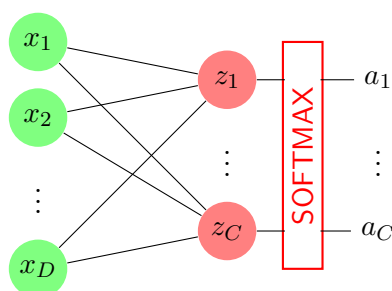


# Problem Sheet 4

## 1 Backpropagation for Multiclass Logistic Regression

We can view multiclass logistic regression as a simple neural network with no hidden layer. Let us suppose that the inputs are  $D$  dimensional and that there are  $C$  classes. Then the input layer has  $D$  units, simply the inputs  $x_1, \dots, x_D$ . The output layer has  $C$  units; there is full connection between the layers, *i.e.*, every input unit is connected to every output unit. Bias terms for each unit in the output layer are added separately. Finally, there is a non-linear activation function on the output layer, the softmax function. Pictorially, we may represent the neural network as follows:



Dropping superscripts to denote layers, as they are unnecessary in this case, using the notation in the lectures, we have:

$$\mathbf{z} = \mathbf{W}\mathbf{x} + \mathbf{b} \tag{1.1}$$

$$\mathbf{a} = \text{softmax}(\mathbf{z}) \tag{1.2}$$

Above  $\mathbf{W}$  is a  $C \times D$  matrix,  $\mathbf{b}$  is a column vector in  $C$  dimensions. Recall that the softmax function is defined as:

$$\text{softmax}(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{l=1}^C e^{z_l}} \tag{1.3}$$

For a point  $(\mathbf{x}, y)$  in the training data,  $y \in \{1, \dots, C\}$ . The model parameters to be learnt are  $\mathbf{W}$  and  $\mathbf{b}$ . We will use the cross entropy loss function, so the contribution of the point  $(\mathbf{x}, y)$  to the objective function is given by:

$$\ell(\mathbf{W}, \mathbf{b}; \mathbf{x}, y) = -\log a_y \tag{1.4}$$

Write the expressions for  $\frac{\partial \ell}{\partial \mathbf{a}}$ ,  $\frac{\partial \ell}{\partial \mathbf{z}}$ ,  $\frac{\partial \ell}{\partial \mathbf{W}}$  and  $\frac{\partial \ell}{\partial \mathbf{b}}$ . You should use the backpropagation equations to obtain these derivatives. You should express these derivatives for a single training point  $(\mathbf{x}, y)$ . When optimising, you'll usually average the gradient over a mini-batch before actually updating the parameters using a gradient step.



## 2 Digit Classification using MLPs

In this problem, we will consider a neural network to classify handwritten digits. You will also use this dataset in your practical this week. In the practical, you will implement a convolutional neural network. Here, we consider networks with fully connected layers and a different way to encode the targets (and the outputs of the network). You are strongly encouraged to implement these networks in tensorflow as well, to help you answer this question.

The inputs are vectors of length 784 (obtained from  $28 \times 28$  grey pixel images). Rather than output a class, the network will output a vector  $\hat{\mathbf{y}} \in \mathbb{R}^{10}$ . The target will be represented as a one-hot encoding, *e.g.*, if the label is 3, we will use the vector  $[0, 0, 1, 0, 0, 0, 0, 0, 0, 0]^T$  (by convention the digit 0 will be the last component in the encoding, not the first).

The data fed into the neural network consists of  $\langle (\mathbf{x}_i, \mathbf{y}_i) \rangle$ , where  $\mathbf{x}_i \in \mathbb{R}^{784}$  and  $\mathbf{y}_i \in \{0, 1\}^{10}$  is the one-hot encoding of the digit. The output of the neural network is also a vector  $\hat{\mathbf{y}} \in \mathbb{R}^{10}$ . The parameters in the network are the weights  $\mathbf{W}^i$  and biases  $\mathbf{b}^i$  for layers in the network. We'll train the networks by minimizing the squared loss objective (with respect to the parameters of the neural network):

$$\mathcal{L}(\mathbf{W}^i, \mathbf{b}^i) = \sum_{i=1}^N (\hat{\mathbf{y}}_i - \mathbf{y}_i)^2$$

The neural network to be used is shown in Figure 1. There are three layers, one input layer, one *hidden* layer and one output layer. The gradient of the loss function with respect to the parameters can be computed using the backpropagation algorithm as we've seen in the lectures. Answer the following questions:

1. To reduce the size of the network and increase efficiency, your colleague suggests using a binary encoding for the outputs instead of one-hot encoding; so 0 is encoded as 0000, 1 as 0001, and so on. In this case, the output layer only has four neurons rather than 10. Perhaps you could also reduce the number of neurons in the middle layer. What do you think about this suggestion?
2. Show that if you have an already trained neural network that has high accuracy with the *one-hot* encoding, you can design a network that uses binary encoding and achieves roughly the same error. (*Hint: You may need to add an additional layer.*)
3. What do you think would happen if you tried to train the neural network you suggested in the previous part directly (rather than adding the last layer by design)?

