

Multi-layer Perceptrons

4/27/2008

Analogy to Neural Networks in the Brain.
 - probably misguided.

Perception.

$$y = \sum_{j=1}^d w_j x_j + w_0$$

threshold function.

hard $S(a) = \begin{cases} 1, & \text{if } a > 0 \\ 0, & \text{otherwise.} \end{cases}$

soft. $y = \bar{\sigma}(\underline{w}^T \underline{x}) = \frac{1}{1 + e^{-\underline{w}^T \underline{x}}}$. $\bar{\sigma}(\cdot)$ sigmoid function.

There are a variety of different algorithms to train a perceptron from labelled examples

Example: Quadratic error.

$$E(\underline{w} | \underline{x}^t, r^t) = \frac{1}{2} (r^t - y^t)^2 \quad y = \underline{w}^T \underline{x}$$

update rule $\Delta w_j^t = \eta (r^t - y^t) x_j^t$.

$$E(\{\underline{w}_i\}, \underline{x}^t, r^t) = -\sum_i (r_i^t \log y_i^t + (1-r_i^t) \log (1-y_i^t))$$

y^t = sigmoid $(\underline{w}^T \underline{x}^t)$.

update rule $\Delta w_j^t = \eta (r^t - y^t) x_j^t$.

Update = Learning factor. (Desired Output - Actual Output)
 x Input.

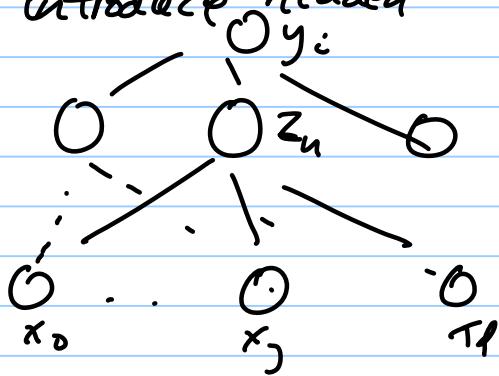
(2)

Multilayer Perceptron.

A single layer perceptron can only approximate linear functions of the input — i.e. the set of perceptrons has limited capacity.

Multilayer perceptrons were invented to increase the capacity — introduce hidden units (nodes)

$$z_h = \sigma(\underline{\omega}_h^T \underline{x})$$
$$= \frac{1}{1 + \exp\left\{-\sum_{j=1}^d \omega_{hj} x_j + \omega_{h0}\right\}},$$
$$h = 1, \dots, H.$$



$$\text{Output } y_i = \underline{v}_i^T \underline{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}.$$

Other output function — e.g. $y_i = \sigma(\underline{v}_i^T \underline{z})$.

Many levels can be specified.

What do the hidden units represent?
Unclear, but many people have tried to explain them.
Any input-output function can be represented
as a multilayer perceptron with enough hidden units
— infinite capacity.

(3)

How to train a multilayer perceptron?

Unknown parameters - the weights ω_{kj} , v_{ij} .

Define an error function:

e.g. $E[\omega, v] = \sum_i |y_i - \sum_h v_{ih} \text{sigmoid}(\sum_j \omega_{hj} x_j)|^2$

update $\Delta \omega_{hj} = -\frac{\partial E}{\partial \omega_{hj}}$

$$\Delta v_{ik} = -\frac{\partial E}{\partial v_{ik}}$$

These gradients can be computed by the chain rule of differentiation

$$\frac{\partial E}{\partial \omega_{hj}} = \frac{\partial E}{\partial y_i} \cdot \frac{\partial y_i}{\partial z_h} \cdot \frac{\partial z_h}{\partial \omega_{hj}}$$

For historical reasons, this is known as the backpropagation algorithm.

→ the errors would be projected back to the hidden units. It was thought hoped that this might be biologically plausible.

(4)

More details.

quadratic error.

For the 2nd layer.

$$\Delta \omega_{ih} = \eta \sum_t (r^t - y^t) z_h^t$$

For the 1st layer

$$\begin{aligned}\Delta \omega_{ij} &= -\eta \frac{\partial E}{\partial \omega_{ij}} = -\eta \sum_t \frac{\partial E}{\partial y^t} \cdot \frac{\partial y^t}{\partial z_h^t} \cdot \frac{\partial z_h^t}{\partial \omega_{ij}} \\ &= -\eta \sum_t -(r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t \\ &= \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t.\end{aligned}$$

error term for
the hidden units

Batch learning versus stochastic learning.
→ advantages of stochastic learning.

+ epoch.

Training procedures:

momentum

$$\Delta \omega_i^t = -\eta \frac{\partial E^t}{\partial \omega_i} + \alpha \Delta \omega_i^{t-1}$$

OVERTRAINING - .

(5) Multilayer Perceptrons and Nonlinear Regression

multilayer perceptrons can be thought as non linear regression — and forget about the several interpolations.

Output function

$$y_i(\underline{x}) = \sum_h v_{ih} \sigma\left(\sum_j w_{hj} x_j\right)$$

regression:

$$y_i(\underline{x}) = \sum_h v_{ih} \bar{o}\left(\sum_j w_{hj} x_j\right) + c$$

\bar{o} zero mean
 Gaussian noise.

Cost function: $\sum_{t=1}^N \sum_i |y_i^t - \sum_h v_{ih} \bar{o}\left(\sum_j w_{hj} x_j^t\right)|^2$

(-log of) probability

Gauss score Since this is nonlinear, we cannot obtain an analytic solution for the coefficients $\{v_{ih}\}, \{w_{hj}\}$.

Instead, we have to minimize the cost function algorithmically.

