

# Levenberg-Marquardt learning

- Multi-layer neural networks
- MLP (Multilayer Perceptrons)
- RBF (Radial Basis Functions)

# Taylor expansion

$$E(\theta + \Delta\theta)$$

$$\approx E(\theta) + E'(\theta)^T \Delta\theta$$

$$E(\theta + \Delta\theta)$$

$$= E(\theta) + E'(\theta)^T \Delta\theta +$$

$$\frac{1}{2} \Delta\theta^T E''(\theta) \Delta\theta$$

# Function Approximation

- High-dimensional nonlinear function approximation
- Linear combination of basis functions
  - Projective basis function
  - Radial basis function

# Function approximation

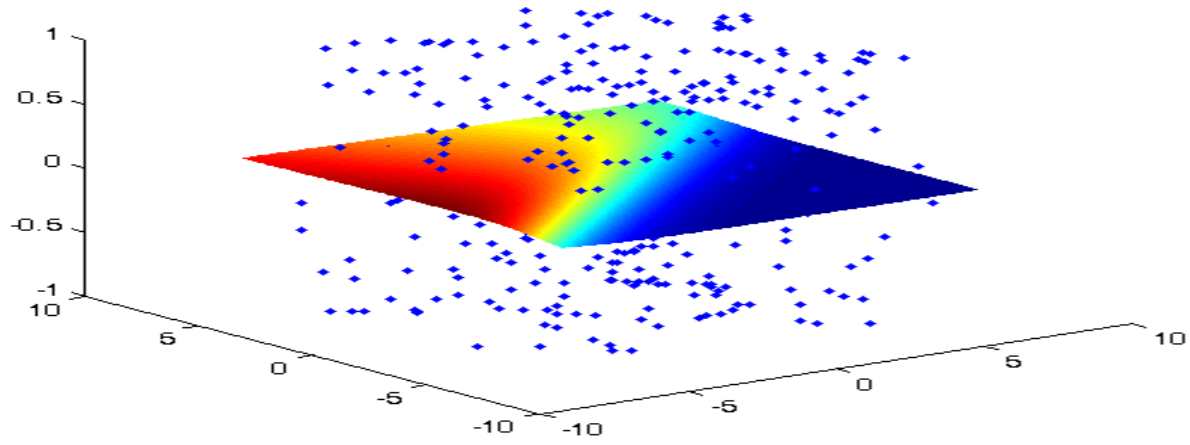
- Data oriented function approximation
- Paired data sampled from an unknown

