

- Computational algorithm inspired by the parallelism of human brain
- Tries to mimic the cognitive ability of the human brain to identify patterns and to make decisions and forecasts based on past experience

How NN works?

- Non-linear models that can be trained to map the past and the future values of a time series
- extract hidden structure and relationship governing the data
- powerful pattern recognition capabilities

NN in the context of statistical methods

- Multivariate, non-linear, non-parametric inference technique
- Completely data-driven and assumption free
- ability to model non-linear dynamics.

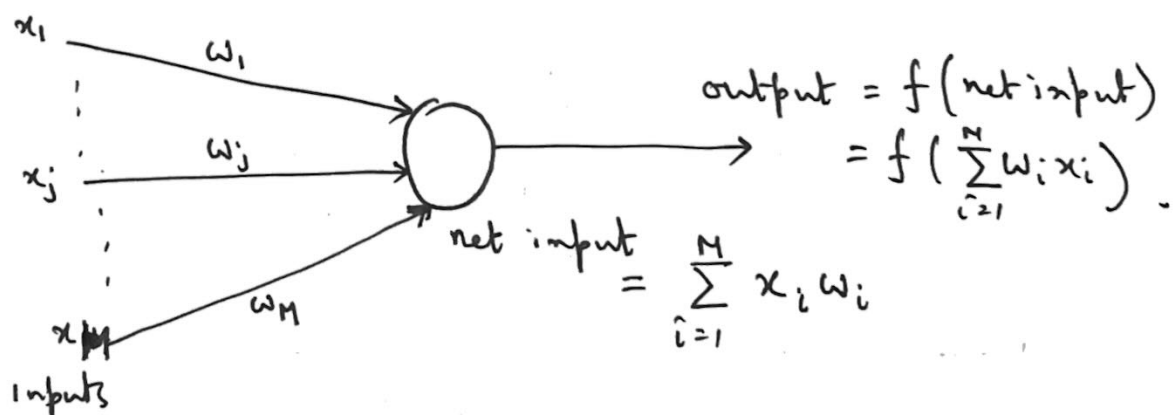
Some application areas

Stock market prediction, foreign exchange trading, bankruptcy prediction, volatility prediction, economic indicator forecasting, credit rating, prediction of international airline passenger traffic, nowcasting in aviation industry, ozone level prediction

Human brain:

Hundred billion neurons - interconnected structure

Basic computational unit of ANN - a neuron - an artificial neuron



x_j : input value from the j^{th} input node

w_j : weight on the connection from the j^{th} input node to the PU.

PU - collects the inputs from all the inputs to get the net input value

$$\sum_{i=1}^M w_i x_i \quad ; \quad w_i \text{ denoting the relative importance of } i^{\text{th}} \text{ input}$$

PU processes the net input value and generates the output

$$\text{Out} = f \left(\sum_{i=1}^M w_i x_i \right)$$

↑
transfer function / squashing f^n .

A The simplest NN

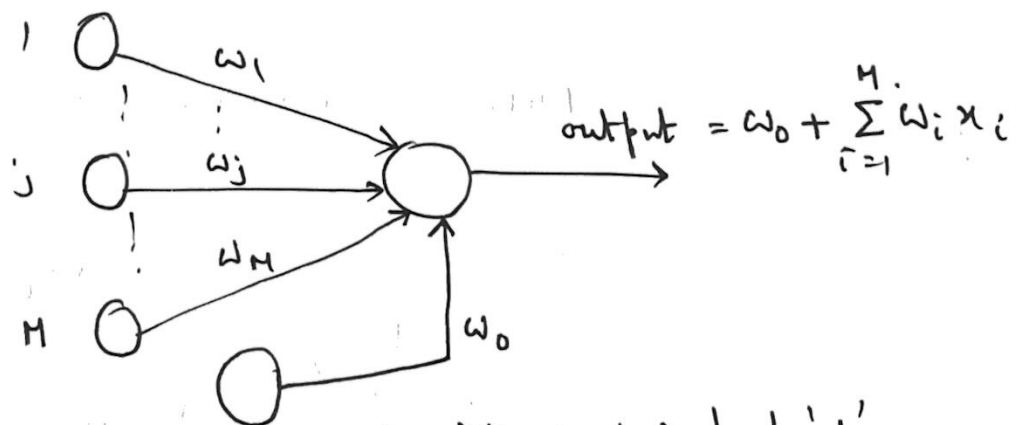
A single PU network.

