

Virtual Reality handwritten character recognition using a Convolutional Neural Network

Problem statement: Recognize characters written in a virtual reality environment using a CNN.

VR project at: <https://github.com/MadalinaMircea/VRCharacterRecognition>

CNN project at: <https://github.com/MadalinaMircea/MLCharacterRecognition>

Flow description:

The virtual reality environment is developed in Unity using C#.

The CNN is developed in C#.

In VR, the user is presented with a whiteboard and a pen. They can pick up the pen and write a letter on the whiteboard. When the user has not written anything for 2 seconds, the whiteboard is cleared and the letter is sent to the CNN for recognition.

The Unity project and the CNN communicate through a folder called “UnityCNN” created on the computer desktop. Unity saves the texture from the whiteboard to the file “image.jpg”. The CNN continuously waits for such a file to exist. When it exists, the CNN reads it and analyses the image, deletes this file and creates a file called “response.txt” where the recognized character is written. The Unity project then reads this file and displays the response in a 3D text object in the VR environment.

The Unity VR project:

Developed for Oculus Rift S using SteamVR.

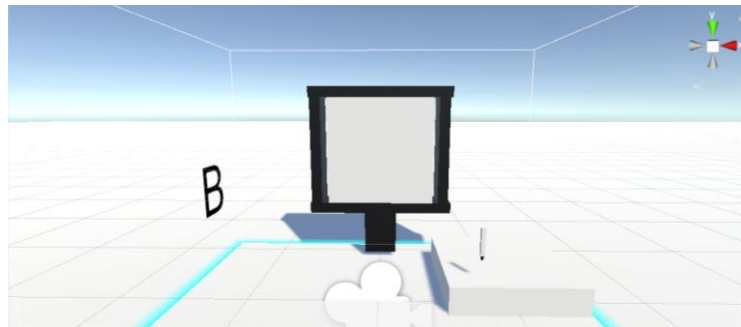
The “Hand” script is added to both controllers. This script makes the controllers able to grab an object.

The “Interactable” script is added to the pen. This allows the pen to be grabbed.

The “Whiteboard” script is attached to the whiteboard object. This controls the writing on the texture of the object by setting the pixels around a collision point to a certain color. This also counts the seconds from the last time the pen touched the board to see if it is time to send the texture to the CNN.

The “WhiteboardMarker” script creates a raycast from the tip of the pen upwards to check if the pen is close to an object tagged with the “Whiteboard” tag. If so, the pen calls methods in the “Whiteboard” script of this object and transmits the position of the collision.

The “CNNController” script controls when the texture image is sent to the CNN and when the response character is displayed.



Screen capture from the Unity VR project

The Convolutional neural network:

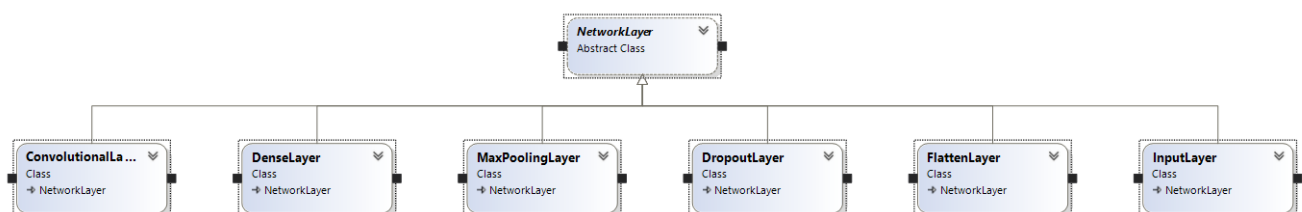
Developed from scratch (the Keras model is used for demonstrative purposes).

The CNNController file creates a ConvolutionalNeuralNetwork model and controls the training, testing and overall usage of this model.