target function :  $\mathbb{R}^d \longrightarrow \mathbb{R}$

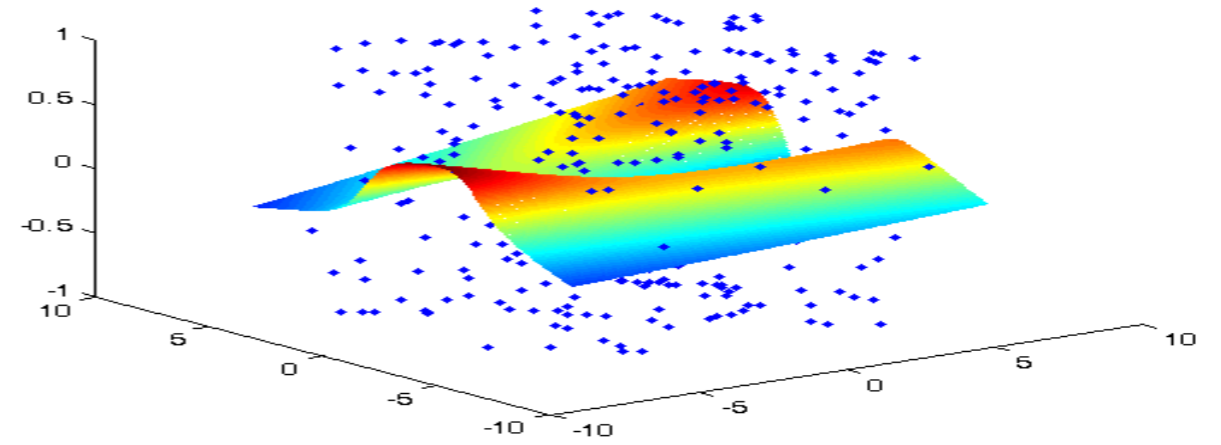
$(\mathbf{x}[t], y[t]), t=1, \dots, N$



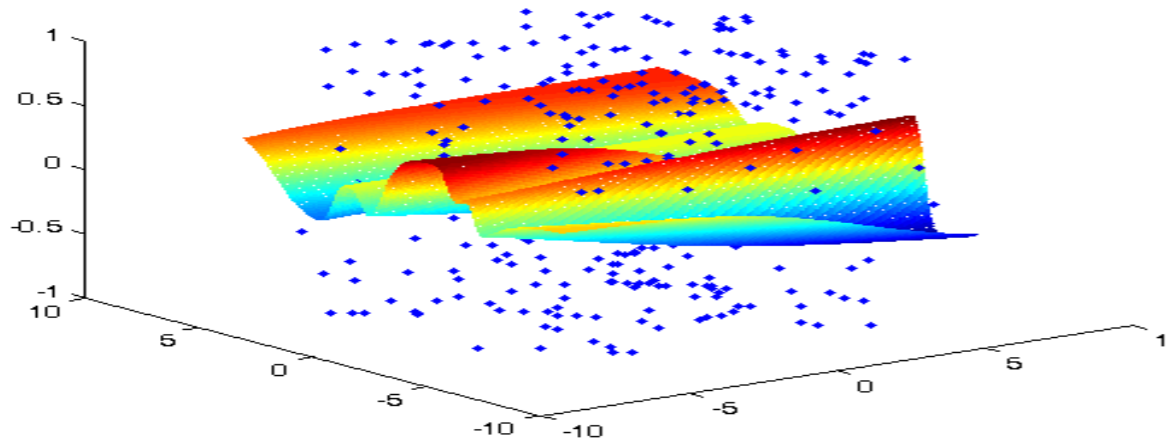
NRBF(3) by annealed FE



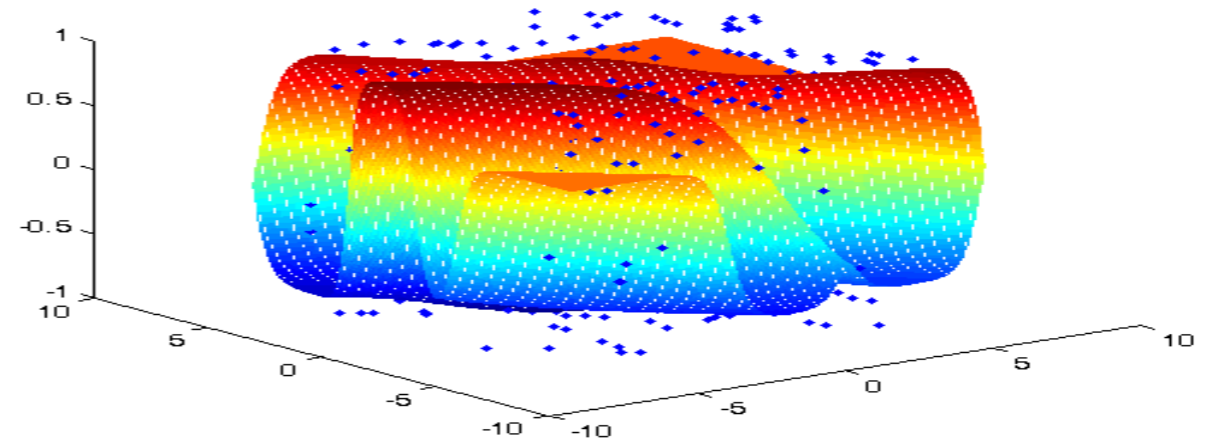
NRBF(6) by annealed FE



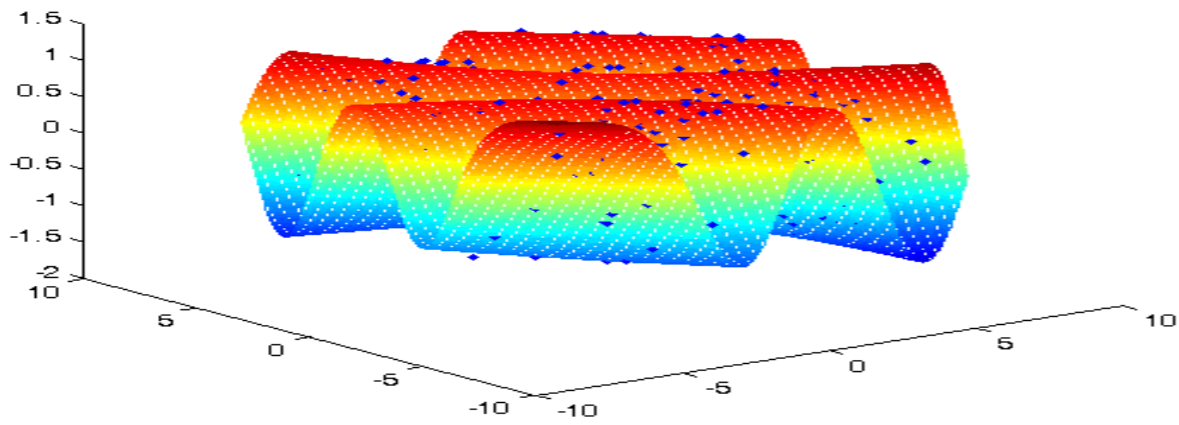
NRBF(9) by annealed FE



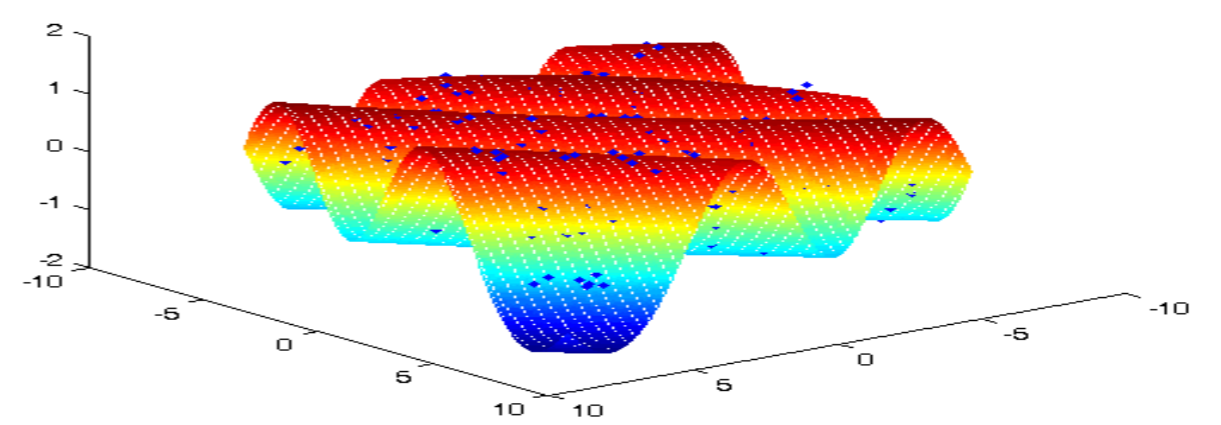
NRBF(12) by annealed FE



NRBF(15) by annealed FE



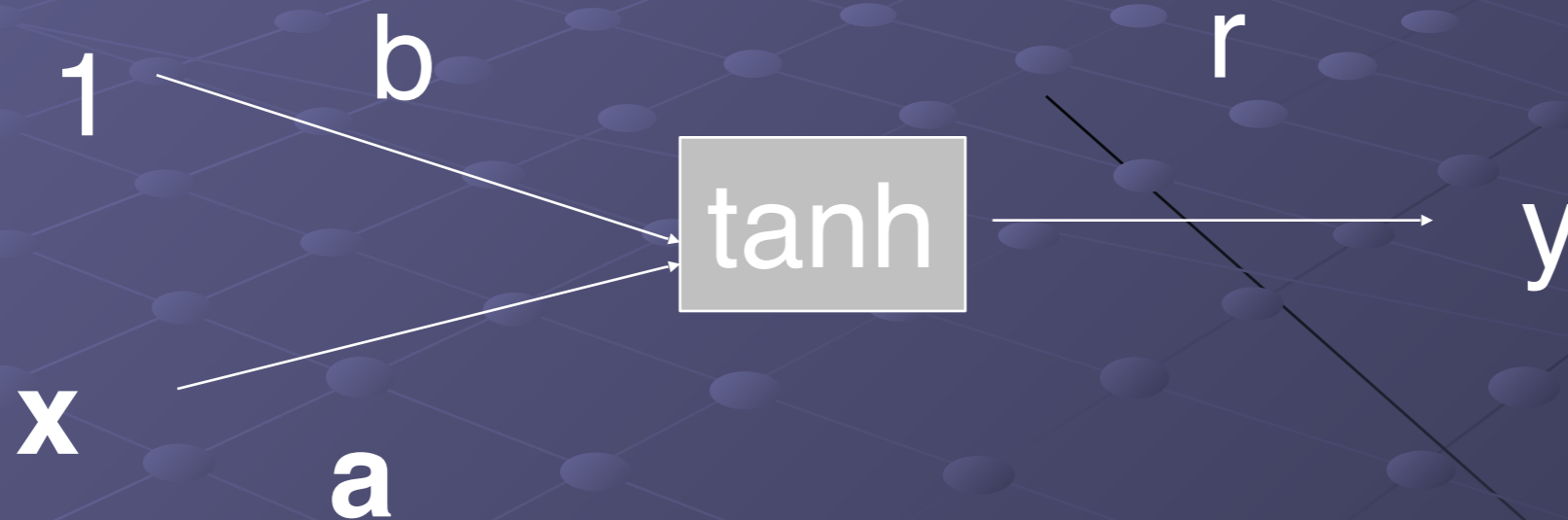
NRBF(18) by annealed FE



# Perceptron

## Projective basis function

$$y = r \tanh(\mathbf{a}^T \mathbf{x} + b)$$



# Post-tanh projection

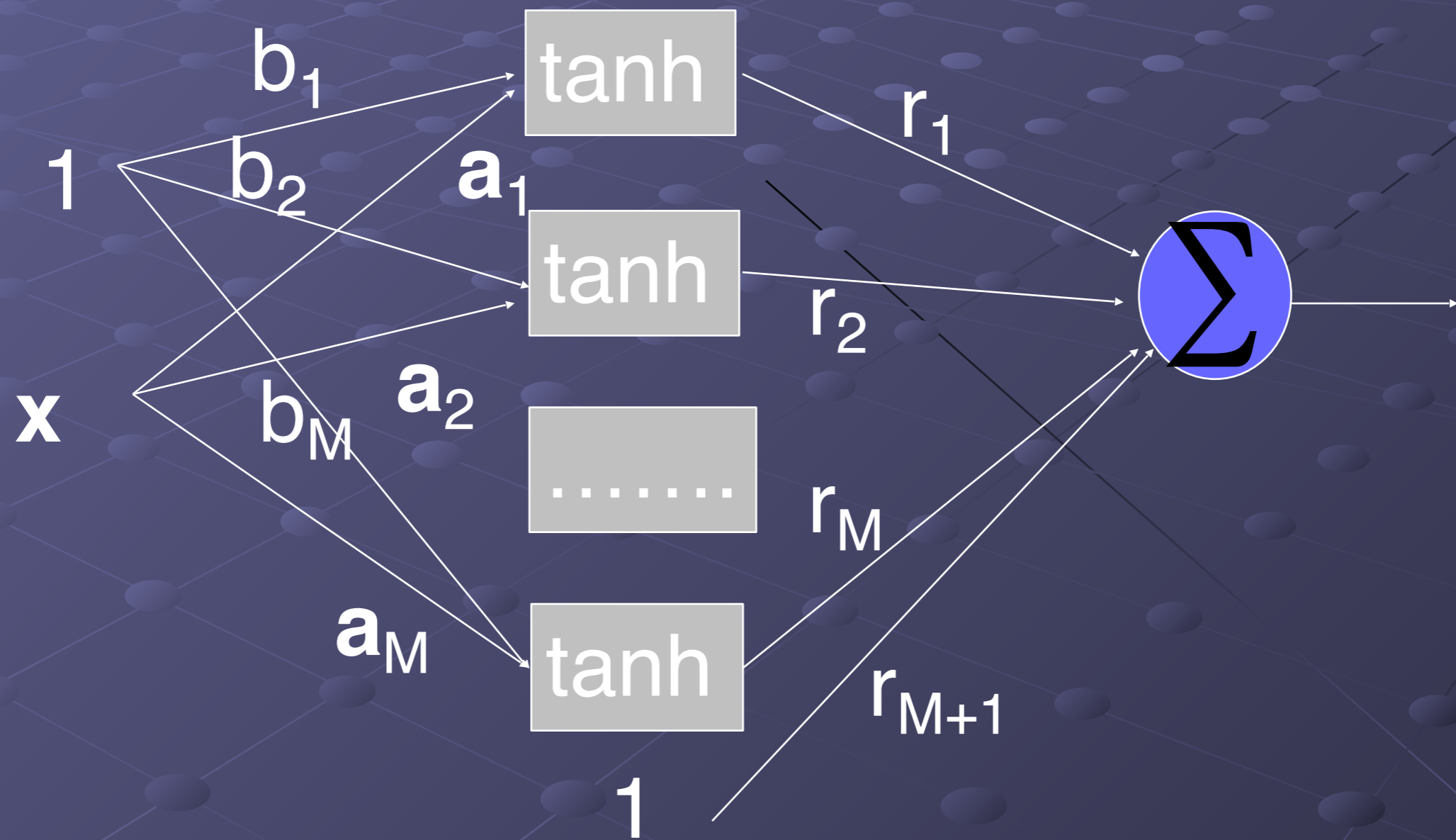
