

# Neural Networks

## 1 Preprocessing

### 1.1 Input Normalization

Linear rescaling to avoid scaling problems. Useful for radial basis function and multilayer networks. Each input treated independently for scaling. Calculate mean and variance of training set  $\mathbf{x}^{(m)}$  with  $m \in [0, M-1]$  and normalize input vectors to  $\mu_i = 0$  and  $\sigma_i = 1$ .

### 1.2 Feature Extraction

For **speech recognition** use AnaFB, power estimation ( $|\cdot|^2$ ), mel-filterbank (15-30 normalized overlapping triangular filters), cepstrum estimation (IDFT/IDCT of log power spectrum - to decorrelate input features - symmetric), temporal features (combine successive features vectors - delta or delta-delta features), dimension reduction (feature space transformation, e.g. LDA - dimensionality reduction by preserving class discriminatory information - variance of features corresponding to one class is minimized, distance between classes maximized).

#### 1.2.1 LDA

Dataset of  $M$  observations of  $N$ -dimensional Euclidian variable  $\mathbf{x}$ . Project  $\mathbf{x}$  to one dimension using a projection vector  $\mathbf{w}$ , such that  $y^{(m)} = \mathbf{w}^T \mathbf{x}^{(m)}$ . Cost function  $\max_{\mathbf{w}} \mathbf{w}^T (\mu_1 - \mu_0)$  to separate class means. **Fisher's idea:** Large separation between projected class means while granting small variance within each class:

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}} \Rightarrow \mathbf{w}_{\text{opt}} = \mathbf{S}_w^{-1} (\mu_1 - \mu_0)$$

with  $\mathbf{S}_b$  being between class covariance matrix,  $\mathbf{S}_w$  within class covariance matrix.

## 2 Threshold Logic Units - Single Perceptrons

Neural Network is a blackbox/mapping machine/network of functions with  $N$ -dim. input  $\mathbf{x}$  and  $M$ -dim. output  $\mathbf{y}$ . A computing element/unit/node has unlimited fan-in and is a primitive function  $f$  of  $N$  arguments.

### 2.1 McCulloch-Pitts Neuron

Binary signals transmitted over edges produce binary result, driven by threshold  $\Phi$ .  $N$  excitatory edges  $x_i$  and  $K$  inhibitory edges  $x'_j$ . Inhibitory edges can inactivate the whole unit if at least one is active. Else it works as a threshold gate. Neuron is activated if  $\sum_{k=0}^{N-1} x_k \geq \Phi$ .

### 2.2 Perceptron

With threshold  $\Phi$  and input weights  $w_k$ , so the output is active if  $\sum_{k=0}^{N-1} w_k x_k \geq \Phi$ . Simple Perceptron and McCulloch-Pitts Unit are equivalent. Geometric interpretation can be used to check if a function can be computed by a Perceptron using a separating line/plane. XOR function with 2 vars can not be solved by one perceptron - not linearly

separable ( $0 < \Phi$ ,  $w_0 > \Phi$ ,  $w_1 > \Phi$ ,  $w_0 + w_1 < \Phi$ ). Two sets of  $A$  and  $B$  with different output values are linearly separable if  $\sum_{x_k \in A} x_k w_k \geq \Phi$  and  $\sum_{x_k \in B} x_k w_k < \Phi$ . Solve XOR problem using network of perceptrons  $(x_0 \wedge \bar{x}_1) \vee (\bar{x}_0 \wedge x_1)$ . First layer labels the region, output layer decodes the classification.

### 2.3 Training

**2.3.0.1 Supervised learning** Weights of the network are initialized randomly. On misclassification network parameters are adjusted. Desired output is always known, thus it is called learning with a teacher. In **reinforcement learning** only input vector is used for the weight adjustment. In **corrective learning** the magnitude of the error together with the input vector are used for the weight correction. Training rule:

$$\Phi_{\text{new}} = \Phi_{\text{old}} + \Delta\Phi, \quad w_{\text{new},i} = w_{\text{old},i} + \Delta w_i$$

**2.3.0.2 Unsupervised learning** For a given input the output is unknown - learning without a teacher.

### 2.4 Perceptron Learning Algorithm

Transform perceptron into zero-valued threshold by turning  $\Phi$  into a weight  $w_N$  (bias) such that  $\sum_{k=0}^{N-1} x_k w_k - \Phi \geq 0$ .