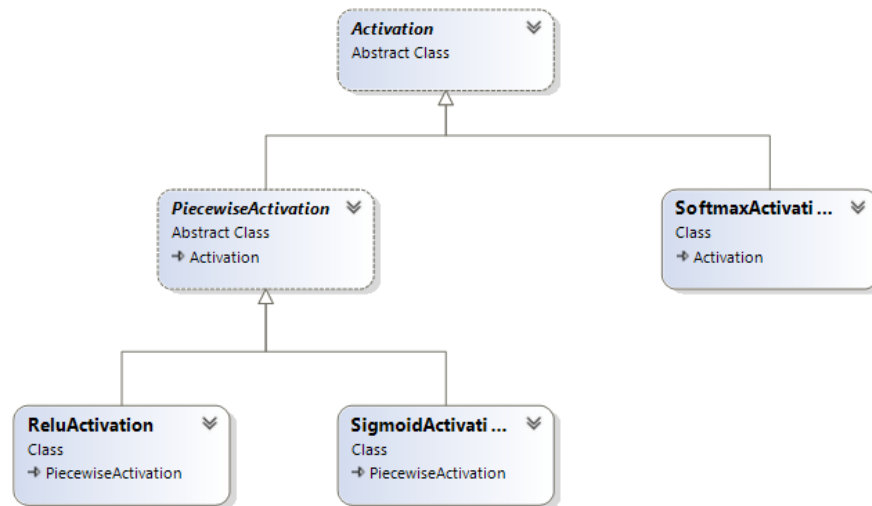
The model architecture is written to file as a json. The weights are written in a folder, with separate folders for each layer containing weights (convolutional and dense). Each layer folder contains separate files for each filter/unit. The filters and units are written to file as json. The same file-folder structure is used when reading a model.

The ImageProcessing folder contains the files needed for processing images. The ImageProcessing file contains methods that convert a JPG image to a matrix of numbers, normalize the matrix and so on. The LayerOutput class is a parent class for the FilteredImage and FlattenedImage classes. FilteredImage is an image with multiple channels, each channel being represented by a matrix. FlattenedImage is an image that has been flattened to one array. These classes are used as output in the layers of the network. The ImageRepository and ImageController classes are used to control the flow of the dataset (training, testing and validation) from the folders.

The CNN folder contains all the needed files. The Layers folder contains the layer classes (Convolutional, Dense, Flatten, Dropout, Maxpooling), as well as the units needed in these layers (Filter, Kernel, Unit). All of the layers are children of the NetworkLayer class.



The Activations folder contains all of the activation files. The Activation class is the parent class. Softmax is a child of this class. PiecewiseActivation is also a child of the Activation class, but this is in itself a parent of all the other activations (Relu, Sigmoid). This parent class is needed for the activations that are applied element by element on a layer, not on the whole layer, like Softmax. The parsing of the elements is done in this parent class, while the element activation is overridden in the children classes.



The Utils folder contains helper classes (JsonHelper), as well as classes with mathematical computations used throughout the algorithm (Convolution, Matrix rotation).

Creating a model:

Before creating a model, the `GlobalRandom.InitializeRandom()` static method needs to be called. This is a static class used to generate all of the random numbers needed in a model. A global random was necessary because this static class will only request one seed for the random number generator, thus insuring that the kernels/weights will not be equal to each other (if each kernel has its own Random variable, there is a chance the seed will be equal among kernels and the kernels will end up equal to each other).

To create a model, a `ConvolutionalNeuralNetwork` object is created. Then, the “Add” method can be used to add any type of layer to it, since this method receives a “`NetworkLayer`” parameter, which is a parent of all the layer classes.

Before using the model, the `Compile` method needs to be called. This calls the `Compile` method of each layer with the previous layer as parameter, creating a link between layers and ensuring that the kernel/unit number is correct for all layers (the kernel number of a convolution layer should be equal to the filter number of the previous layer; the output matrix size of a

maxpooling layer will be equal to the size of the previous layer matrix divided by the pool size; etc).