Posterior Weight

$$y = r \tanh(\mathbf{a}^T \mathbf{x} + b)$$

Linear projection

Post-nonlinear function

# Multiple (Multi-layer) perceptrons





# MLP (multilayer perceptrons)

● Weight sum of post-tanh projections

$$F(\mathbf{x};\theta) = \sum_{m=1}^M r_m \tanh(\mathbf{a}_m^T \mathbf{x} + b_m) + r_0$$

$$\theta = \{\mathbf{a}_m\} \cup \{b_m\} \cup \{r_m\}$$

# Why Projection?

## ● Why projection?

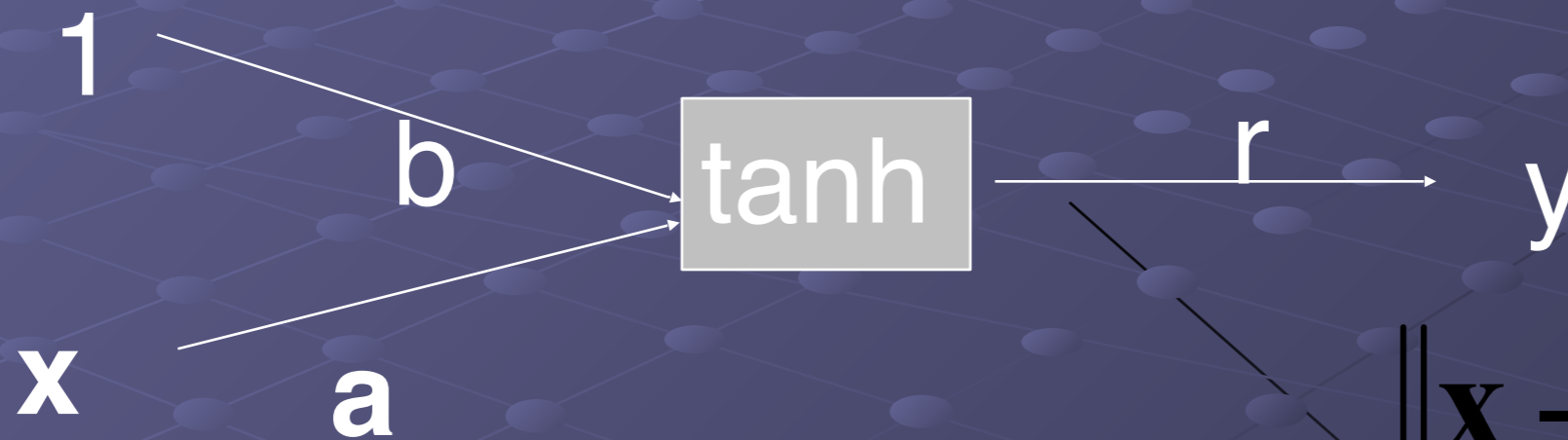
- Linear projection versus Radial basis function

## ● Why hyper-tangent

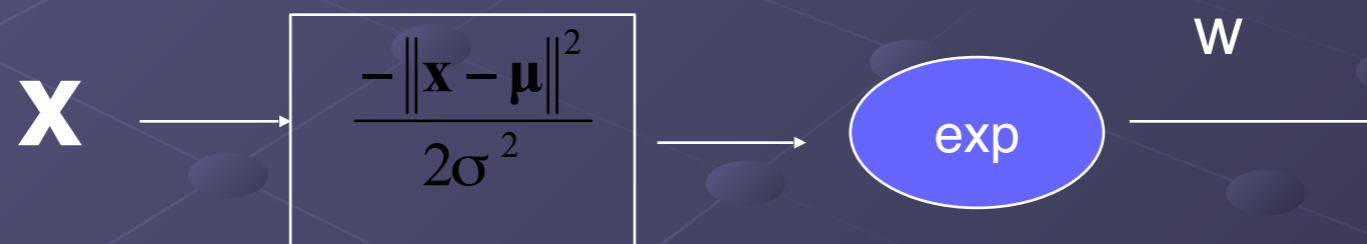
- Hyper-tangent function versus exponential function

