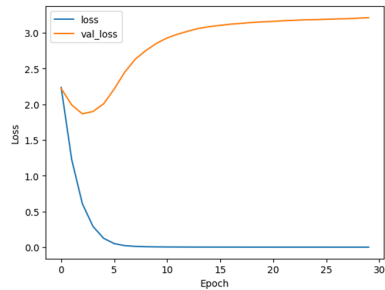The model correctly predicted all of our tests.

### 4.3 Experiment 3

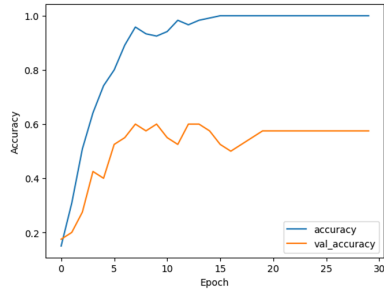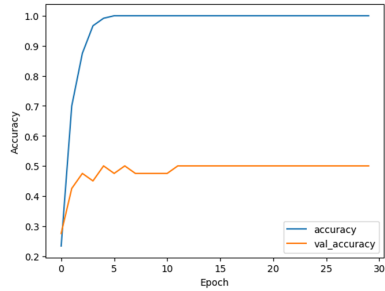We swapped out to a smaller model. The Advantage of using a smaller model is that it is less computational expensive, but at the sacrifice of it making worse predictions. We also reduced the Epoch to 6, reduced training data down to 100 in length and valuation data down to 40 in length.

We see that the current model runs worse than our original model.

### 4.4 Experiment 4

We increased the Epoch to 30.

```
Test Loss:  3.109323501586914
Test Accuracy:  0.6489999890327454
Test Precision:  0.6561224460601807
Test Recall:  0.6430000066757202
Test F1 Score:  0.6494949514702161
```
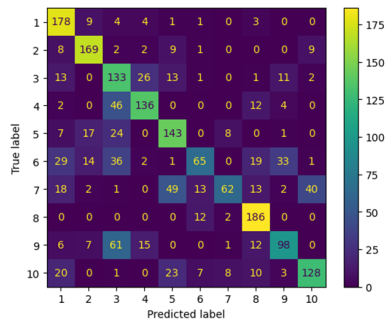
```
Test Loss:  2.7096049785614014
Test Accuracy:  0.5210000276565552
Test Precision:  0.5653631091117859
Test Recall:  0.5059999823570251
Test F1 Score:  0.5340369208419844
```

We can observe over fitting, and slight gradient explosions in accuracy when valuating.

### 4.5 Experiment 5

We went back to the previous model.



We see the same issue with over fitting again, so it means we are lacking data. We also see exploding gradients.

## 5 Limitations

### 5.1 Augmentation

When training, we do not take account into rotation nor mirroring, which would be a smart thing to do to have a more robust model.

### 5.2 Improved datasets

The dataset we are currently using is sampled in infrared and has ideal conditions. In the future, we should use samples with more noise, so we get a more robust model.

## 6 Acknowledgments

### Acknowledgments

To Leap Motion for providing the datasets for this study[1][2]

### References

[1] Leap Motion 2016. *Hand Gesture Recognition Database*. Leap Motion, Vienna, Austria. Retrieved November 25, 2024 from https://www.kaggle.com/datasets/gti-upm/leapgestrecog

[2] 2016. *Hand Gesture Recognition using Infrared Imagery Provided by Leap Motion Controller*. https://doi.org/10.1007/978-3-319-48680-2_5

[3] nasope. 2024. *dat305-2024*. github. https://github.com/nasope/dat305-2024