



Non-Bayesian classifiers.

**Linear discriminants,
neural networks.**



Discriminant functions(1)

Bayes classification rule:

$$P(\omega_1 | x) - P(\omega_2 | x) > 0? \quad \omega_1 : \omega_2$$

Instead might try to find a function:

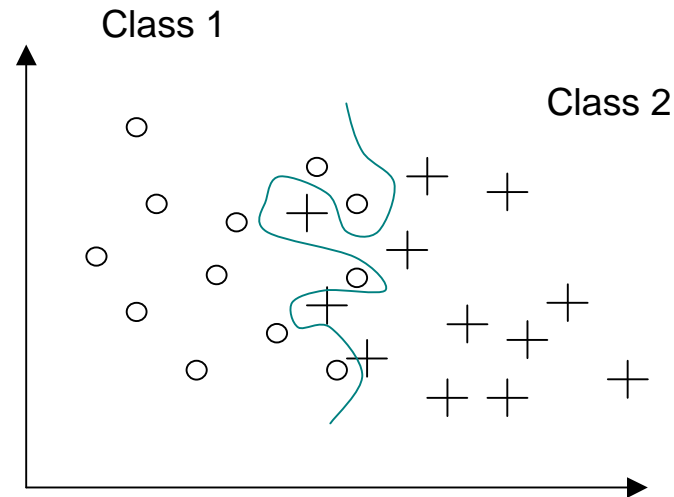
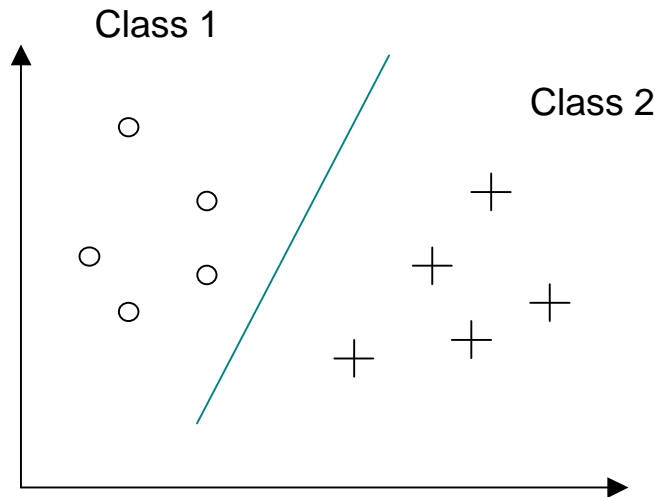
$$f_{\omega_1, \omega_2}(x) > 0? \quad \omega_1 : \omega_2$$

$f_{\omega_1, \omega_2}(x)$ is called discriminant function.

$\{x | f_{\omega_1, \omega_2}(x) = 0\}$ - decision surface



Discriminant functions (2)



Linear discriminant function:

$$f_{\omega_1, \omega_2}(x) = w^T x + w_0$$

Decision surface is a hyperplane $w^T x + w_0 = 0$



Linear discriminant – perceptron cost function

Replace $w \leftarrow \begin{bmatrix} w \\ w_0 \end{bmatrix}$ and $x \leftarrow \begin{bmatrix} x \\ 1 \end{bmatrix}$

Thus now decision function is $f_{\omega_1, \omega_2}(x) = w^T x$
 and decision surface is $w^T x = 0$

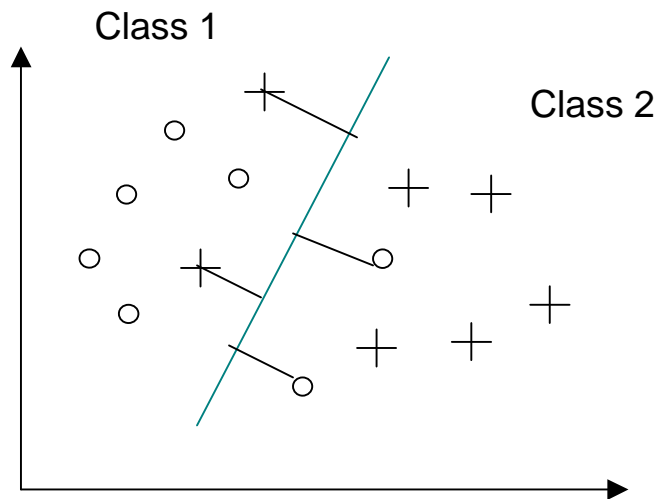
Perceptron cost function:

$$J(x) = \sum_x \delta_x w^T x$$

where $\delta_x = \begin{cases} -1, & \text{if } x \in \omega_1 \text{ and } x \text{ is } w^T x < 0 \\ 1, & \text{if } x \in \omega_2 \text{ and } x \text{ is } w^T x > 0 \\ 0, & x \text{ is correctly classified} \end{cases}$



Linear discriminant – perceptron cost function



Perceptron cost function:

$$J(x) = \sum_x \delta_x w^T x$$

Value of $J(x)$ is proportional to the sum of distances of all misclassified samples to the decision surface.

If discriminant function separates classes perfectly, then $J(x) = 0$
 Otherwise, $J(x) > 0$ and we want to minimize it.

$J(x)$ is continuous and piecewise linear. So we might try to use gradient descent algorithm.