# Radial basis function

$$y = r \tanh(\mathbf{a}^T \mathbf{x} + b)$$



$$y = w \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}\right)$$



# RBF Network function

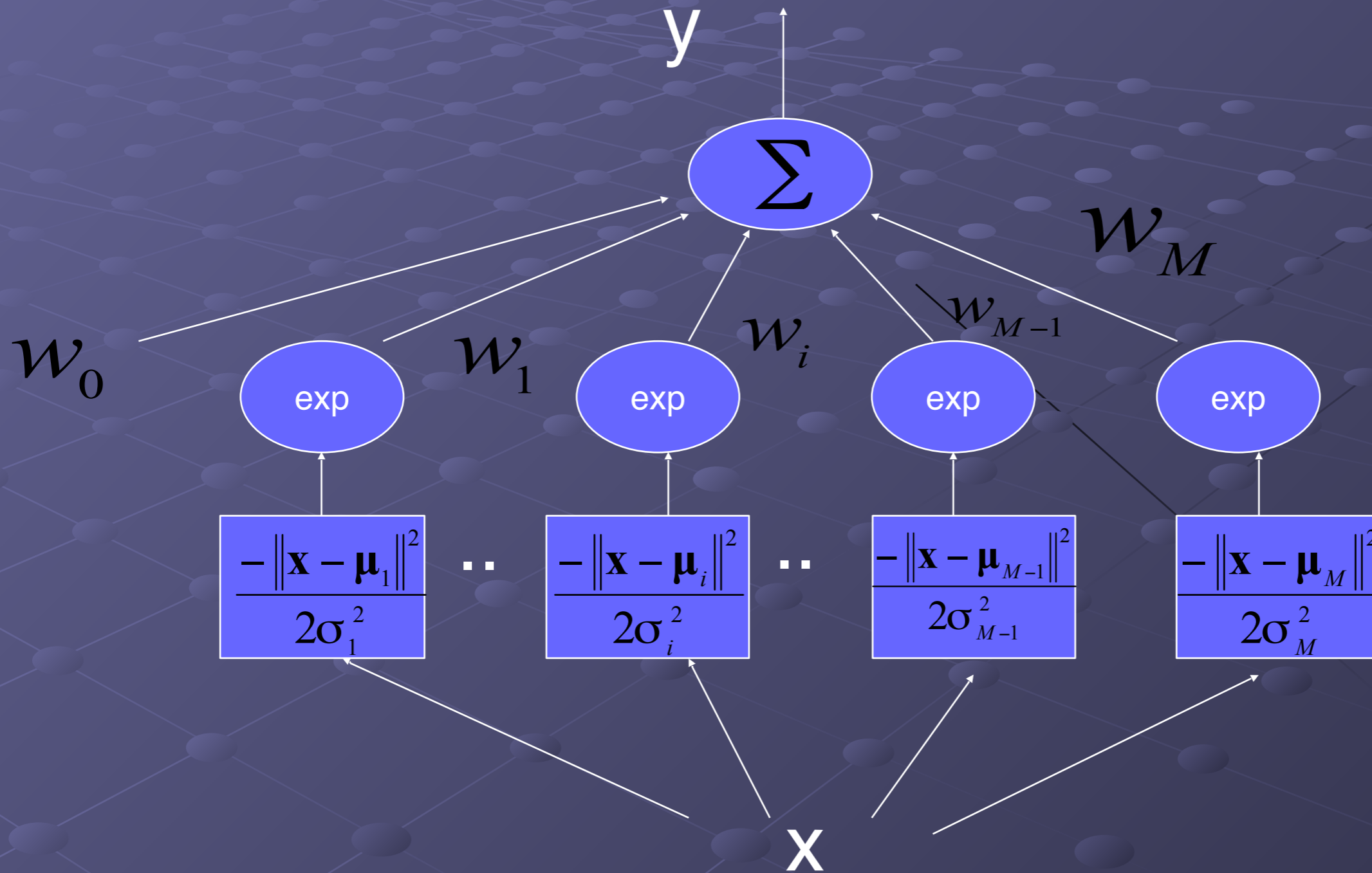
$$y(t | \theta) = G(\mathbf{x}[t] | \theta)$$

$$= w_0 + \sum_{m=1}^M w_m \exp\left(-\frac{\|\mathbf{x}[t] - \boldsymbol{\mu}_m\|^2}{2\sigma_m^2}\right)$$

*Network* parameter

$$\theta = \{\mathbf{w}_i\}_i \cup \{\boldsymbol{\mu}_i\}_i \cup \{\sigma_i\}_i$$

# RBF Network



# Posterior weights

- Two stages
- Calculation of radial basis functions
- A linear combination

Stage I

Radial basis functions

$$f_m(x | \mu_m, \delta_m)$$

$$= \exp\left(-\frac{\|x - \mu_m\|^2}{2\delta_m^2}\right)$$

$$O[t] = (O_1[t], \dots, O_M[t])^T$$

$$O_m[t] = f_m(x[t])$$

Stage II

Posterior weighting

$$y[t] = r^T o[t] + v_0$$

-----

$$x[t] \rightarrow o[t] \rightarrow y[t]$$

Given  $\{u_m\}, \{\delta_m\}$

Optimize  $r$



$$D = \begin{bmatrix} 0 & [1]^T \\ 0 & [2]^T \\ \vdots & \\ 0 & [N]^T \end{bmatrix}, \quad N \times d$$

$$D = [0 \quad \text{ones}(N, 1)]$$

$$Dy = y$$

$$Or = y$$

$$O^T Or = O^T y$$

$$C(r) = r^T O^T Or - O^T y$$

$C(r)$

$$= r^T O r - r^T y + \lambda r^T r$$

Regularization

$$\frac{dc}{dr}$$

$$= O^T O r - O^T y + \lambda r = 0$$

$$(O^T O + \lambda I) r = O^T y$$

$$Y = \text{pinv}(0 + \lambda I) y$$


optimization of

$Y$  for given

$\{N_m\}$  and  $\{b_m\}$

# Ratsch Algorithm

## Algorithm RBF-Net( $K, \lambda, O$ )

### Input:

Sequence of labeled training patterns  $\mathbf{Z} = \langle (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l) \rangle$

Number of RBF centers  $K$

Regularization constant  $\lambda$

Number of iterations  $O$

### Initialize:

Run  $K$ -means clustering to find initial values for  $\mu_k$  and determine  $\sigma_k, k = 1, \dots, K$ , as the distance between  $\mu_k$  and the closest  $\mu_i$  ( $i \neq k$ ).

### Do for $o = 1 : O$ ,

1. Compute optimal output weights  $\mathbf{w} = (G^T G + 2\frac{\lambda}{l} \mathbf{I})^{-1} G^T \mathbf{y}$
- 2a. Compute gradients  $\frac{\partial}{\partial \mu_k} E$  and  $\frac{\partial}{\partial \sigma_k} E$  as in (28) and (27) with optimal  $\mathbf{w}$  and form a gradient vector  $\mathbf{v}$
- 2b. Estimate the conjugate direction  $\bar{\mathbf{v}}$  with Fletcher-Reeves-Polak-Ribiere CG-Method (Press et al., 1992)
- 3a. Perform a line search to find the minimizing step size  $\delta$  in direction  $\bar{\mathbf{v}}$ ; in each evaluation of  $E$  recompute the optimal output weights  $\mathbf{w}$  as in line 1
- 3b. Update  $\mu_k$  and  $\sigma_k$  with  $\bar{\mathbf{v}}$  and  $\delta$

### Output: Optimized RBF net

# K Means

- Given  $\{x[t]\}$ , find centers and variances
- Matlab [kmeans.m](#)

[http://134.208.26.59/MathProgramming2010/Lecture10/  
Lecture10II.pdf](http://134.208.26.59/MathProgramming2010/Lecture10/Lecture10II.pdf)

$\{x[t], y[t]\}$   
 $M, \lambda$

$\mu = \text{kmeans}(x, M)$   
 $\delta = \text{kvar}(x, \mu)$

$O = g(x | \mu, \delta)$   
 $\Gamma = \text{inv}(O^T O + \lambda I) O^T$   
 $v = \begin{bmatrix} -\frac{dE}{d\mu} & -\frac{dE}{d\delta} \end{bmatrix}$

$p = v$   
 $\delta_{\text{new}} = v^T v$

$k = 0$   
 $n = \text{size}([\mu \delta])$

not HC

EXIT

$\delta = \text{linmin}(\mu, \delta, r, v)$   
 $[\mu, \delta] = [\mu, \delta] + \delta \times$   
 $P$

$O = g(x | \mu, \delta)$   
 $\Gamma = \text{inv}(O^T O + \lambda I) O^T$   
 $v = \begin{bmatrix} -\frac{dE}{d\mu} & -\frac{dE}{d\delta} \end{bmatrix}$

$\delta_{\text{old}} = \delta_{\text{new}}$   
 $\delta_{\text{new}} = v^T v$   
 $\beta = \delta_{\text{new}} / \delta_{\text{old}}$   
 $P = v + \beta P$   
 $k = k + 1$

if  $k = n \mid v^T p \leq 0$   
 $p = v$   
 $k = 0$



# Relation between RBF and MLP

- Why Euclidean distance?
- Why exponential function?
- Why hyper-tangent function?

