

# ACSE-8 Coursework |

## Question 1.

Because in our simplified version of Alexnet, we don't split the data into two GPUs, means that we need the input size  $n$  to satisfy:

$$\frac{n+2p-f}{s} + 1 = 55, \text{ where } p=0, f=11, s=4,$$

55 is the first layer of Alex net (on one GPU), this gives us  $n=227$ .

For the blanks in the spreadsheet, please see the Excel file in the ~~the~~ repository.

Question 2.

As the author argued in the Introduction, he claims that CNNs have much fewer connections and parameters, so they are easier to train comparing to standard feed forward neural networks. As we know, the parameters in a standard feed forward neural <sup>number of</sup> network strongly depend on the size of the input data, while CNNs not.

The main remaining obstacle to using even larger CNNs is that they are prohibitively expensive (in time) and ~~also~~ the limitation of the amount of memory available on current GPUs.

## Question 3

According to Section 3.1, as demonstrated in the plot, the deep convolutional neural networks work with ReLUs train several times faster than equivalents with tanh units (using less epochs to achieve same training error rate.) To conclude,

ReLU (non-saturating non linearity) does much better job than  $\tanh(x)$  or sigmoid (saturating nonlinearities) in the sense of fast learning, which is very important for training large datasets.

## Question 4.

The last layer is a fully connected <sup>convolutional</sup> neuron network with output size is  $1000 \times 1$ , 1 channel. It just convert the input to a length 1000

vector, with an appropriate activation function, represents the probability of the input for each category (because the Original ILSVRC dataset has 1000 categories). ~~We~~ Thus we could make a prediction for a certain input as the category which is the largest in this length 1000 <sup>corresponding</sup> vector.

Question 5.

The loss function should be a <sup>V-net</sup> cross-entropy,

$$J(w) = \sum_{x \in \text{Image}} w(x) \log p_L(x).$$

Where  $x$  is a pixel in the image and  $w(x)$  is the weighting function. also  $p_L(x)$  is a probability provided by Softmax function, and  $L$  is the true category where  $x$  is (our  $y$ ).

## Question 6.

Because the ILSVRC data set has 1000 categories, of which is just the labels for our supervised learning (similar to constraints in solving PDEs or functional extrema), and also  $10 \text{ bit} = 2^{10} \text{ bytes} = 1024 \approx 1000$ , so after compressing, the constraints has size of 10 bits.

## Question 7

We initialize the weights in each layer to ~~form~~ form a Gaussian distribution with mean of 0 and standard derivation of 0.01. Also we initialize the neuron bias in the 2nd, 4th and 5th convolutional layers with constant 1, initialize the neuron bias in the remaining layers with constant 0.

Question 8.

Not quite, ~~by~~ by our formula,

Output size =  $\frac{n + 2p - f}{s} + 1$ , after extracting

$224 \times 224$  from  $256 \times 256$ , we have output size

$$= \frac{256 + 0 - 224}{1} + 1 = 33, \text{ so the pre-filter}$$

has a size of  $33 \times 33$ , 2 channels (because

adding horizontal reflections), ~~the~~ then we have

the total factor is  $33^2 \times 2 = 2178 \neq 2048$ .

Question 9.

From Section 5, the author stated "We trained the network for roughly 90 cycles through the

training set of 1.2 million images." So the

epochs used is just 90 and the batch size

is  ~~$\frac{1.2 \times 10^6}{90} \approx 13333$~~  128 as he stated

in the beginning of Section 5.

Question 10.

The top-1 and top-5 errors on the dataset are 67.4% and 40.9% respectively attained by the net in Fall 2009 version of ImageNet with 10.084 categories and 8.9 million images, but with an additional, sixth convolutional layer over the last pooling layer.