*start* :  $\mathbf{w}_0$  is generated randomly

$n = 0$

*test* : Select  $\mathbf{x} \in P \cup F$  randomly

if  $\mathbf{x} \in P$  and  $\mathbf{w}_t \mathbf{x} > 0$  go to test

if  $\mathbf{x} \in P$  and  $\mathbf{w}_t \mathbf{x} \leq 0$  go to add

if  $\mathbf{x} \in N$  and  $\mathbf{w}_t \mathbf{x} < 0$  go to test

if  $\mathbf{x} \in N$  and  $\mathbf{w}_t \mathbf{x} \geq 0$  go to subtract

*add* : set  $\mathbf{w}_{n+1} = \mathbf{w}_t + \mathbf{x}$  and  $n = n + 1$ , goto test

*subtract* : set  $\mathbf{w}_{n+1} = \mathbf{w}_t - \mathbf{x}$  and  $n = n + 1$ , goto test

### 2.5 Fast learning algorithm - delta rule

If  $\mathbf{x} \in P$  is classified erroneously we have  $\mathbf{w}_n^T \mathbf{x} \leq 0$  so we get the error  $\Delta(n) = -\mathbf{w}_n^T \mathbf{x}$  so that the weight vector gets corrected:  $\mathbf{w}_{n+1} = \mathbf{w}_n \frac{\Delta(n) + \epsilon^+}{\|\mathbf{x}\|^2} \mathbf{x}$ . So that

$\mathbf{w}_{n+1}^T \mathbf{x} = \left( \mathbf{w}_n \frac{\Delta(n) + \epsilon^+}{\|\mathbf{x}\|^2} \mathbf{x} \right)^T \mathbf{x} = \epsilon^+ > 0$ .  $\epsilon$  guarantees that the new weight vector barely skips over the border of the region with higher error. For  $\mathbf{x} \in F$  use  $\epsilon^-$ . A variant is using  $\gamma(\Delta(n) + \epsilon)\mathbf{x}$  with  $\gamma$  as learning factor.

### 2.6 Convergence

If  $P$  and  $N$  are finite and linearly separable.

$$\begin{aligned} \cos(\phi) &= \frac{\mathbf{w}_{\text{des}}^T \mathbf{w}_{n+1}}{\|\mathbf{w}_{n+1}\|} \\ &= \frac{\mathbf{w}_{\text{des}}^T \mathbf{w}_0 + (n+1)\delta}{\sqrt{\|\mathbf{w}_0\|^2 + (n+1)}} \end{aligned}$$

## 3 Multilayer perceptrons

ANN is a  $(L, C)$  tuple with a set of nodes  $L$  and a set of directed edges  $C$ .  $c = (v, l) \in C$  are directed edges from node  $v$  to node  $l$ .  $L_{in}$  input layer subset,  $L_{hid}$  hidden layer subset,  $L_{out}$  output layer. Set  $L_l^{(pre)}$  consists of prior nodes of node  $l$  (predecessors). Set  $L_l^{(suc)}$  consists of subsequent nodes of node  $l$  (successors). Every edge has a weight  $w_v^{(l)}$ .

**3.0.0.3 Generalized neuron** Consists of network input function  $f_{net}^{(l)}$ , activation function  $f_{act}^{(l)}$  and network output function  $f_{out}^{(l)}$  and three states respectively  $y_{net}^{(l)}$ ,  $y_{act}^{(l)}$ ,  $y_{out}^{(l)}$ , also external input  $x_{ref}^{(l)}$ .

**3.0.0.4 Properties of multilayer perceptrons** A NN with layered architecture doesn't contain cycles. A MP is a feed-fwd NN with strictly layered structure. Normally all units in a layer connected to all other units of the next layer.

**3.0.0.5 Weight matrices** Given  $L_{in} = \{v_0, \dots, v_{N-1}\}$  and  $L_{out} = \{l_0, \dots, l_{M-1}\}$  the connection weights are gathered in a matrix

$$\mathbf{W} = \begin{bmatrix} w_{v_0}^{(l_0)} & w_{v_1}^{(l_0)} & \dots & w_{v_{N-1}}^{(l_0)} \\ w_{v_0}^{(l_1)} & w_{v_1}^{(l_1)} & \dots & w_{v_{N-1}}^{(l_1)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{v_0}^{(l_{M-1})} & w_{v_1}^{(l_{M-1})} & \dots & w_{v_{N-1}}^{(l_{M-1})} \end{bmatrix} \in \mathbb{R}^{M \times N}$$

so that  $\mathbf{y}_{net}^{(L_1)} = \mathbf{W} \mathbf{y}_{out}^{(L_0)}$

Copyright © 2014 Major Ring Ding Ding Dong feat. Jingjong Ba-Dingdong