

# Course 02457 Non-Linear Signal Processing, Exercise 6

This exercise is based on C. M. Bishop: *Pattern Recognition and Machine Learning*, Chapter 5. Print and comment on the figures produced by the software as outlined below at the **Check-points**. You are going to use the ‘Neural Classification’ Matlab toolbox programmed at Digital Signal Processing section at IMM, DTU.

## Pattern recognition with neural networks

The neural network is designed to approximate posterior probabilities from a set of examples. This is achieved using the transformation from the neural network outputs  $\phi$  to the probabilities  $y$  by the normalized exponential or so-called *softmax* function:

$$y_k = \frac{\exp(\phi_k)}{\sum_j^c \exp(\phi_j)} . \quad (1)$$

This function has a built-in redundancy: If there are  $c$  classes and we use  $c$  outputs we can calculate the output of the  $c$ 'th unit from the sum of the others. A way to overcome this redundancy is via a *modified softmax* which only has  $c - 1$  outputs from the neural network (before the softmax). The  $c$ 'th output is then calculated from the others:

$$y_k = \frac{\exp(\phi_k)}{\sum_j^{c-1} \exp(\phi_j) + 1} , \quad k = 1, 2, \dots, c - 1 \quad (2)$$

$$y_c = 1 - \sum_k^{c-1} y_k . \quad (3)$$

## Pima indian data set

This is a data set where the task is to classify a population of women according to the risk of diabetes (binary classification). There are 7 input variables, 200 training examples and 332 test examples. 68 (34%) in the training set and 109 (32.82%) in the test set have been diagnosed with diabetes. In Brian Ripley's textbook *Pattern Recognition and Neural Networks* he states that his best method obtains about 20% misclassification on this data set so this is what we can use for reference. The input variables are:

1. Number of pregnancies
2. Plasma glucose concentration
3. Diastolic blood pressure
4. Triceps skin fold thickness
5. Body mass index (weight/height<sup>2</sup>)
6. Diabetes pedigree function
7. Age

The target output is 1 for examples having diabetes (“positive” diagnosis), and 0 for healthy subjects (“negative” diagnosis). The variables are not directly feed into the neural network. Instead the mean and standard deviation is calculated on the training set and these values are then used to normalize both the training and the test set, by first subtracting the mean and the dividing by the standard deviation. This procedure is called standardization, see Bishop page 567.

It is not possible to view all 7 dimension in one plot. A way to plot the data is to select to variables and plot them against each other. An other way is to use principal component decomposition on the data and plot one principal component against another.

### Checkpoint 6.1

Use the program `main6a.m` to optimize a neural network and plot the decision boundary and conditional probabilities. The heavy blue line is for locating the contour line corresponding to the equal probability of the two classes.

Run the network with different weight decays and comment on the difference in the decision surfaces. Create a flow chart of `main6a.m` and its functions.

## Entropic cost function

One could use the squared error function to optimize the neural network for classification:

$$E_{\text{square}} = 1/2 \sum_n^N \sum_k^c (y_k^n - t_k^n)^2 \quad (4)$$

The square error cost function is linked to ‘normal noise’ hypothesis. However, as  $y$  is suppose to model a probability and  $\mathbf{t} = (0, 1, 0, \dots, 0)$  is a position coded target vector, the Gaussian error is inappropriate. We can write the conditional likelihood function:

$$p(\mathbf{t}^n | \mathbf{x}^n) = \prod_k^c [y_k^n]^{t_k^n} \quad (5)$$

This will give an more appropriate *cross-entropy* cost function. The negative log-likelihood of this is for an entire data set of independent samples:

$$E_{\text{entropy}} = - \sum_n^N \sum_k^c t_k^n \ln y_k^n \quad (6)$$

### Checkpoint 6.2

The program `main6b.m` calculates the cost function value associate with one example for a two-class problem. It sweeps the output of the neural network  $\phi_1$  from -10 to 10.

What was the target  $t_1$ ? How much do the two cost functions penalize a “wrong” decision.

## Posterior probabilities

The posterior probability that the neural network outputs can be written using Bayes theorem (here with a one-dimensional input  $x$ ):

$$P(C_k|x) = \frac{p(x|C_k) p(C_k)}{p(x)} \quad (7)$$

### Checkpoint 6.3

Use the program `main6c.m` to plot posterior probabilities. The data is very simple: two classes each Gaussian distributed with  $\sigma_A^2 = \sigma_B^2 = 1$ ,  $\mu_A = -2$  and  $\mu_B = 2$ , thus  $p(x|C_A) \sim \mathcal{N}(-2, 1)$ . There are 100 samples in  $\mathcal{A}$  and 12 samples in  $\mathcal{B}$ , initially, change this to larger values (e.g. 500 and 60, to keep the prior probabilities the same) and observe how the posterior approaches the true posterior.

## Pruning

The program `main6d.m` will train and prune connections to optimize a neural network on the Pima indian set. The network is a multi-layer perceptron with softmax normalization.

We base pruning on the so-called ‘optimal brain damage’ (OBD) approach. The saliency of each weight is estimated using a second order Taylor expansion of the cost function around its minimum. This estimates the saliency (cost of removing the weight) of weight number  $q$  as  $S_q = \frac{1}{2} \frac{\partial^2 E}{\partial w_q^2} w_q^2$ .

The pruning procedure trains the initial configuration which fully connected between inputs and hidden layer and between hidden layer and outputs. The saliencies of all weights are computed and the least salient weight is pruned (put permanently to zero). The pruned network is retrained and the procedure repeated until a minimal configuration is reached. During pruning the test error is measured and after pruning the best network in test is selected.

### Checkpoint 6.4

Derive the saliency expression as the approximate cost of removing a single weight while keeping the remaining weights fixed at their ‘optimal’ values.

The best net is selected as the net with the lowest cost function value on the test set.

Note that not all inputs are used, is the number of pregnancies important for the diagnosis?

The target output is 1 for positive diagnosis (diabetes), and 0 for negative diagnosis (healthy).

What does the neural network say about the relation between age and diabetes (positive/negative correlation)?