# Network parameter

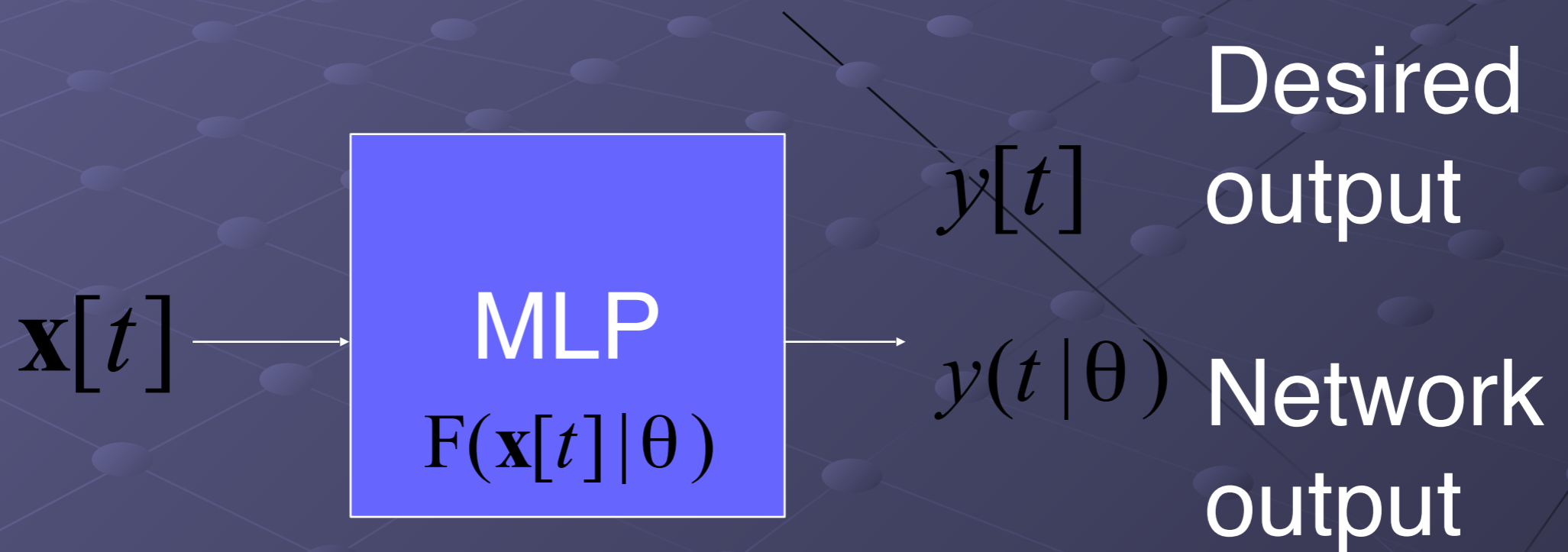
a vector of  $d$  elements

scalar

$$\theta = [\mathbf{u}_1^T \mathbf{u}_2^T \cdots \mathbf{u}_M^T \sigma_1 \sigma_2 \cdots \sigma_M w_0 w_1 w_2 \cdots w_M]^T$$
$$= [\theta_1, \dots, \theta_{M*d+2M+1}]$$

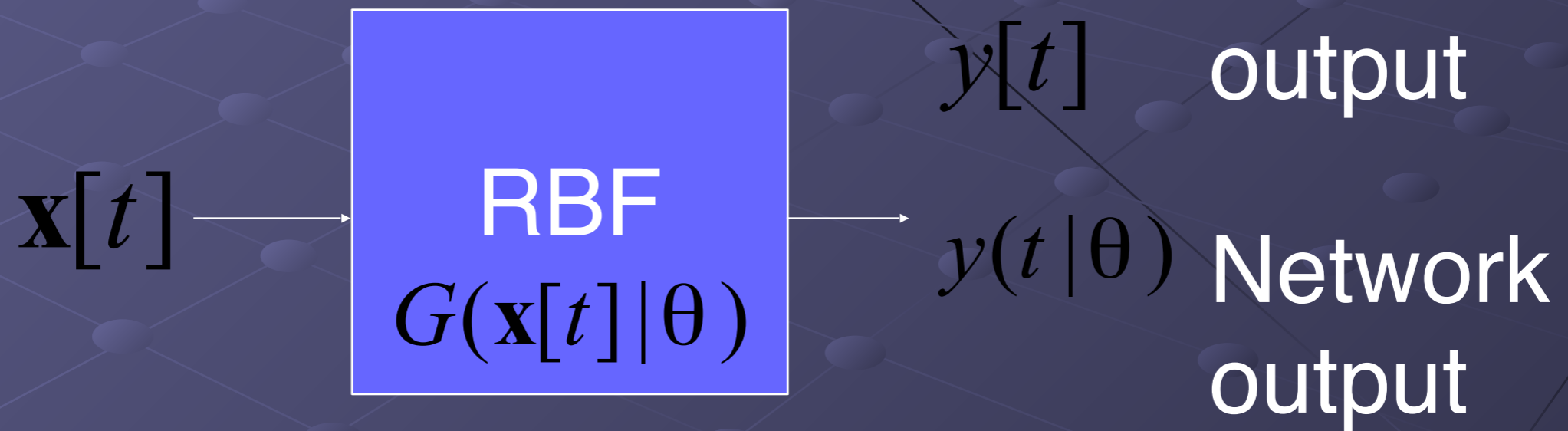
# Training set

$$S = \{(\mathbf{x}[t], y[t])\}_{t=1}^N$$



# Training set

$$S = \{(\mathbf{x}[t], y[t])\}_{t=1}^N$$



# Mean square errors

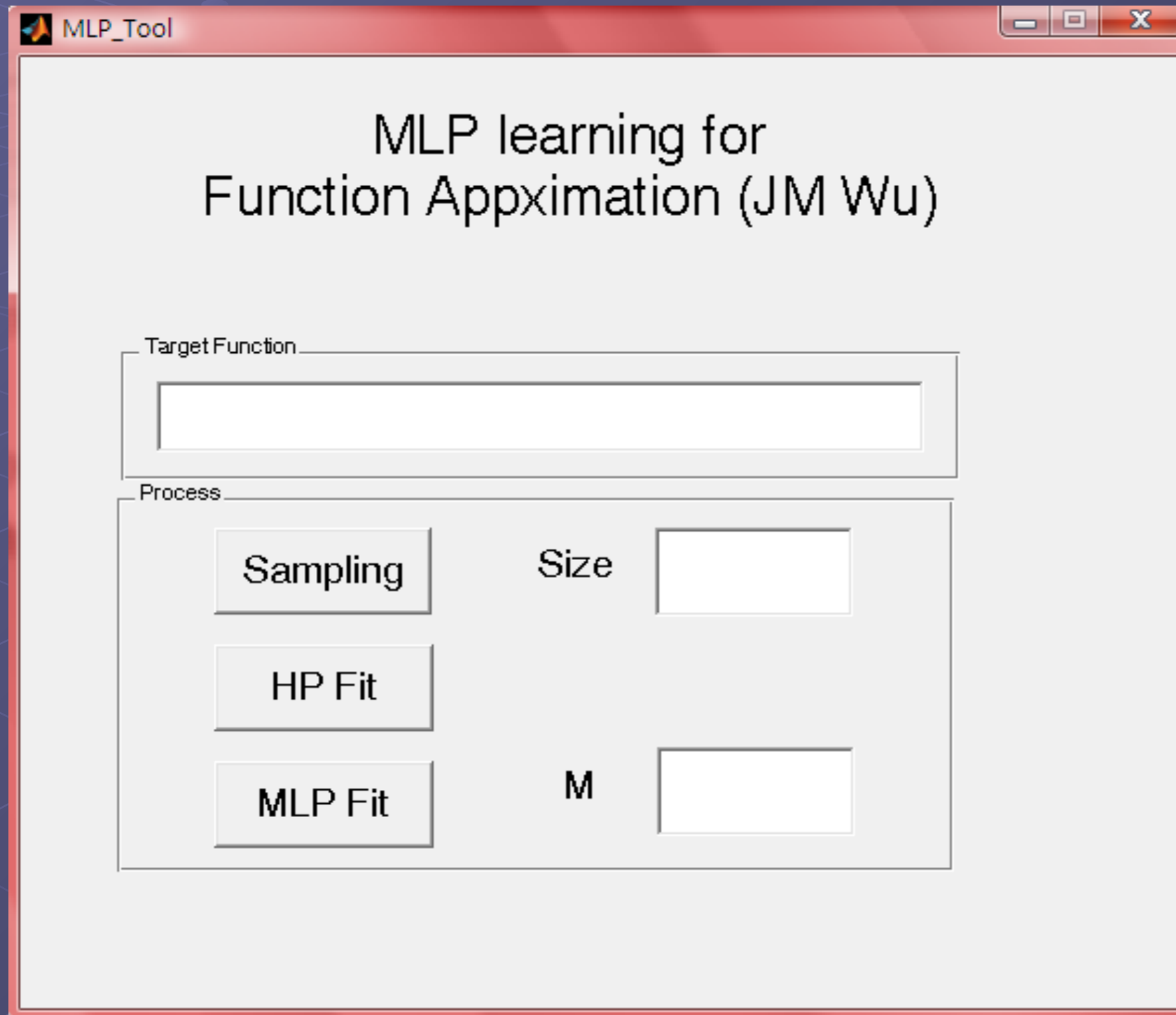
$$E_S(\theta) = \frac{1}{2N} \sum_{t=1}^N (y[t] - y(t|\theta))^2$$