MLR setup $\mathcal{A} = \{ (x_1, y_1), \dots, (x_n, y_n) \}$.

Consider a transfer f^n as identity

$$\text{output} = f \left(\sum_{i=1}^M w_i x_i \right) = \sum_{i=1}^M w_i x_i$$

MLR model without a constant term



bias unit with const input '1'.

Unknown parameters : weight connections (w_1, \dots, w_n)
or (w_0, w_1, \dots, w_n)

optimum values of (w_0, w_1, \dots, w_M) can be obtained by minimizing residual ssq.

Note: Injecting non-linearity into the model -

consider a non-linear transfer f^n .

e.g (i) A sigmoidal transfer f^n .

$$f(x) = \frac{1}{1 + e^{-x}}$$

(ii) A tanh transfer f^n

$$f(x) = \tanh(x).$$

For (i), output = $\frac{1}{1 + e^{-(w_0 + \sum_{i=1}^M w_i x_i)}} = g(x)$

A non-linear ^{model} of response in terms of the values of the feature vector components

Let (\tilde{x}_p, y_p) denote the p th pattern in \mathcal{L} .

Network output: $O_p = f\left(w_0 + \sum_{j=1}^M x_{pj} w_j\right) / f\left(\sum_{j=1}^M x_{pj} w_j\right).$

Actual output: y_p

Error for the p th pattern: $y_p - o_p$

Error square: $E_p = \frac{1}{2}(y_p - o_p)^2 \leftarrow \text{instantaneous error for } p\text{th pattern}$

Total error over \mathcal{L}

$$E = \frac{1}{2} \sum_{p=1}^n E_p = \frac{1}{2} \sum_{p=1}^n (y_p - o_p)^2.$$

Training of the network: determination of optimum weight vector.

Modes of learning : instantaneous mode or single pattern adaptation mode
or Batch mode or multiple pattern adaptation mode.

Method for obtaining the weight updation : Gradient descent algorithm

Instantaneous mode!

$$E_p = \frac{1}{2} (y_p - o_p)^2$$

$$E_p = \frac{1}{2} (y_p - f(\sum_{j=1}^M w_j x_{pj}))^2$$

Gradient of the error E_p w.r.t. the weights

$$\begin{aligned} \frac{\partial E_p}{\partial w_j} &= -(y_p - o_p) \cdot \frac{\partial f(\sum w_j x_{pj})}{\partial w_j} \\ &= -(y_p - o_p) \frac{\partial f(\sum w_j x_{pj})}{\partial (\sum w_j x_{pj})} \cdot \frac{\partial (\sum w_j x_{pj})}{\partial w_j} \\ &= -(y_p - o_p) f'(\sum w_j x_{pj}) x_{pj} \end{aligned}$$

e.g. for sigmoidal transfer f^* .

$$f(\sum w_j x_{pj}) = \frac{1}{1 + e^{-\sum w_j x_{pj}}} = o_p$$

$$\left[f(x) = \frac{1}{1 + e^{-x}} ; f'(x) = \frac{+e^{-x}}{(1 + e^{-x})^2} = \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}} \right) \right]$$

$$f'(\sum w_j x_{pj}) = o_p(1 - o_p)$$

$$\Rightarrow \frac{\partial E_p}{\partial w_j} = -(y_p - o_p) o_p(1 - o_p) x_{pj}$$

