

Intuition

Neural Networks are great for estimating behaviour of unknown input on base of training data. They can be thought of as Linear (or Logistic) Regression with adaptive feature selection, what makes them more suitable to describe the problem. Structure and adaptiveness of Neural Networks resemble the human brain, therefore they are often compared to processes running in our brains.

This flyer will introduce you to working with Neural Networks and coding the fundamentals in Octave. With this you will be able to run more sophisticated algorithms using e.g. Octave's Neural Network Package ¹. Besides the following topics are covered:

- Understanding the network architecture
- Calculating results with Forward Propagation
- Setting up cost functions for multiclass networks
- Calculating weights with Back Propagation
- Initialization and suggested defaults

Neural Networks are a way of supervised learning (i.e. learning from training data with known result), as is Linear (or Logistic) Regression. To get the whole picture we encourage you to also read the [Linear Regression Flyer](#). Choosing one over the other is mainly a question of computational cost and accuracy. While Regression is cheaper, Networks are often more precise. For further reading on Neural Networks you can use Google or one of these sites as a starting point:

- Wikipedia ²
- Machine Learning @ The Stanford OpenClassroom ³
- NN FAQ from comp.ai.neural-nets ⁴

The flyer can be distributed free of charge, as long as it is not altered in any way. If you have questions or comments contact the flyered science project.

fsadmin@flyered-science.org

Numerical Solution (in Octave)

Like in Linear Regression we will use a numerical solver to compute $\Theta^{(l)}$. We already defined the cost function and the gradients needed, to pass it along with an initial guess of parameters to a suitable linear solver.

The procedure needs the cost function parameters to be vectors. Therefore the weight matrices $\Theta_{ji}^{(l)}$ have to be "unrolled" to a vector, to pass it along, and then "reshaped" again before calculating cost and gradient:

```
ThetaVec = [Theta1 (:); Theta2 (:); ...  
  
function [J, grad] = neuralCF(ThetaVec ...  
    Theta1 = reshape(ThetaVec(1:sizeL1* ...  
        (sizeL0+1),sizeL1 ,sizeL0+1)  
    (...)  
endfunction
```

With an initial $\Theta_{ji}^{(l)}$ you calculate the forward propagation $a_l = \text{sigmoid}(a_{l-1}\Theta_l')$. (With optimized $\Theta^{(l)}$ this is how the network is used later on.

Cost function and gradients are then summed over all test samples, as explained before. The gradient sum in Octave can be written as follows:

```
%compute backpropagation  
delta_output = a2(counter,:) - yvec;  
delta_hidden = (delta_output*Theta2) ...  
    .* ((a1 .* (1-a1))(counter,:));  
  
%sum gradients  
T1grad = T1grad+(delta_hidden(2:end)' ...  
    .* X(counter,:));  
T2grad = T2grad+(delta_output' ...  
    .* a1(counter,:));
```

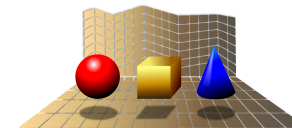
One can optimize Θ using the cost function with `fminunc` as depicted in Linear Regression flyer (Math.1).



Neural Networks

An introduction to representation and algorithm, with implementation in Octave

The Flyered Science Series provides condensed factsheets on scientific topics. Persons new to the subject will find the intuitive introduction with links to further reading quite useful, while the experienced user benefits from the structured representation, formulary and implementation examples. Fits into every pocket, serves perfectly as bookmark in your textbooks.



¹http://www.plexso.com/61_octave/index.html

²http://en.wikipedia.org/wiki/Artificial_neural_network

³<http://openclassroom.stanford.edu>

⁴<ftp://ftp.sas.com/pub/neural/FAQ.html>

Theory

Forward Propagation

Neural networks are used to estimate the result $y = h^{(\Theta)}(x)$ of unknown x on base of training data $x^{(l)}$ with known regression. The problem is comparable to linear (or logistic) features x , the neural network introduces (hidden) layers with adaptive features (activations) $a^f_{(l)}$, that are calculated from the preceding layer using the weight matrix $\Theta^{(l-1)}$.

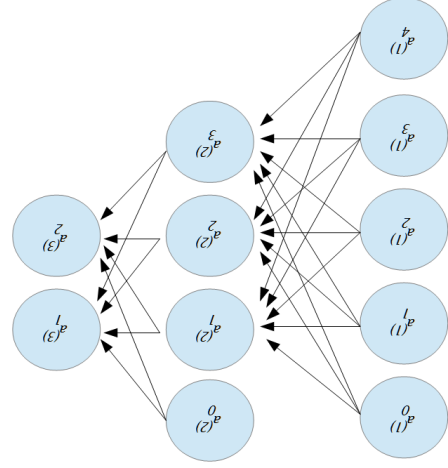


Figure 1: Neural network layout with one hidden layer.