# Unconstrained optimization

$$\hat{\theta} = \arg \min_{\{\theta\}} E_S(\theta)$$

*Find*  $\theta$  to minimize  $E_S(\theta)$

# Revisit hyper-plane Fitting



# Sampling

Blue : training data  
Red : testing data

MLP learning for  
Function Appximation (JM Wu)

Target Function

$\tanh(x_1+x_2)$

Process

Sampling

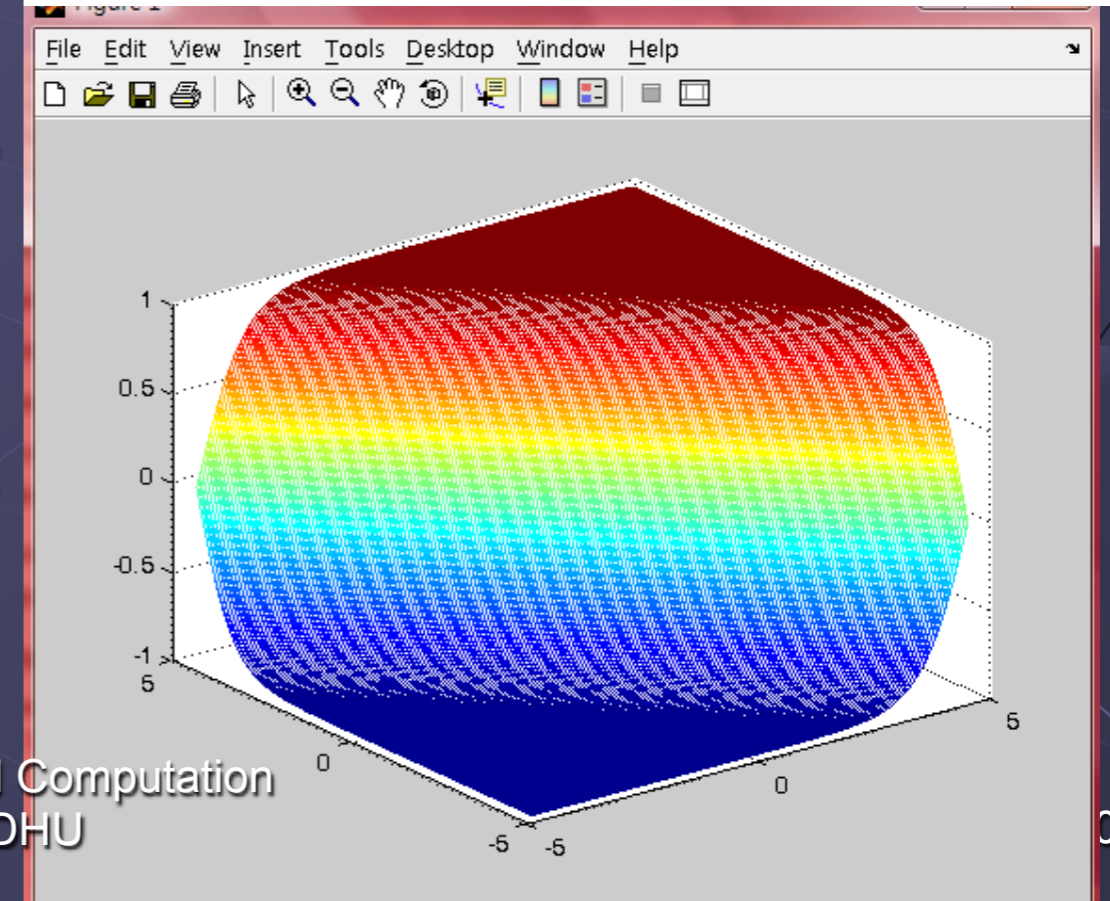
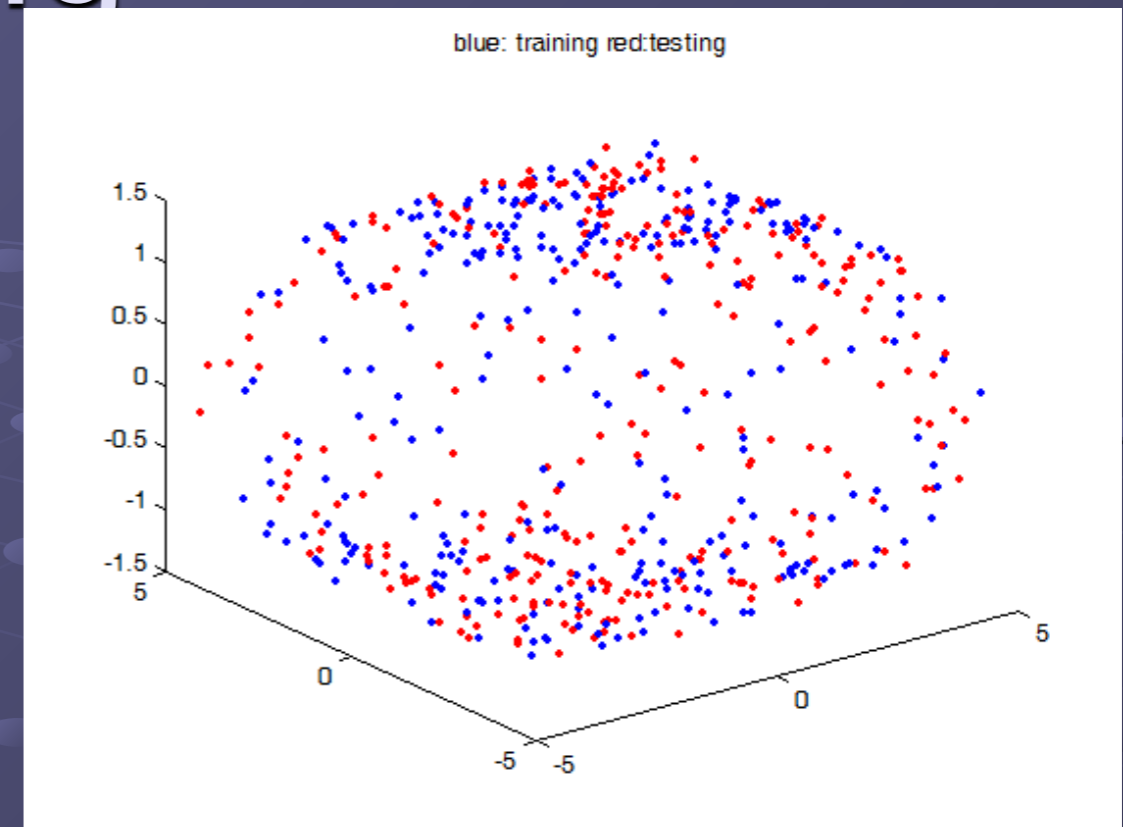
Size