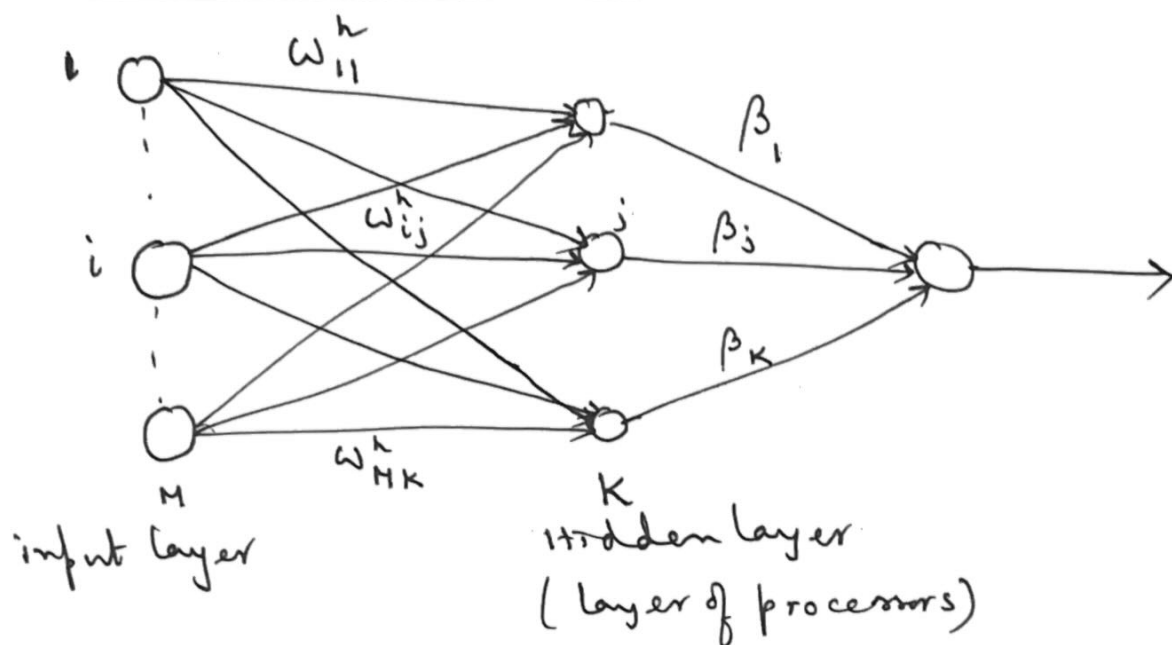
Gradient descent algorithm weight updation

$$w_j^{(k+1)} = w_j^{(k)} + \eta \left(- \frac{\partial \epsilon_p}{\partial w_j} \right) \Big|_{w_j^{(k)}}$$

$$w_j^{(k+1)} = w_j^{(k)} + \eta (y_p - o_p) o_p (1 - o_p) x_{pj}$$

Present the patterns at random and use the above equation for weight updation untill stopping criterion is reached.

Multi-layer Perceptron (MLP)



M input, K hidden layer neurons and one output
 $M-K-1$ MLP.

p th pattern $\rightarrow (x_p, y_p)$

$$\text{net}_{pj}^h = \sum_{i=1}^M w_{ij}^h x_{pi} \quad (\text{may have a bias unit})$$

net input value at j th neuron of hidden layer

output from the j th neuron

$$f(\text{net}_{pj}^h) = o_{pj}$$

i_{pj} : j^{th} input to the output unit

net input value at output

$$\sum_{j=1}^K \beta_j i_{pj} = \text{net}_p^o$$

Final network output

$$O_p = f(\text{net}_p^o) \quad \left(\text{this 'f' can be diff from the f used at hidden layer neurons} \right)$$

$$\begin{aligned} \text{i.e. } O_p &= f(\text{net}_p^o) \\ &= f\left(\sum_{j=1}^K \beta_j i_{pj}\right) \\ &= f\left(\sum_{j=1}^K \beta_j f\left(\sum_{i=1}^M \omega_{ij}^h x_{pi}\right)\right) \\ \text{or } &= g\left(\sum_{j=1}^K \beta_j f\left(\sum_{i=1}^M \omega_{ij}^h x_{pi}\right)\right) \end{aligned}$$

Note: We can have different 'f' at different hidden layer neurons, in that case

$$O_p = g\left(\sum_{j=1}^K \beta_j f_j\left(\sum_{i=1}^M \omega_{ij}^h x_{pi}\right)\right)$$

Note: Suppose we consider sigmoid transfer f^n at all neurons (hidden and output) then

$$\begin{aligned} O_p &= f\left(\sum_{j=1}^K \beta_j \left\{1 + \exp\left(-\sum_{i=1}^M \omega_{ij}^h x_{pi}\right)\right\}^{-1}\right) \\ &= \left\{1 + \exp\left(-\sum_{j=1}^K \beta_j \left[1 + \exp\left(-\sum_{i=1}^M \omega_{ij}^h x_{pi}\right)\right]^{-1}\right)\right\}^{-1} \\ \text{i.e. } O_p &= f^*\left(x_{p1}, \dots, x_{pM}\right) \rightarrow \text{a non-linear model.} \end{aligned}$$

Note: The network above is called a feed-forward network.

unknown parameters for an $M-K-1$ network: weight connections.

input layer \rightarrow hidden layer: $M \times K$ weight matrix

$$\begin{pmatrix} w_{11}^h & w_{12}^h & \dots & w_{1K}^h \\ - & - & - & - \\ w_{M1}^h & w_{M2}^h & \dots & w_{MK}^h \end{pmatrix}$$

hidden layer \rightarrow output node: K -dimensional weight vector
 $(\beta_1, \dots, \beta_K)'$

Total # of unknown weights: $(MK + K)$.