Linear discriminant – Perceptron algorithm

Gradient descent:

$$w(t+1) = w(t) - \rho_t \left. \frac{\partial J(w)}{\partial w} \right|_{w=w(t)}$$

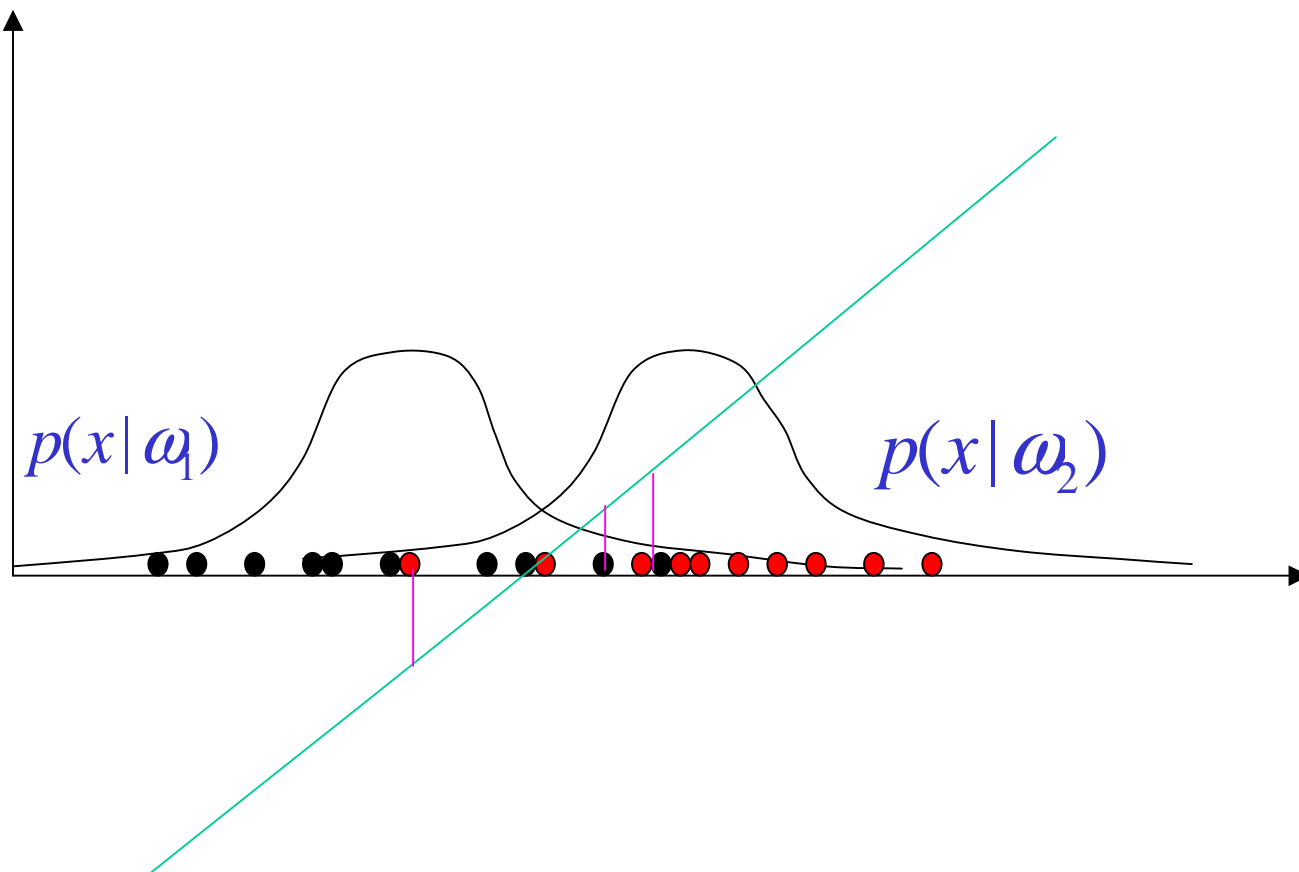
At points where $J(x)$ is differentiable $\frac{\partial J(w)}{\partial w} = \sum_{\substack{\text{misclassified} \\ x}} \delta_x x$

Thus $w(t+1) = w(t) - \rho_t \sum_{\substack{\text{misclassified} \\ x}} \delta_x x$

Perceptron algorithm converges when classes are linearly separable with some conditions on ρ_t



Linear discriminant with perceptron cost





Sum of error squares estimation

Let denote $y(\mathbf{x}) = \pm 1$ as desired output function, 1 for one class and -1 for the other.

Want to find discriminant function $f_{\omega_1, \omega_2}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ whose output is similar to $y(\mathbf{x})$