Training a model:

To train a model, a list of InputOutputPair is given to the Train method. The InputOutputPair contains the path of a file, the character this image file represents, as well as the one-hot-encoded version of this character. This list has already been shuffled. The Train method parses the list, takes the image file and calls ImageProcessing to request a normalized FilteredImage version, the feeds this image to the network. The output is taken, the error between the output and the expected output is computed, then this error is passed to the Backpropagate method of each layer, starting with the last one.

Datasets used:

- Chars74K
- NIST
- original characters generated in the Windows Forms application

Dataset split:

- Training: $40 * \text{Chars74K/letter} + 1160 * \text{NIST / letter} (+ 30 * \text{original / letter}) \Rightarrow 1230 \text{ images/letter} \Rightarrow 31980 \text{ images}$
- Testing: $10 * \text{Chars74K/letter} + 150 * \text{NIST / letter} (+ 10 * \text{original / letter}) \Rightarrow 170 \text{ images/letter} \Rightarrow 4420 \text{ images}$
- Validation: $5 * \text{Chars74K/letter} + 190 * \text{NIST / letter} (+ 10 * \text{original / letter}) \Rightarrow 205 \text{ images/letter} \Rightarrow 5330 \text{ images}$

Convolution operation:

For the zero-padded version, the matrix was not padded, but if statements were used to ensure that the indices were inside the bounds. This saves the computation time that would have been required by the padding.

The outer for-loops parse the output matrix. The inner for-loops parse the kernel. For the version with “same” padding, the following multiplications are summed:

$\text{matrix}[i + \text{kernelI} - \text{halfSize}, j + \text{kernelJ} - \text{halfSize}] * \text{kernel}[\text{kernelI}, \text{kernelJ}];$

For the version with “valid” padding, the following multiplications are summed:

$\text{matrix}[i + \text{kernelI}, j + \text{kernelJ}] * \text{kernel}[\text{kernelI}, \text{kernelJ}];$

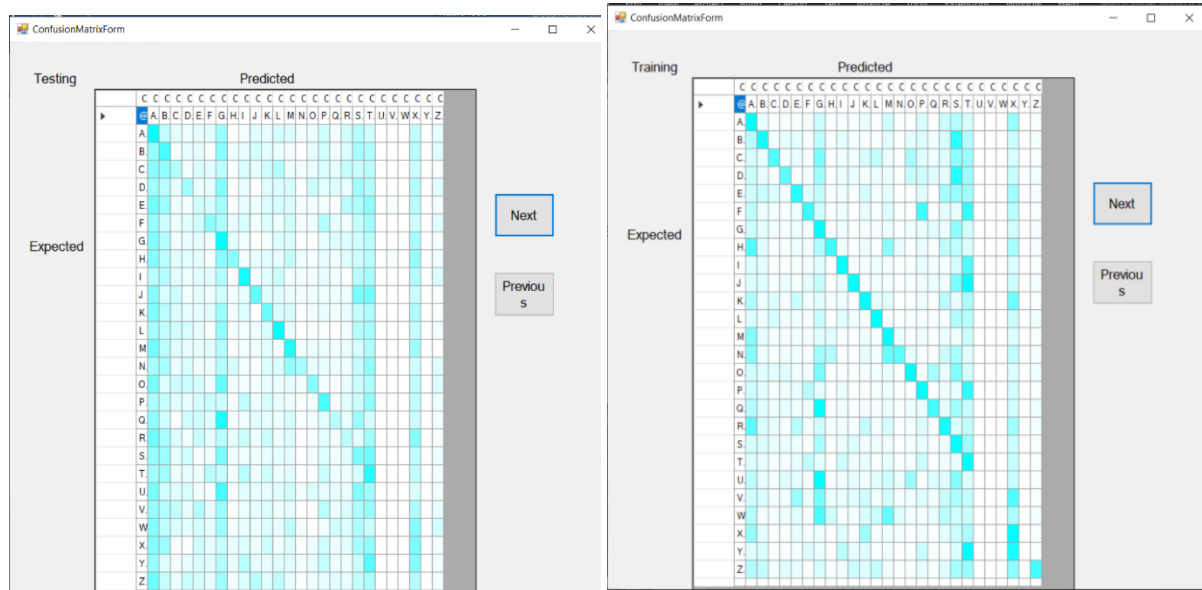
Performance:

Training for one image takes around 4 seconds.