## 3 Neural Nets with Linear Activation Functions

Consider a neural network for multiclass classification. Let us assume that the input is 784 dimensional, we'll use one hidden layer, and there are 10 classes, *e.g.*, digits. For the hidden unit we'll use the "identity" activation function, *i.e.*,  $f(x) = x$ . Answer the following questions:

1. Suppose you use 25 hidden units and no regularization. For the output layer, we'll use a softmax non-linearity as we did for multiclass logistic regression in Problem 1. An

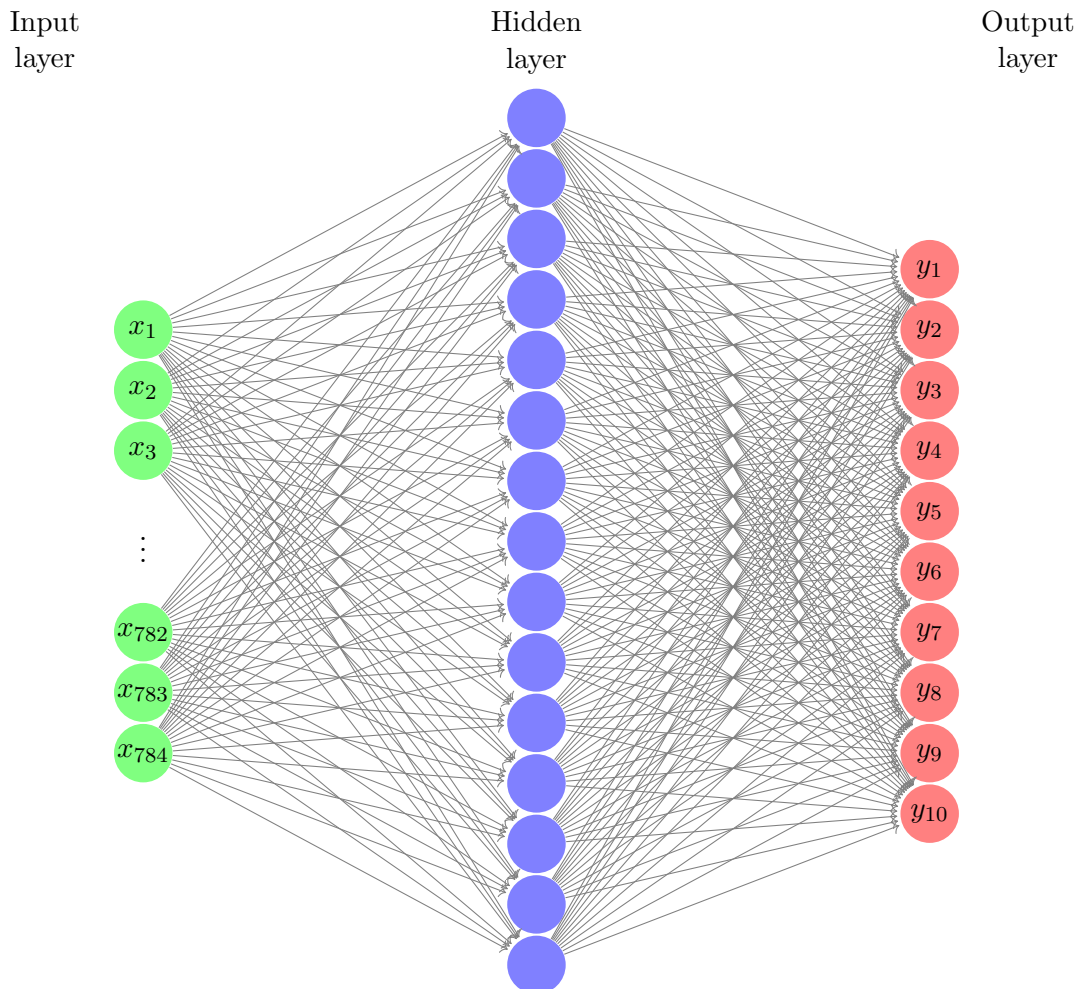


Figure 1: Three-layer neural network for handwritten digit classification.

alternative to such a neural network would be to just use the multiclass logistic regression approach directly. Which of the two models is more powerful, *i.e.*, able to represent possibly more complex relationships between inputs and outputs? (For this part, ignore the difficulty of training the models.)

2. Now suppose you only use 4 hidden units. What can you say about the relative power of the two models? Which model would you prefer?
3. For the network in Part 2, if you use the cross entropy loss function on the output layer, do you get a convex optimisation problem?