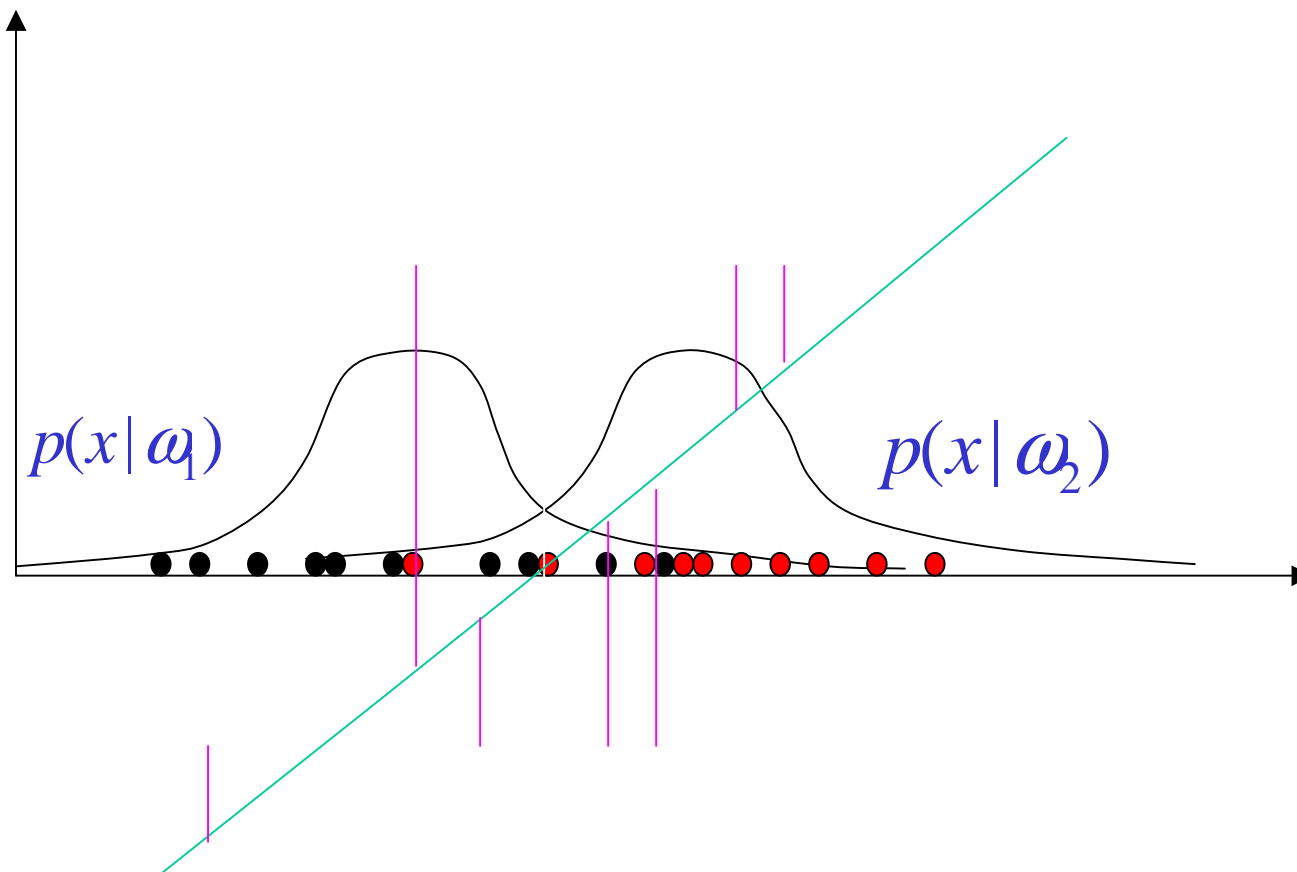
Use sum of error squares as similarity criterion:

$$J(\mathbf{w}) = \sum_{i=1}^N |y_i - \mathbf{w}^T \mathbf{x}_i|^2$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} J(\mathbf{w})$$



Linear discriminant with sum of error squares





Sum of error squares estimation

Minimize mean square error:

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 2 \sum_{i=1}^N \mathbf{x}_i (y_i - \mathbf{w}^T \mathbf{x}_i) = 0$$

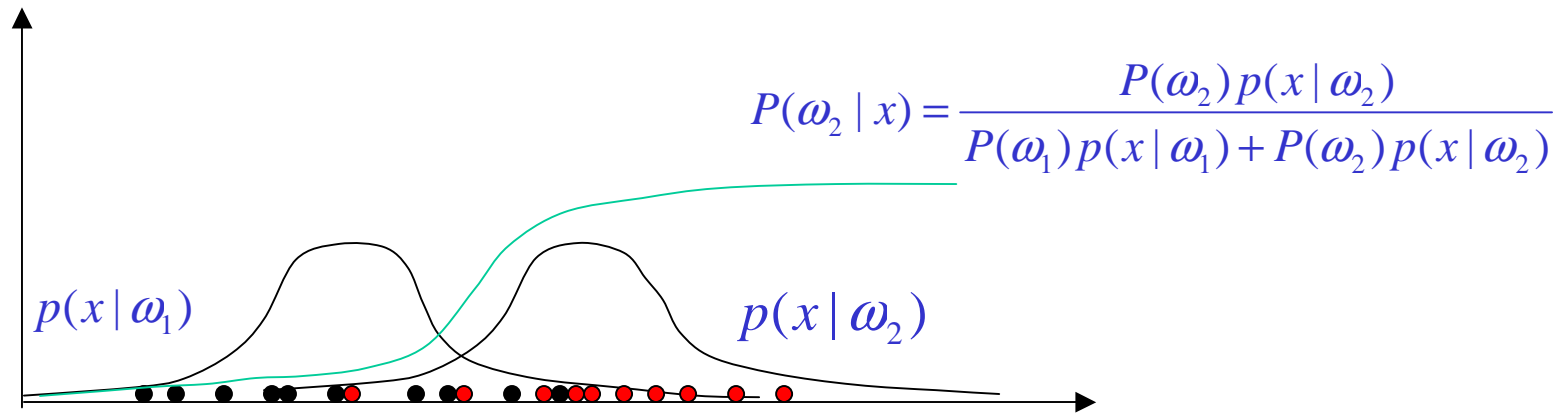
$$\left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right) \hat{\mathbf{w}} = \sum_{i=1}^N \mathbf{x}_i y_i$$

Thus

$$\hat{\mathbf{w}} = \left(\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \left(\sum_{i=1}^N \mathbf{x}_i y_i \right)$$