(6)

Multilayer Perceptron and Steepest Descent

Strategy: Batch Mode

use all the data, i.e. full cost function.

steepest descent

$$\frac{d v_{ih}}{dt} = -\frac{\partial E}{\partial v_{ih}} \quad \text{E cost function}$$

$$\frac{d w_{hj}}{dt} = -\frac{\partial E}{\partial w_{hj}}$$

discrete formulation

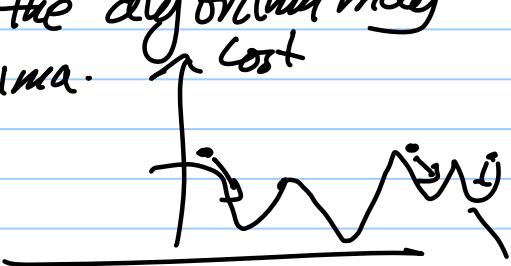
$$v_{ih}(t+\Delta t) = v_{ih}(t) - \Delta t \frac{\partial E}{\partial v_{ih}}$$

$$w_{hj}(t+\Delta t) = w_{hj}(t) - \Delta t \frac{\partial E}{\partial w_{hj}}$$

(page 324).

The cost function will usually have many local minima, so the algorithm may get stuck in a local minimum.

Final result depends on initial conditions.



(7) Multilayer Perceptron & Stochastic Descent

Perform one iteration from each data sample to

Current state: $\{v_{ih}, w_{hj}\}$

Select data sample y_i^t, x_j^t .

$$\underline{\text{set}} \quad v_{ih} \rightarrow v_{ih} + \Delta v_{ih}$$

$$w_{hj} \rightarrow w_{hj} + \Delta w_{hj}$$

$$\text{with} \quad \Delta v_{ih} = -\frac{\partial}{\partial v_{ih}} \sum_i \left(y_i^t - \bar{v}_{ih} \sigma(\sum_j w_{hj} x_j^t) \right)^2$$

$$\Delta w_{hj} = -\frac{\partial}{\partial w_{hj}} \sum_i \left(y_i^t - \bar{v}_{ih} \sigma(\sum_j w_{hj} x_j^t) \right)^2$$

This is "stochastic" because the algorithm has a random element - the choice of which data sample to use.

Stochastic descent is able to escape from some local minima (theoretical results).

(7) Critical Issues for multilayer perceptrons

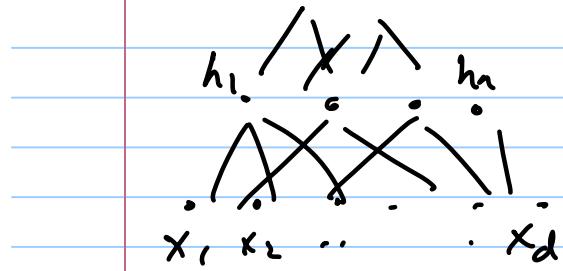
— how many hidden units to use?

Book describes several techniques for dealing with this → for example, having more hidden units than you need and then penalizing the weights.

E.g. Add term $\sum_{h,j} (\omega_{hj}^2) + \sum_{i,h} (v_{ih}^2)$ to the cost function → intuition, if the weights are small to a hidden unit, then the hidden unit is not used.

In practice, some of the most effective multi-layer perceptions are those which the structure was hand designed.

(9) Multilayer Perceptrons /SVM/ AdaBoost (next lecture).



$$y_i = \sum_k v_{ik} h_k$$

$$h_k = \bar{o}(\sum_j w_{kj} x_j)$$

SVM can also be represented in this way.

$$y = \text{sign}\left(\sum_{\mu} \alpha_{\mu} y_{\mu} \underline{x}_{\mu} \cdot \underline{x}\right)$$

hidden units response $\underline{x}_{\mu} \cdot \underline{x} = h_{\mu}$.

$$y = \text{sign}\left(\sum_{\mu} \alpha_{\mu} y_{\mu} h_{\mu}\right)$$

Advantage of SVM — number of hidden units is given by the no. of support vectors.
 $\rightarrow (\alpha_{\mu})$ specified by minimizing the primal problem (well defined algorithm to perform this minimization).

(1b) Multilayer Perceptrons / SVM / AdaBoost (next lecture).

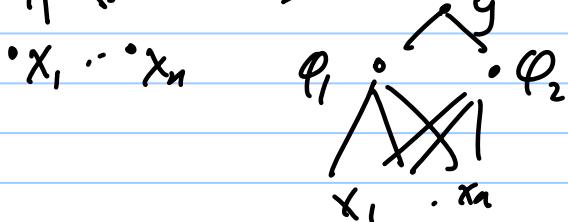
AdaBoost (next lecture)

Build a multilayer perceptron incrementally

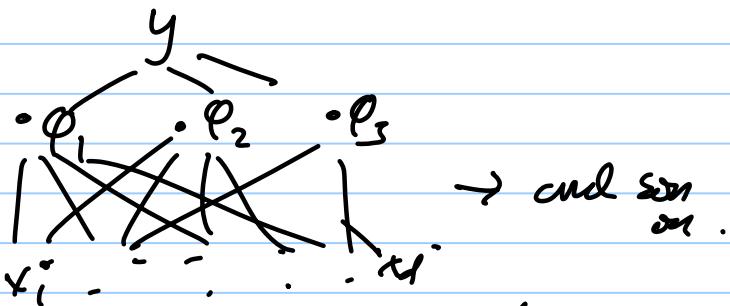


one hidden unit. $\Phi_1(x)$

→ select next hidden unit $\Phi_2(x)$



→ new unit $\Phi_3(x)$



— AdaBoost selects the hidden units incrementally by minimizing a criterion.