

Assignment 3

Group 30

John Walker

Adam Fiksen

Giovanni Charles

November 20, 2012

1 Results

2 Fold Performance

3 Implementation Details

4 Questions

- 4.1 Discuss how you obtained the optimal topology and optimal values of network parameters. Describe the performance measure you used (and explain why you preferred it over other measures) and the different topologies / parameters you experimented with**

We experimented with the following parameters

1. Number of hidden layers
2. Number of hidden neurons in each layer
3. Transfer function
4. Training Function Parameters

We decided to make the search space for 1. between 1 - 2 we were able to limit the search space like this through research. For instance in the lecture notes for neural networks it says "Networks with many hidden layers are prone to overfitting and harder to train, for most problems one hidden layer should be enough, 2 hidden layers can sometimes lead to improvement". Through some rudimentary testing we established that as we had quite a large problem 2 layers did usually perform better although it took quite a lot longer so we decided that was the maximum we could feasibly have. Thus we decided the search space was 1 - 2.

We decided to make the search space for 2. between 5 - 30 layer 1, 0 - 17 layer 2. Layer 2 was dictated mainly through time constraints we found that if layer 1 was large enough any size of layer 2 much greater than 17 would take far too long to train. The constraints from layer 1 were chosen through a mixture of experience and research. We found that any number of neurons in layer 1 below 5 gave a horribly low classification rate so we set the lower bound as 5. We then did some research and found a few 'rules of thumb' two given by Jeff Heaton author of Introduction to Neural Networks In Java were "the optimal size of the hidden layer is usually between the size of the input and size of the output layers" "A good number of neurons is likely to be around the average of the inputs and outputs" we tried out some networks and noticed past 30 training time simply far surpassed any improvement and therefore we decided to make this the ceiling.

We limited our search space to the following transfer functions:

1. trainbfg
2. traincgb
3. traingda
4. traingdp
5. traingda

7. traingdx
8. trainlm
9. trainoss
10. trainrp
11. trainscg

These were selected from a list of around 20 gathered from the Matlab NNToolbox documentation. We then removed the training functions which weren't suitable for batch training and we were left with 13. We then read further into the documentation and decided to not use 'trainb', 'traingd' and 'traigdm' the reason being "The previous section presented two backpropagation training algorithms: gradient descent, and gradient descent with momentum. These two methods are often too slow for practical problems. In this section we discuss several high performance algorithms that can converge from ten to one hundred times faster than the algorithms discussed previously."

4.2 Explain what strategy you employed to ensure good generalisation ability of the networks and overcome the problem of overfitting

We optimised our networks based on the f1 measure for predictions of unseen target values, this would mean that overfitted networks should perform worse and be discarded by our genetic algorithm. To make our genetic algorithm faster we discarded networks which were likely to overfit the data by bounding some of the parameters. We made sure the number of neurons per level to be above 5, so it could understand the data sufficiently, and below 30 so that it is not overfitted. We also restricted the number of epochs to 100 so that the network does not become too familiar to the data.

4.3 In Part VIII you used the optimal parameters that you found in part VI to train your networks. However, there is a problem with this approach, the data you used for validation at some point will be used for testing in cross-validation. Can you explain why this is a problem? Ideally how should you optimise the parameters in cross-validation?

We were given a number of examples and targets we then divided these into 2/3 training and 1/3 validation data. The validation data then influences the neural network which is produced. Essentially you are training the neural network to perform well for the validation data and you hope that the validation data is large and representative enough so that the neural network you produce generalises well and can classify unseen examples. When you come to perform cross validation you then use a portion of the data as testing data and the rest as training/validation data. This means that eventually you will be testing a section of data which was the network was geared for this will produce a very low error and thus will skew the average error rate affecting your cross validation result. To combat this the parameters (data) given to cross validation should be rearranged into a random order this will reduce the effect of any 'seen' data and any effects will be evenly spread and the average will not be skewed, although it could be raised very slightly.

4.4 Is there any difference in the classification performance of the two different classification approaches. Discuss the advantages / disadvantages of using 6 single-output NNs vs. 1 six-output NNs

Six single output neural networks perform much better than a single network.

The six networks have the advantage of collectively forming a larger, more complicated and understanding, network without overfitting individually. The isolation between the networks allows for parallelism and for greater parameter specialisation for different emotions. One disadvantage is that having multiple networks may introduce merging conflicts when resolving an answer. This could be solved by limiting outputs to continuous output so that results are very unlikely to match or by ranking networks based on reliability, for example f1 measure, classification rate, most general/simple (i.e. fewer neurons), which also adds to the processing.

The single network however forms a smaller network which is less memory intensive and quicker to train and run. This means more aggressive optimisations can be performed resulting in parameters which are more finely tuned. This network would be more suited for machines where computing resources are limited.

5 Code Flow Chart