Logistic function discrimination

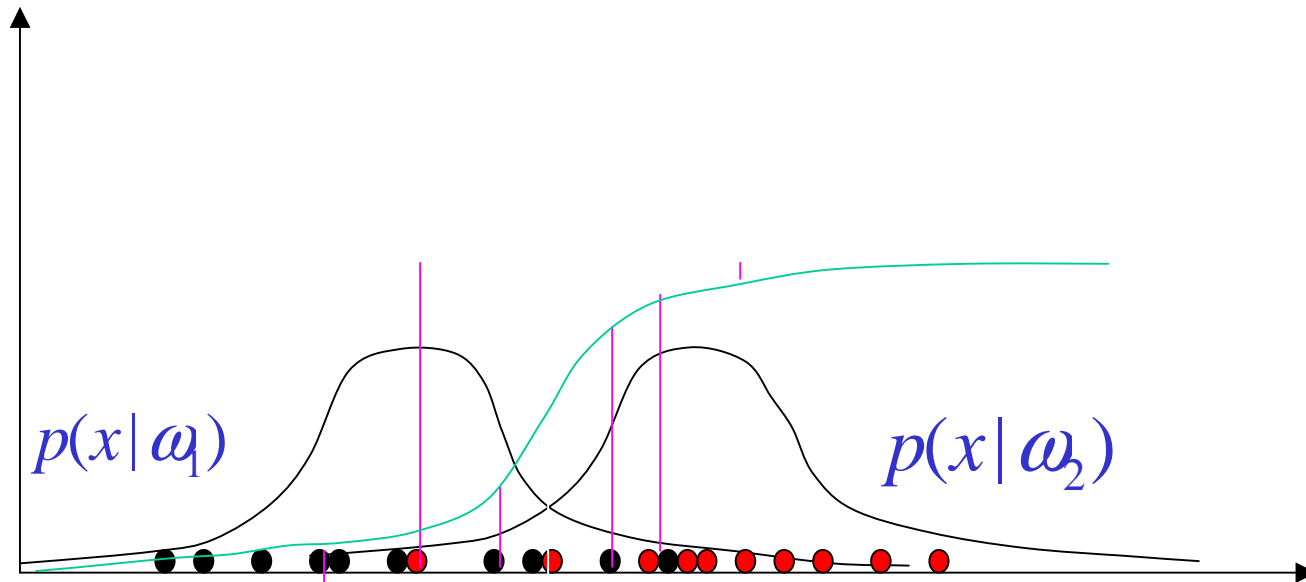


If class densities are gaussians, then optimal discrimination function (posterior class probability) is logistic:

$$P(\omega_2 | x) = \frac{P(\omega_2)C_2 e^{-\frac{(x-\mu_2)^2}{\sigma^2}}}{P(\omega_1)C_1 e^{-\frac{(x-\mu_1)^2}{\sigma^2}} + P(\omega_2)C_2 e^{-\frac{(x-\mu_2)^2}{\sigma^2}}} = \frac{1}{\frac{P(\omega_1)C_1}{P(\omega_2)C_2} e^{-\frac{(x-\mu_1)^2}{\sigma^2} + \frac{(x-\mu_2)^2}{\sigma^2}} + 1} = \frac{1}{C_3 e^{-ax} + 1}$$