Learning of the network: determination of weights.

Either instantaneous mode OR batch mode. using
 Back-propagation algorithm (b-p of errors).

Consider for example learning on an instantaneous mode.

$$\mathcal{E}_p = \frac{1}{2} (y_p - o_p)^2 \quad \leftarrow \text{instantaneous sq. error for } p^{\text{th}} \text{ pattern}$$

$$\mathcal{E}_p = \frac{1}{2} (y_p - f^o(\text{net}_p^o))^2 \quad \left| \quad o_p = f^o\left(\sum_j \beta_j f_j^o(\text{net}_{pj}^h)\right)\right.$$

Gradient of \mathcal{E}_p w.r.t. output hidden to output layer wts.

$$\begin{aligned} \frac{\partial \mathcal{E}_p}{\partial \beta_j^o} &= -\frac{1}{2} (y_p - o_p) \frac{\partial f^o(\text{net}_p^o)}{\partial \beta_j^o} \\ &= -\frac{1}{2} (y_p - o_p) \frac{\partial f^o(\text{net}_p^o)}{\partial \text{net}_p^o} \cdot \frac{\partial \text{net}_p^o}{\partial \beta_j^o} \end{aligned}$$

$$\text{net}_p^o = \sum_{j=1}^K \beta_j^o i_{pj}$$

$$i_{pj} = f_j^h\left(\sum_{i=1}^M x_{pi} w_{ij}^h\right)$$

$$= -\frac{1}{2} (y_p - o_p) f^{o'}(\text{net}_p^o) i_{pj}$$

$$\text{i.e. } \frac{\partial \mathcal{E}_p}{\partial \beta_j^o} = -\frac{1}{2} (y_p - o_p) o_p (1 - o_p) i_{pj}$$

for a sigmoid transfer f^o at the
 output node

(176)

By gradient of E_p w.r.t. input layer \rightarrow hidden layer.

$$\frac{\partial E_p}{\partial w_{ij}^h} = \frac{\partial}{\partial w_{ij}^h} \frac{1}{2} \left(y_p - f^o \left(\sum_{j=1}^K \beta_j f_j \left(\sum_{i=1}^M w_{ij}^h x_{pi} \right) \right) \right)^2$$

$$\begin{aligned} o_p = f^o(\text{net}_p^o) & \quad \left| \quad = -(y_p - o_p) \frac{\partial f^o(\text{net}_p^o)}{\partial \text{net}_p^o} \cdot \frac{\partial \text{net}_p^o}{\partial i_{pj}} \cdot \frac{\partial i_{pj}}{\partial \text{net}_{pj}^h} \cdot \frac{\partial \text{net}_{pj}^h}{\partial w_{ij}^h} \right. \\ \text{net}_p^o = \sum_{j=1}^K \beta_j i_{pj} & \quad \left| \quad = -(y_p - o_p) \frac{\partial f^o(\text{net}_p^o)}{\partial \text{net}_p^o} \cdot \beta_j^o \cdot f_j'(\text{net}_{pj}^h) \cdot x_{pi} \right. \\ i_{pj} = f_j \left(\sum_{i=1}^M w_{ij}^h x_{pi} \right) & \quad \left| \quad \delta_p^o \right. \\ i_{pj} = f_j(\text{net}_{pj}^h) & \quad \left| \quad \right. \\ \text{net}_{pj}^h = \sum_{i=1}^M w_{ij}^h x_{pi} & \quad \left| \quad \right. \end{aligned}$$

Define, δ output layer delta as.

$$\delta_p^o = (y_p - o_p) f^{o'}(\text{net}_p^o)$$

and hidden layer delta as

$$\delta_{pj}^h = \delta_p^o \beta_j^o f_j'(\text{net}_{pj}^h)$$

Weight updation eqⁿs

output layer weights

$$\beta_j^o(k+1) = \beta_j^o(k) + \eta \delta_p^o i_{pj} \Big|_{\beta(k), w(k)}$$

hidden layer weights

$$w_{ij}^h(k+1) = w_{ij}^h(k) + \eta \delta_{pj}^h x_{pi} \Big|_{\beta(k+1), w(k)}$$

η : learning rate parameter

Iteration continues till stopping criterion is reached

Either loss error below threshold OR change in wt vector is lower than a threshold.