Recognising one image takes around 1 second.

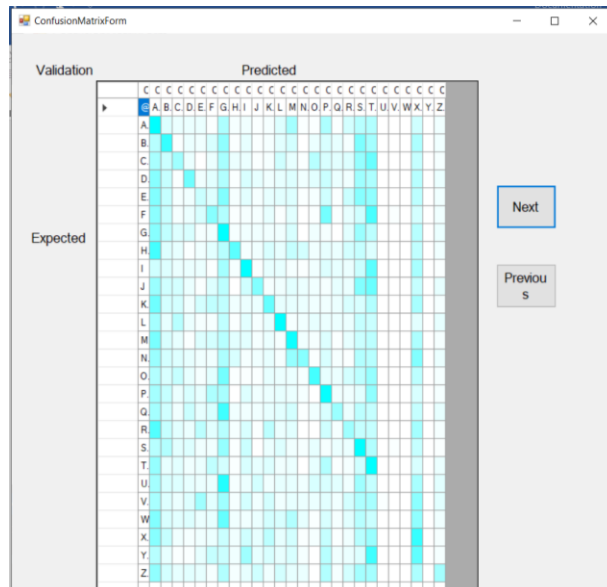
These statistics were obtained after introducing parallelism in the computations. The C# Task class was used to parallelize computations. Each filter is computed in parallel, and each kernel from one filter is also computed in parallel. Also, for all of the 4-piece for-loops, the first for loop created tasks for each of the steps, then these tasks were waited for outside of the for loops.

Confusion matrix



Confusion matrix for the testing set.

Confusion matrix for the training set.