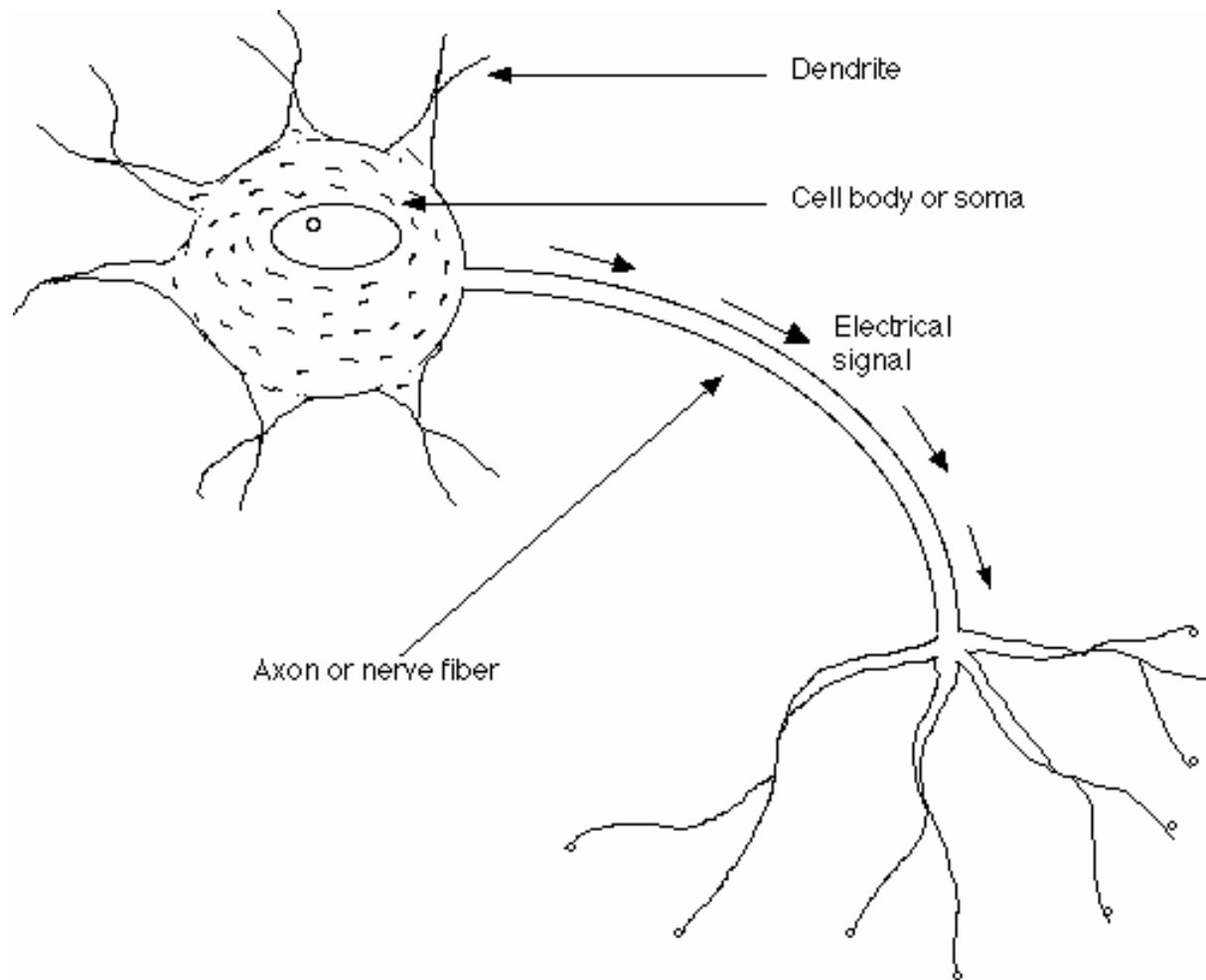


Logistic discriminant with sum of error squares



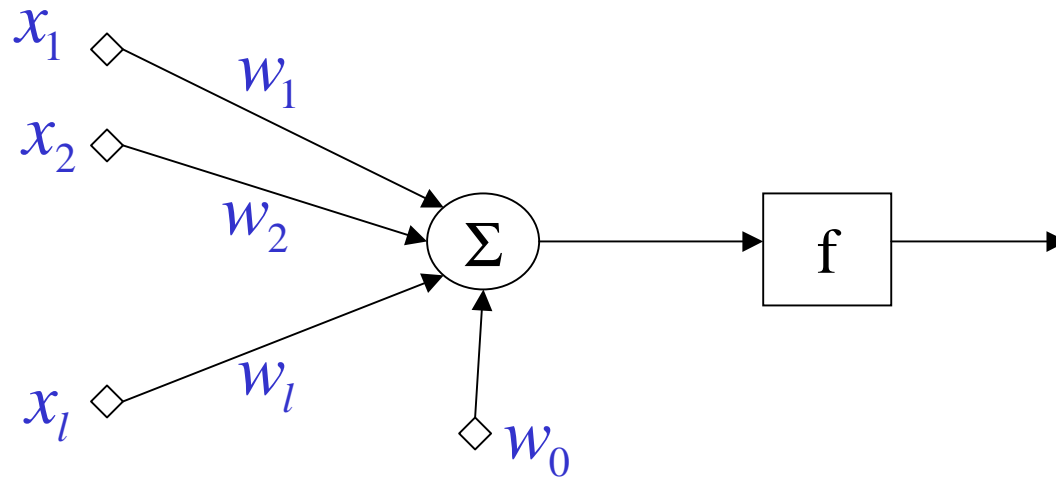


Neurons





Artificial neuron.



Above figure represent artificial neuron calculating:

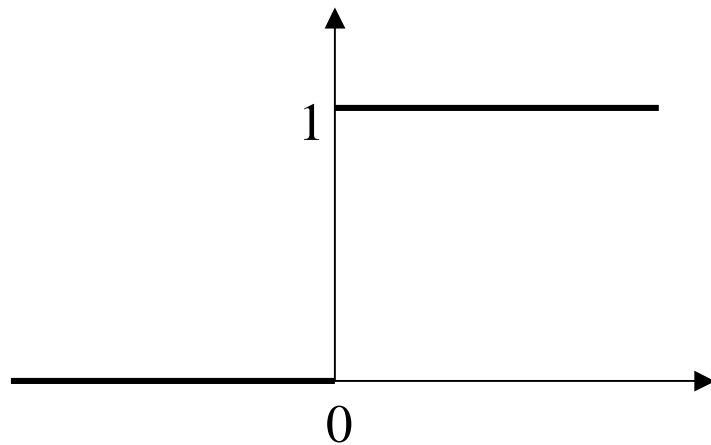
$$y = f\left(\sum_{i=1}^l w_i x_i\right)$$



Artificial neuron.

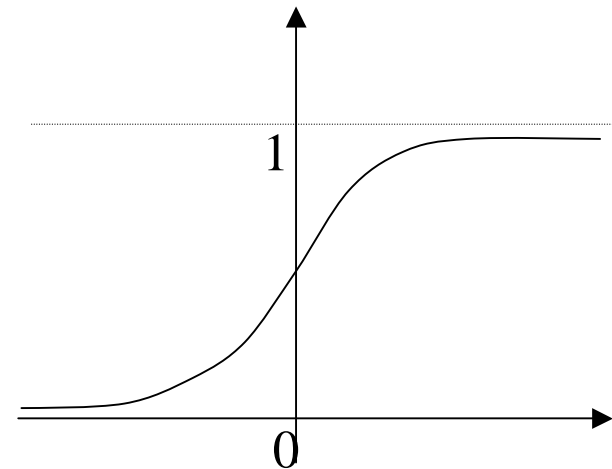
Threshold functions f :