300

HP Fit

MLP Fit

M

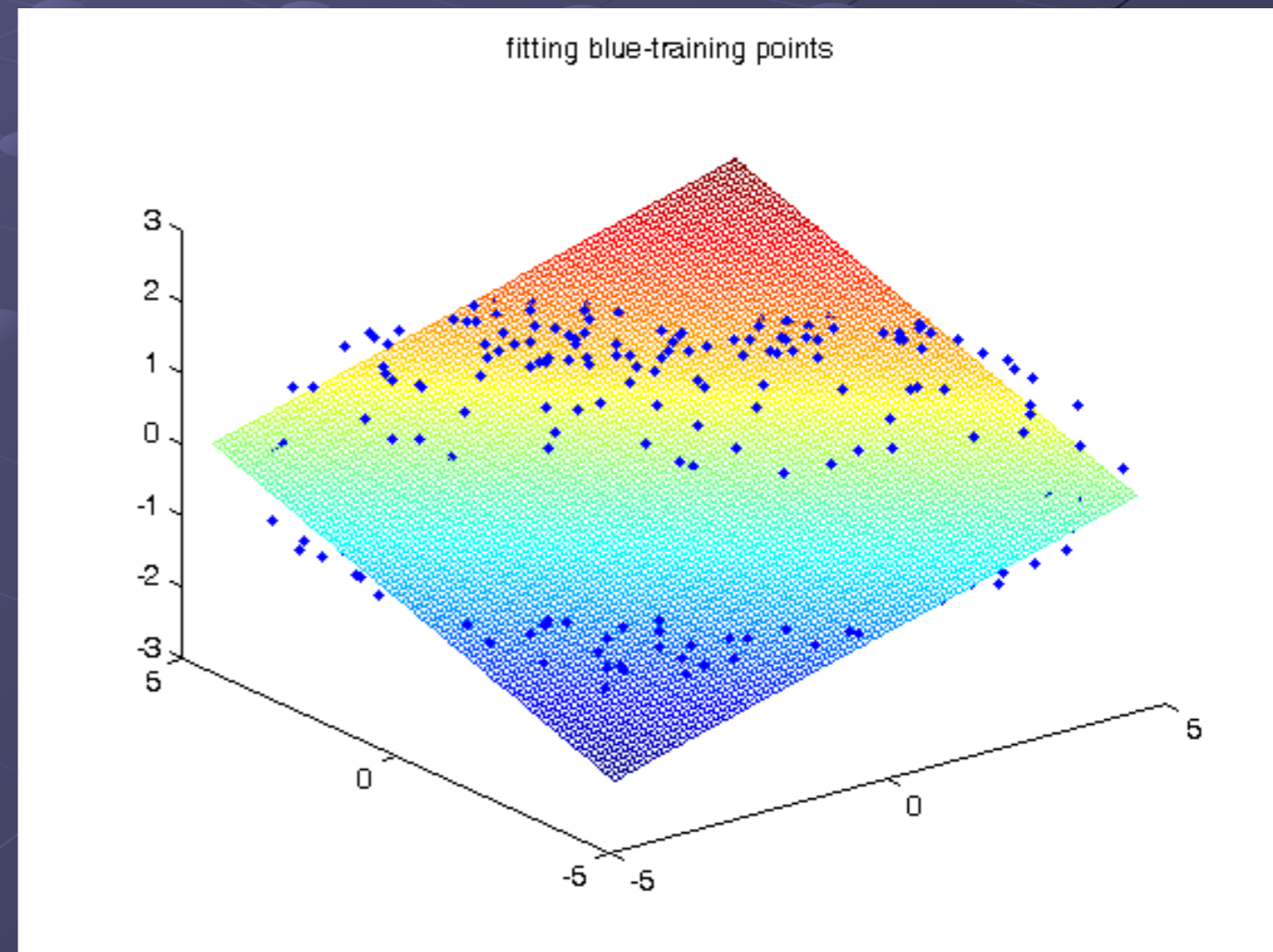




# Fitting

1. Fitting blue training data
2. The performance is reflected by the mean square testing error

mean square training error 0.159647  
mean square testing error 0.160200



# Verification of Generalization

## ● Training phase

- Given a training set, *find*  $\theta$  to minimize  $E_S(\theta)$

## ● Testing phase

- Verify  $F(x; \theta_{\text{opt}})$  by a testing set, which is assumed unknown during training phase

## ● Both training and testing sets are assumed oriented from the same generating process

# Performance evaluation

- Mean square error of training
  - Substitute each predictor in a training set to a network function
  - Determine the mean square error of approximating the desired target by the network output
- Mean square error of testing indicates generalization capability of a network function

# Mean square training error

$$E_S(\theta) = \frac{1}{2N} \sum_{t=1}^N (y[t] - F(\mathbf{x}[t]; \theta))^2$$

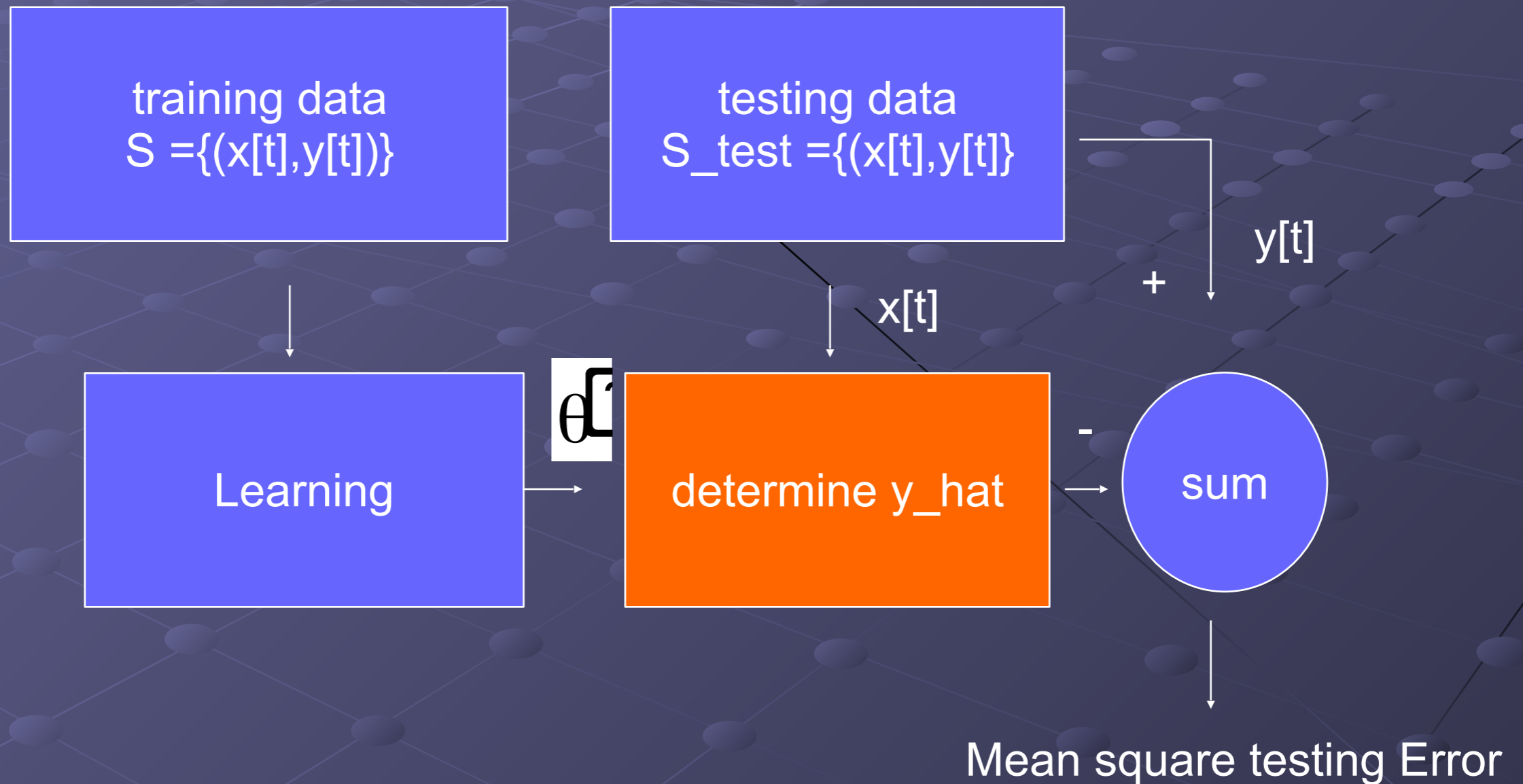
$S = \{(\mathbf{x}[t], y[t])\}_t$  : training set