Confusion matrix for the validation set

Performance JSONs

```
TrainingPerformance - Notepad
File Edit Format View Help
[59,38,10,20,27,87,33,3,24,5,6,9,21,6,7,494,8,19,28,264,0,0,0,56,0,6],
[92,38,19,29,15,16,365,8,14,9,11,9,23,5,22,26,238,29,119,68,0,0,0,56,0,19],
[205,45,28,7,88,17,58,10,14,12,71,15,48,3,11,45,9,291,47,76,0,0,0,116,0,14],
[51,46,28,26,22,4,105,4,25,42,16,21,20,1,14,24,12,18,558,121,0,0,0,63,0,9],
[39,27,12,11,27,76,31,2,60,20,29,9,29,7,12,63,6,16,27,652,0,0,0,70,0,5],
[81,30,15,40,32,13,331,11,20,32,15,36,58,23,125,9,46,39,158,39,0,0,0,70,0,7],
[57,46,18,15,125,23,173,23,29,38,59,16,55,13,18,43,13,66,100,74,0,0,0,218,0,8],
[89,17,12,20,38,9,227,60,13,14,46,69,198,37,14,13,28,87,62,61,0,0,0,102,0,14],
[131,35,4,5,29,10,23,3,40,12,121,12,22,11,5,17,5,31,65,102,0,0,0,539,0,8],
[39,44,5,7,12,40,26,2,83,31,42,3,19,2,5,66,8,9,45,450,0,0,0,287,0,5],
[96,75,18,38,29,7,51,9,14,65,41,75,16,4,8,28,9,36,113,123,0,0,0,141,0,234]],
"TP": [582,298,217,201,268,246,636,246,637,314,418,605,656,203,374,494,238,291,558,652,0,0,0,539,0,234],
"TN":
[30102,29818,29737,29721,29788,29766,30156,29766,30157,29834,29938,30125,30176,29723,29894,30014,29758,29811,30078,30172,29520,29
520,29520,30059,29520,29754],
"FN": [648,932,1013,1029,962,984,594,984,593,916,812,625,574,1027,856,736,992,939,672,578,1230,1230,1230,691,1230,996],
"FP": [648,932,1013,1029,962,984,594,984,593,916,812,625,574,1027,856,736,992,939,672,578,1230,1230,1230,691,1230,996],
"Sensitivity":
[0.47317073170731705,0.24227642276422764,0.17642276422764228,0.16341463414634147,0.21788617886178863,0.2,0.51707317073170733,0.2,
0.51788617886178867,0.25528455284552848,0.33983739837398375,0.491869918699187,0.5333333333333333,0.16504065040650406,0.304065040
65040652,0.40162601626016259,0.19349593495934958,0.23658536585365852,0.45365853658536587,0.53008130081300808,0.0,0.0,0.0,0.438211
38211382116,0.0,0.19024390243902439],
"Specificity":
[0.97892682926829266,0.96969105691056912,0.9670569105691057,0.96653658536585363,0.96871544715447155,0.968,0.98068292682926828,0.9
68,0.98071544715447156,0.97021138211382119,0.97359349593495936,0.97967479674796742,0.98133333333333328,0.96660162601626021,0.9721
6260162601631,0.97606504065040656,0.967739837398374,0.96946341463414631,0.97814634146341461,0.98120325203252035,0.96,0.96,0.96,0.
9775284552845529,0.96,0.967609756097561],
"ClasswiseAccuracy":
[0.95947467166979361,0.9417135709818637,0.93664790494058792,0.9356472795497186,0.93983739837398372,0.93846153846153846,0.96285178
23639775,0.93846153846153846,0.96291432145090683,0.9427141963727329,0.94921826141338339,0.9609130706691682,0.96410256410256412,0.
93577235772357725,0.94646654158849286,0.95397123202001255,0.93796122576610386,0.94127579737335831,0.95797373358348969,0.963852407
75484682,0.92307692307692313,0.92307692307692313,0.92307692307692313,0.95678549093183241,0.92307692307692313,0.93771106941838644]
,"ClasswiseMissclassification":
[0.040525328330206382,0.058286429018136333,0.063352095059412139,0.064352720450281425,0.060162601626016263,0.061538461538461542,0.
037148217636022517,0.061538461538461542,0.037085678549093185,0.05728580362726704,0.050781738586616637,0.039086929330831771,0.0358
97435897435895,0.064227642276422761,0.053533458411507195,0.046028767979987492,0.062038774233896185,0.058724202626641651,0.0420262
66416510319,0.036147592245153223,0.076923076923076927,0.076923076923076927,0.076923076923076927,0.076923076923076927,0.0432145090681676,0.076923076923
076927,0.062288930581613507],
"ClasswisePrecision":
[0.47317073170731705,0.24227642276422764,0.17642276422764228,0.16341463414634147,0.21788617886178863,0.2,0.51707317073170733,0.2,
0.51788617886178867,0.25528455284552848,0.33983739837398375,0.491869918699187,0.5333333333333333,0.16504065040650406,0.304065040
65040652,0.40162601626016259,0.19349593495934958,0.23658536585365852,0.45365853658536587,0.53008130081300808,0.0,0.0,0.0,0.438211
38211382116,0.0,0.19024390243902439],
"Total": 31980, "ClassNr": 26}
```

Performance JSON for the training test (the JSONs for the testing and validation sets are similar)

Future improvements:

While I am rather confident about the feed-forward section of the network, the backpropagation section is still somewhat of a blur. The network suffers from exploding gradients when using softmax with a learning rate higher than 0.0001 and does not seem to be learning when using sigmoid. The accuracy obtained with a random model is usually around 3%. The highest accuracy I was able achieve was 28% before a NaN error was obtained (which is a common problem with CNNs, even with frameworks like tensorflow). The network is very slow and needs to be made more efficient. The use of optimizers will also help.