Step function



$$f(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

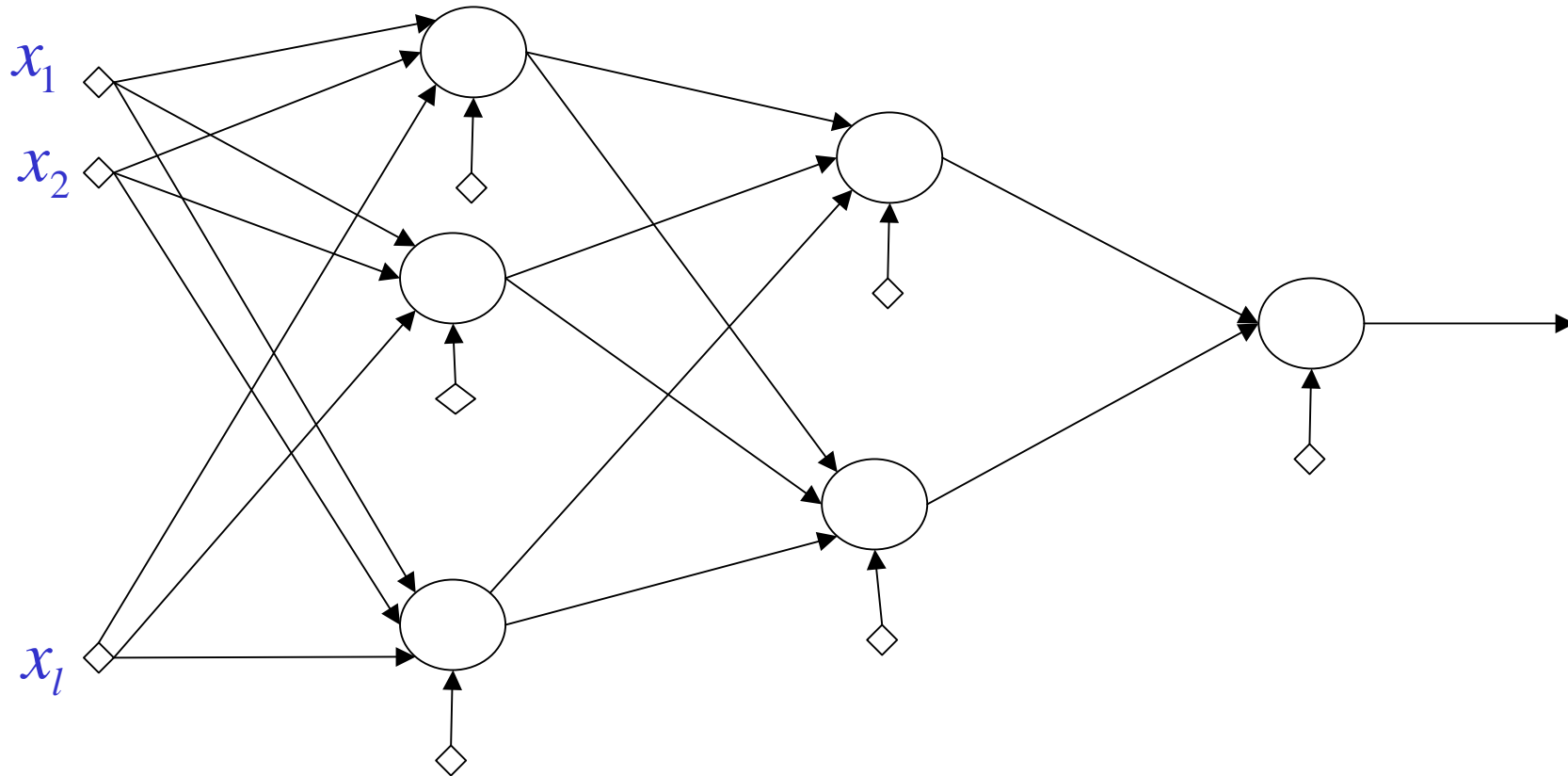
Logistic function



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



Combining artificial neurons



Multilayer perceptron with 3 layers.