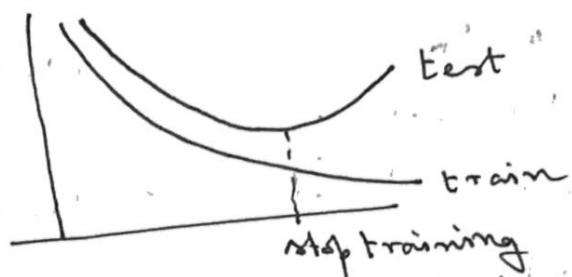
Note: Information flows from input level to output level
 → error calculated on output layer → propagated back for adjustment of hidden to output layer wts
 → with updated hidden to output layer wts error propagates backwards for weight updation of input to hidden layer weights.
 → Back-propagation of error.

Some important points: (after the discussion on Variants)

- (i) Data splitting - training set and test set
- (ii) Preprocessing - normalization of inputs and output
 Along channel & across channel
- (iii) Overfitting - • MSE on training & test set

$$\frac{x_i - \bar{x}_i}{s_{d x_i}}$$

$$\frac{x_i - \min}{\max - \min}$$



• Weight decay - add a penalty term for large weights

$$R(\theta) + \lambda J(\theta)$$

$$R(\theta) = \sum_{p=1}^n (y_p - o_p)^2 \quad \theta: \text{set of weights (MK + K)}$$

$$J(\theta) = \sum_{i,j} w_{ij}^2 + \sum_K \beta_K^2$$

λ : tuning parameter, $\lambda \geq 0$
 larger λ tend to shrink weights towards 0
 (approach similar to Ridge regression)

• Weight elimination

$$J(\theta) = \sum_{i,j} \frac{w_{ij}^2}{1 + w_{ij}^2} + \sum_K \frac{\beta_K^2}{1 + \beta_K^2}$$

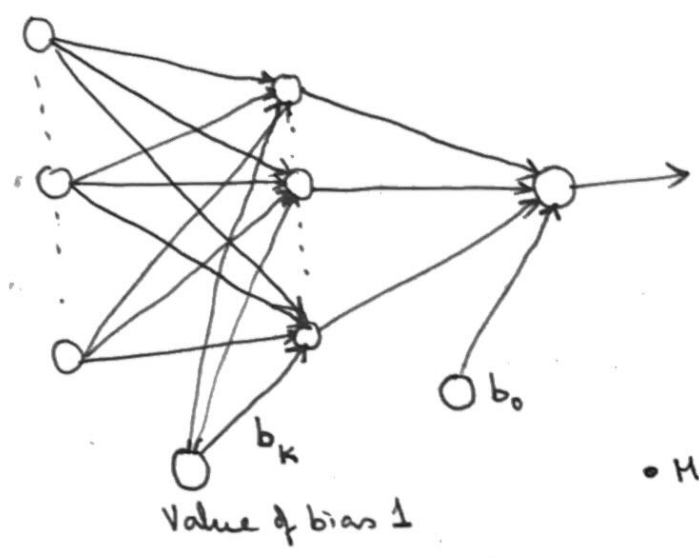
(iv) Number of hidden layers & units

One approach is to resort to experimentation

Systematic approach $\left\{ \begin{array}{l} \text{simulated annealing} \\ \text{Genetic algorithm} \end{array} \right.$