Input and hidden layer(s) can contain a constant offset (bias term) $a^0_{(l)} = 1$. Given L different layers (i.e. $L - 2$ hidden layers) the output layer hypothesis $h^{(\Theta)} = a^{(L)}$ is calculated from input activations $a^i_{(0)} = x_i$ considering the weight matrix Θ and the activation formula $g(\Theta^T a)$. This is called "Forward Propagation".

In this layer we use the sigmoid function as activation formula, i.e. $g(\Theta^T a) = \frac{1}{1 + e^{-\Theta^T a}}$. This is called "Logistic Unit". Using a "Linear Unit" (i.e. weighted sum), a step function or a linear combination is straight forward.

Weight Matrices

We now discuss how to derive the weight matrix Θ from a training set of labeled data (known result) $(x^{(l)}, y^{(l)})$. Suppose that activation $a^{(l-1)}$ in layer $l - 1$ contributes to activation $a^f_{(l)}$ in layer l with the weight $\Theta^{ff}_{(l-1)}$. This yields $a^f_{(l)} = g(\Theta^{ff}_{(l-1)})^T a^{(l-1)}$, as shown below.

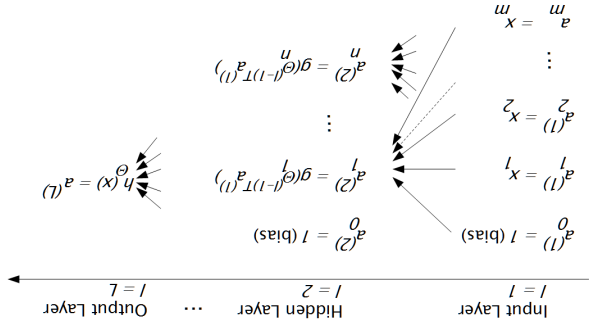


Figure 2: Calculating activations in neural networks.

Like in logistic regression (see also Hyper Math.1) a cost function is calculated for the training set, that considers results $y^{(l)}$ and hypothesis $h^{(\Theta)}(x^{(l)})$ from Forward Propagation. With m training sets and K classes in the output layer the cost function J then becomes

$$J(\Theta) = - \sum_{k=1}^m \sum_{K} \frac{1}{m} \log(h^{(\Theta)}(x^{(l)})) + \sum_{k=1}^m (1 - y^{(l)}) \log(1 - h^{(\Theta)}(x^{(l)})) + \frac{\lambda}{2} \sum_{l=1}^L \sum_{m=1}^m \sum_{f=1}^f (\Theta^{ff}_{(l)})^2$$

The first term applies for results $y^{(l)} = 1$ and is zero, when the hypothesis matches the result. The second term applies for results $y^{(l)} = 0$ in the same manner. The third term is called the regularization term, which prevents overfitting (i.e. small Θ s are preferred, see also Hyper Math.1).

The objective is to solve $\max_{\Theta} J(\Theta)$ for an optimal set of weight matrices $\Theta^{(l)}$. This is done by "Backpropagation" of the deviations between result and hypothesis.

Backpropagation

The idea of Backpropagation is to consider deviations of the trained network hypothesis from the real result. To find the optimal set of $\Theta^{(l)}$, the deviations have to be minimized. In the output layer L the deviation is calculated as

$$\delta^f_{(L)} = a^f_{(L)} - y^f$$

The deviation can be propagated back through the network using the following equation:

$$\delta^{(L-1)} = (\Theta^{(L)} \Theta^{(L-1)})^T \delta^{(L)} \cdot * \cdot \delta^f_{(L)} \cdot * \cdot \delta^{(L-1)}$$

The $*$ operator means elementwise multiplication. Without considering regularization (i.e. $\lambda = 0$) the deviation of the activation function g can be written as

$$\delta^f_{(L-1)} = \delta^{(L-1)} \cdot (1 - a^{(L-1)}) \cdot a^{(L-1)}$$

The task is to compute all $\delta^{(l)}$. Remember: The input parameters in the input layer have zero deviation: $\delta^{(0)} = 0$. Considering only one test example the gradients can be computed from this:

$$\frac{\partial \Theta^{ff}_{(l)}}{\partial a^f_{(l+1)}} J(\Theta) = \delta^f_{(l+1)}$$

Considering the whole test set, the deviations are then the sum of all test set contributions:

$$\Delta^{ff}_{(l)} = \sum_m \delta^f_{(l+1)}$$

The gradient calculates accordingly:

$$\frac{\partial \Theta^{ff}_{(l)}}{\partial a^f_{(l+1)}} J(\Theta) = \frac{1}{m} \Delta^{ff}_{(l)} + \lambda \Theta^{ff}_{(l)} \quad \text{for } f \neq 0$$

$$\frac{1}{m} \Delta^{ff}_{(l)} = 0 \quad \text{for } f = 0$$

Remember: The bias terms $f = 0$ are not included in regularization. Following the gradients using an optimization algorithm then yields the optimal set of $\Theta^{(l)}$ for the given training set.