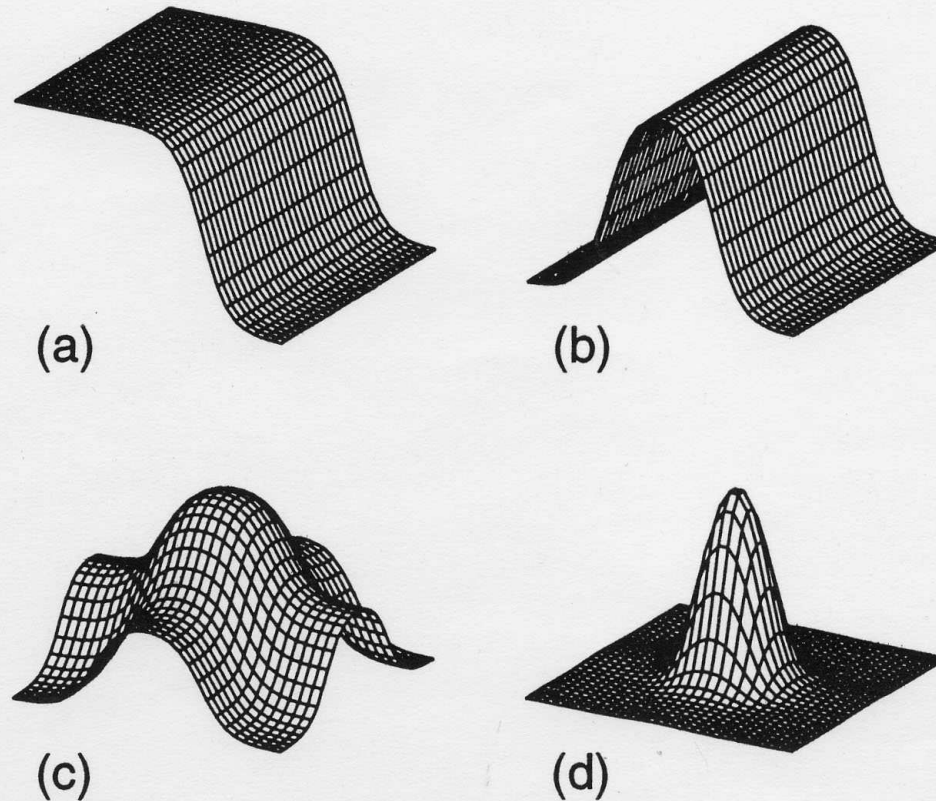
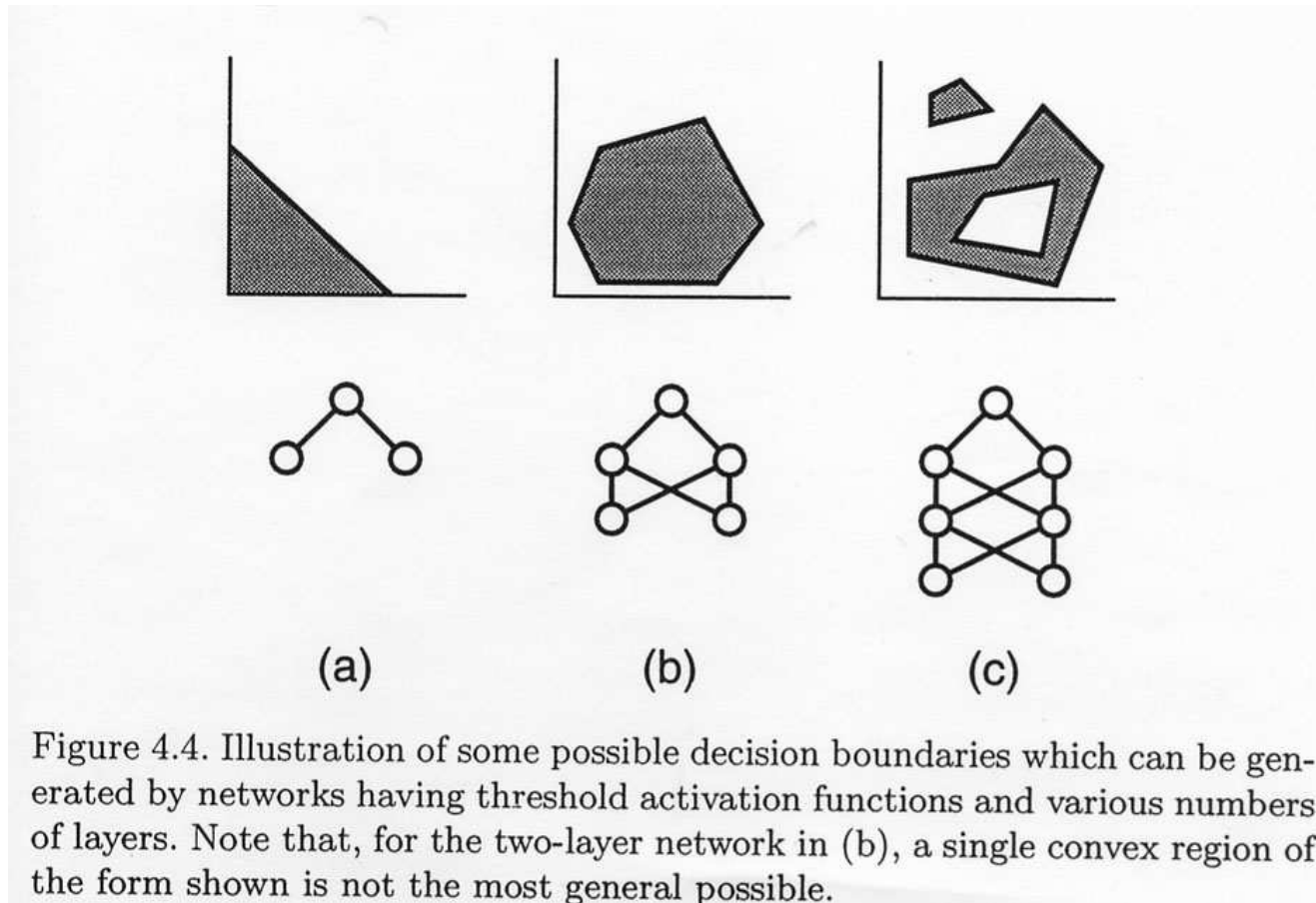


Figure 4.10. Demonstration that a network with three layers of weights, and sigmoidal hidden units, can approximate a smooth multivariate mapping to arbitrary accuracy. In (a) we see the output of a single sigmoidal unit as a function of two input variables. Adding the outputs from two such units can produce a ridge-like function (b), and adding two ridges can give a function with a maximum (c). Transforming this function with another sigmoid gives a localized response (d). By taking linear combinations of these localized functions, we can approximate any smooth functional mapping.