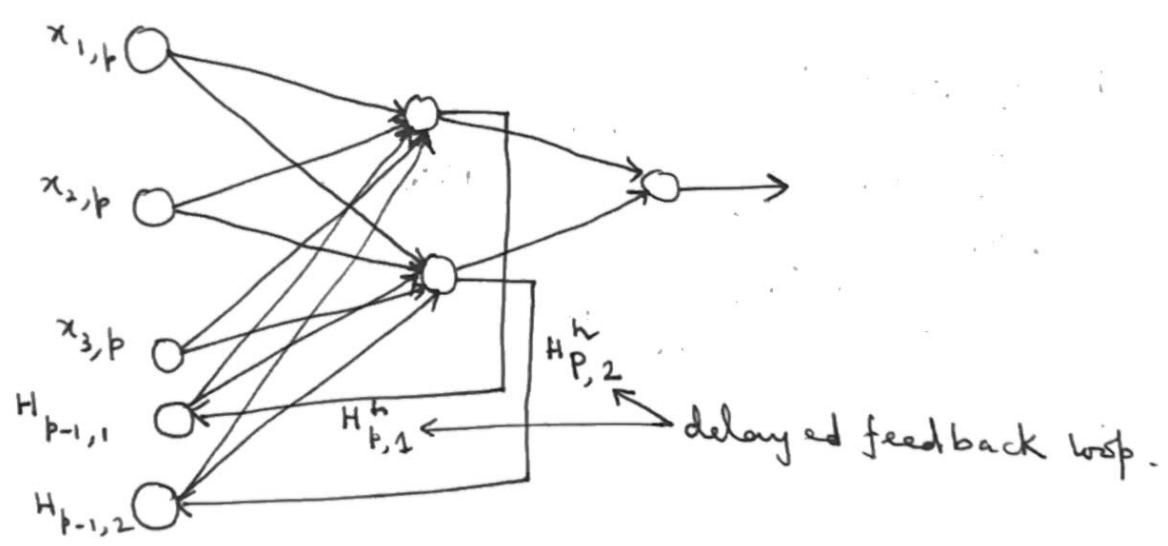
Variants of the feedforward MLP

- MLP with bias units



- Multiple output MLP.

- ~~MLP~~ Recurrent MLP (more appropriate for time series set up)
3-2-1



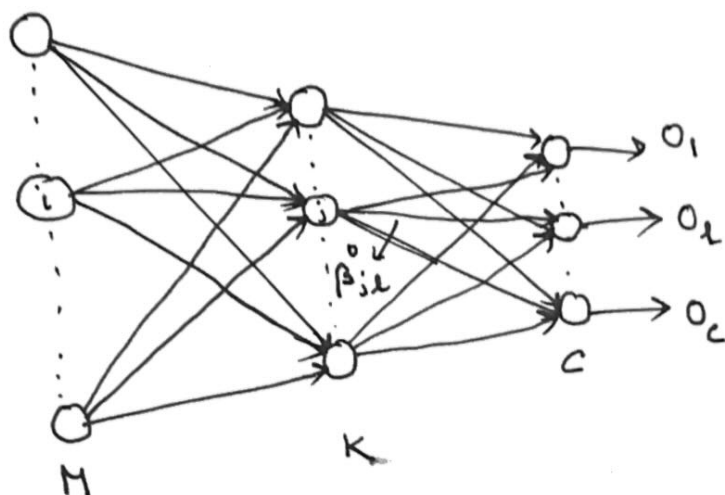
$$\hat{y}_p = \phi \left(\sum_{k=1}^2 \beta_k^o \psi_k^h \left(\sum_{i=1}^3 x_{p,i} w_{i,j}^h + \sum_{k=1}^2 w_k^f \psi_k^h \left(\sum_{i=1}^3 x_{p-1,i} w_{i,j}^h \right) \right) \right)$$

← feed back input reaching hidden units.

ANN classification model

179

c-class problem



$$\beta_{1l}^o i_{p1} + \dots + \beta_{jl}^o i_{pj} + \dots$$

$$net_{pj}^h = \sum_{i=1}^M x_{pi} w_{ij}^h$$

$$i_{pj} = f_j^h (net_{pj}^h) = f_j^h \left(\sum_{i=1}^M x_{pi} w_{ij}^h \right)$$

At l^{th} output unit

inputs are $i_{p1}, i_{p2}, \dots, i_{pK}$
 Net input at l^{th} output $\sum_{j=1}^K \beta_{jl}^o i_{pj} = T_l$
 $l = 1, \dots, c$

output of l^{th} output unit

$$O_l = \frac{e^{i_{pl} T_l}}{\sum_{j=1}^c e^{i_{pj} T_j}} ; l = 1, \dots, c$$

$$\sum_{l=1}^c O_l = 1 ; 0 \leq O_l \leq 1$$

Corresponding classifier is

$$G(x) = \underset{k}{\operatorname{argmax}} O_k$$

Total error : $\sum_{i=1}^n \sum_{k=1}^c (y_{ik} - O_k)^2$

or deviance $-\sum_{i=1}^n \sum_{k=1}^c y_{ik} \log O_k(x_i)$

(*) Discuss how training set is to be constructed.

This will specify what is

y_{ik}

Steps for ANN model building

180

- (1) Creation of Data records
- (2) Splitting of data into training & test set
- (3) Preprocessing of inputs and output
- (4) Decide network architecture - # of inputs, # of hidden layers, # of neurons in hidden layer, transfer f^{ns}
- (5) Training of the network
- (6) ~~the~~ Calculate error measure for training & test set