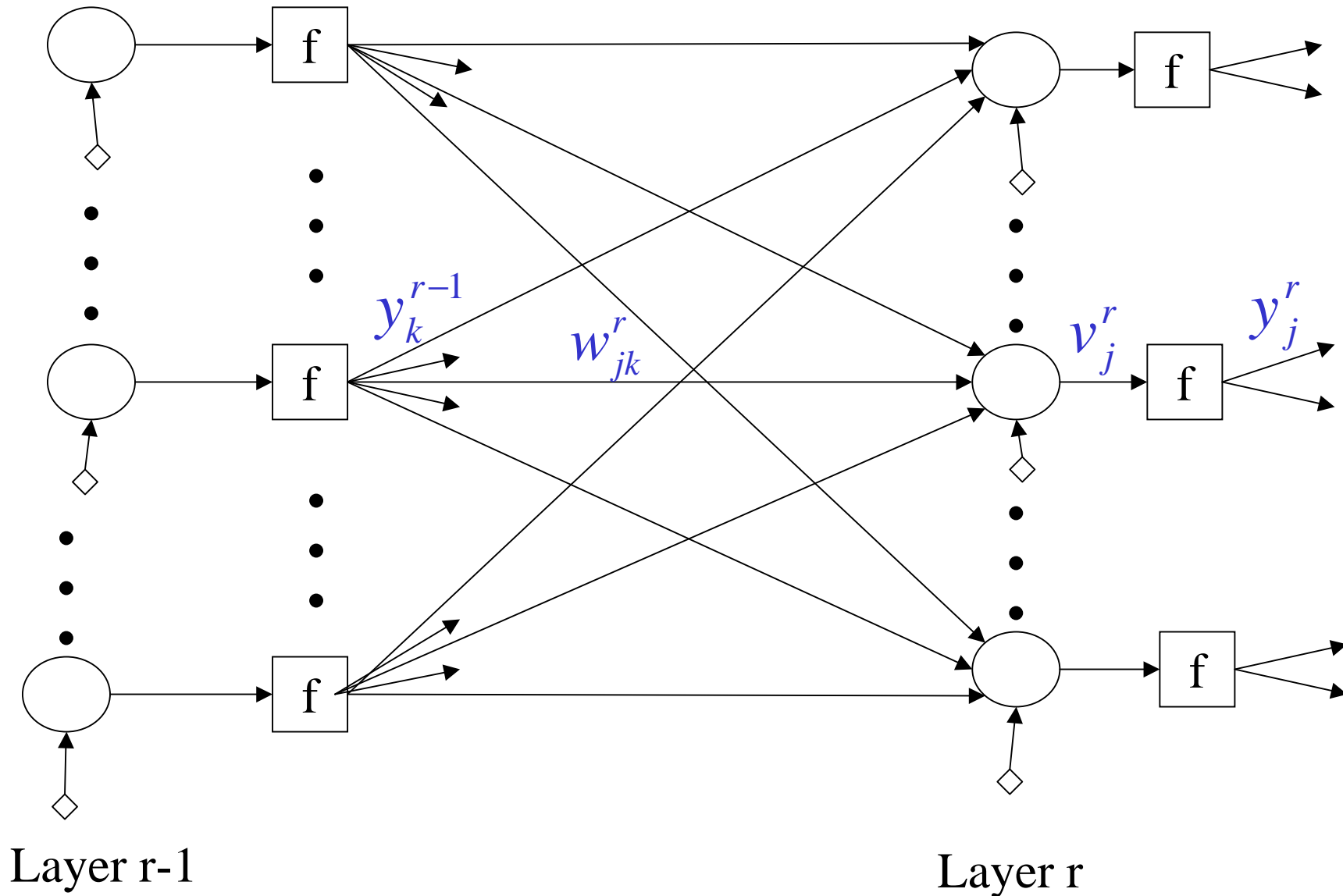
Discriminating ability of multilayer perceptron

Since 3-layer perceptron can approximate any smooth function, it can approximate $F(x) = P(\omega_1 | x) - P(\omega_2 | x)$ - optimal discriminant function of two classes.





Training of multilayer perceptron





Training and cost function

Desired network output: $x(i) \rightarrow y(i)$

Trained network output: $x(i) \rightarrow \hat{y}(i)$

Cost function for one training sample:

$$E(i) = \frac{1}{2} \sum_{m=1}^{k_L} (y_m(i) - \hat{y}_m(i))^2$$

Total cost function: $J = \sum_{i=1}^N E(i)$

Goal of the training: find values of w_{jk}^r which minimize cost function J .



Gradient descent

Denote: $\mathbf{w}_j^r = [w_{j0}^r, w_{j1}^r, \dots, w_{jk_{r-1}}^r]^T$

Gradient descent: $\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \frac{\partial J}{\partial \mathbf{w}_j^r}$

Since $J = \sum_{i=1}^N E(i)$, we might want to update weights after processing each training sample separately:

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \frac{\partial E(i)}{\partial \mathbf{w}_j^r}$$



Gradient descent

Chain rule for differentiating composite functions:

$$\frac{\partial E(i)}{\partial \mathbf{w}_j^r} = \frac{\partial E(i)}{\partial v_j^r(i)} \frac{\partial v_j^r(i)}{\partial \mathbf{w}_j^r} = \frac{\partial E(i)}{\partial v_j^r(i)} y^{r-1}(i)$$

Denote: $\delta_j^r(i) = \frac{\partial E(i)}{\partial v_j^r(i)}$



Backpropagation

If $r=L$, then

$$\begin{aligned}
 \delta_j^L(i) &= \frac{\partial E(i)}{\partial v_j^L(i)} = \frac{\partial}{\partial v_j^L(i)} \left(\frac{1}{2} \sum_{m=1}^{k_L} (f(v_m^L(i)) - \hat{y}_m(i))^2 \right) \\
 &= (f(v_j^L(i)) - \hat{y}_j(i)) f'(v_j^L(i)) = e_j(i) f'(v_j^L(i))
 \end{aligned}$$

If $r < L$, then

$$\begin{aligned}
 \delta_j^{r-1}(i) &= \frac{\partial E(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^{k_r} \frac{\partial E(i)}{\partial v_k^r(i)} \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} \\
 &= \sum_{k=1}^{k_r} \delta_k^r(i) \frac{\partial v_k^r(i)}{\partial v_j^{r-1}(i)} = \sum_{k=1}^{k_r} \delta_k^r(i) w_{kj}^r f'(v_j^{r-1}(i))
 \end{aligned}$$



Backpropagation algorithm

- Initialization: initialize all weights with random values.
- Forward computations: for each training vector $x(i)$ compute all $v_j^r(i)$, $y_j^r(i)$
- Backward computations: for each i , j and $r=L, L-1, \dots, 2$ compute $\delta_j^{r-1}(i)$
- Update weights:

$$\mathbf{w}_j^r(\text{new}) = \mathbf{w}_j^r(\text{old}) - \mu \frac{\partial E(i)}{\partial \mathbf{w}_j^r}$$

$$= \mathbf{w}_j^r(\text{old}) - \mu \delta_j^r(i) y^{r-1}(i)$$



MLP issues

- What is the best network configuration?
- How to choose proper learning parameter μ ?
- When training should be stopped?
- Choose another threshold function f or cost function J ?