# Mean square testing error

$$E_{S\_Test}(\theta) = \frac{1}{2N} \sum_{t=1}^N (y[t] - F(\mathbf{x}[t]; \theta))^2$$

$S\_Test = \{(\mathbf{x}[t], y[t])\}_t$  : testing set

# Prediction



# Over-fitting

- Extremely low training error but high testing error
- Over-fitting implies false generalization

# Table lookup

- Over-fitting
- S: training set

$$f(x) = y[t^*]$$

$$t^* = \arg \min_t \| \mathbf{x} - \mathbf{x}[t] \|^2$$



# Iterative approaches

- Gradient method
- Newton-Gauss method
- Levenberg-Marquardt method

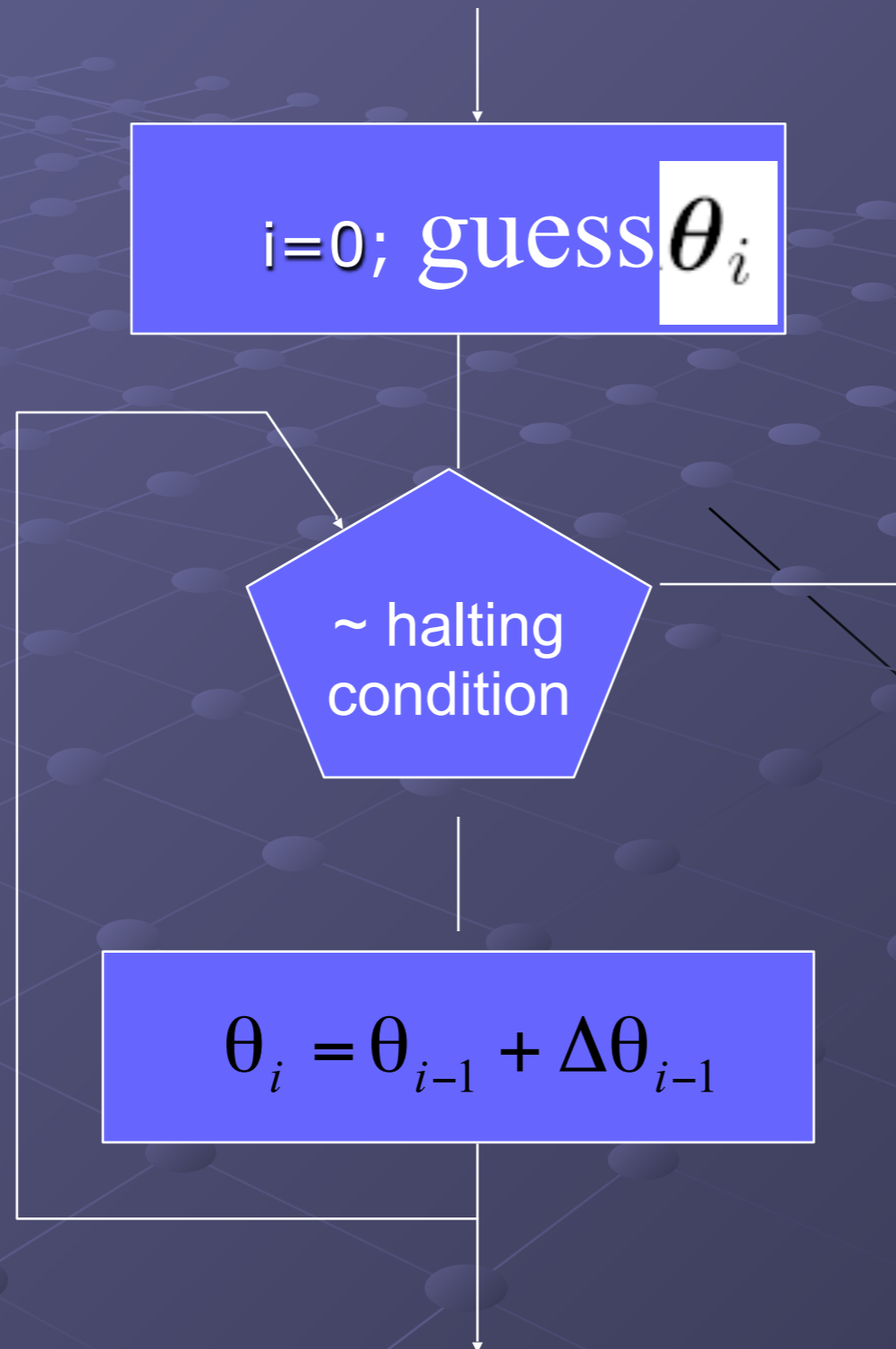
# Iterative approach

1. Initialize  $\theta_i$  with  $i=0$
2. Determine  $\Delta\theta_i$
3. Update network parameters

$$\theta_{i+1} = \theta_i + \Delta\theta_i$$

4. If halting condition holds, exit  
otherwise  $i=i+1$ , go to step 2

# Flow Chart



# Objective function

$$E_S(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{t=1}^N \varepsilon^2(t, \boldsymbol{\theta})$$

$$\varepsilon(t, \boldsymbol{\theta}) = y[t] - y(t|\boldsymbol{\theta})$$

Error of  
the  $i$ th paired data

Desired  
output

Network  
output

# Gradient method

$$\Delta \boldsymbol{\theta}_i \propto - \left. \frac{dE_S(\boldsymbol{\theta})}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_i}$$

- Derive the derivative of  $E_S$  with respect to  $\boldsymbol{\theta}$
- Substitute  $\boldsymbol{\theta}$  to  $\frac{dE_S(\boldsymbol{\theta})}{d\boldsymbol{\theta}}$

# Gradient

$$E_S(\theta) = \frac{1}{2N} \sum_{t=1}^N (y[t] - y(t|\theta))^2$$

$$\frac{dE_S(\theta)}{d\theta} = \frac{-1}{N} \sum_{t=1}^N (y[t] - y(t|\theta)) \frac{dy(t|\theta)}{d\theta}$$

$y(t | \theta)$ 

# Derivation

$$= F(\mathbf{x}[t]; \theta) = \sum_{m=1}^M r_m \tanh(\mathbf{a}_m^T \mathbf{x}[t] + b_m) + r_0$$

$$\theta = [\mathbf{a}_1^T \ \mathbf{a}_2^T \ \cdots \ \mathbf{a}_M^T \ b_1 \ b_2 \ \cdots \ b_M \ r_0 \ r_1 \ r_2 \ \cdots \ r_M]^T$$
$$= [\theta_1, \dots, \theta_{M \cdot d + 2M + 1}]$$

$$\frac{dy(t | \boldsymbol{\theta})}{d\boldsymbol{\theta}}$$

$$= \left( \frac{dy(t | \boldsymbol{\theta})}{d\theta_1} \quad \frac{dy(t | \boldsymbol{\theta})}{d\theta_2} \quad \dots \quad \frac{dy(t | \boldsymbol{\theta})}{d\theta_{Md+2M+1}} \right)^T$$

$$\boldsymbol{\theta} = [\mathbf{a}_1^T \quad \mathbf{a}_2^T \quad \dots \quad \mathbf{a}_M^T \quad b_1 \quad b_2 \quad \dots \quad b_M \quad r_0 \quad r_1 \quad r_2 \quad \dots \quad r_M]^T$$

$$= [\theta_1, \dots, \theta_{M*d+2M+1}]$$



$$\frac{dy(t | \boldsymbol{\theta})}{d\boldsymbol{\theta}}$$

$$= \left( \frac{dy(t | \boldsymbol{\theta})}{d\mathbf{a}_1^T} \quad \frac{dy(t | \boldsymbol{\theta})}{d\mathbf{a}_2^T} \quad \cdots \quad \frac{dy(t | \boldsymbol{\theta})}{db_1^T} \quad \cdots \quad \frac{dy(t | \boldsymbol{\theta})}{dr_1^T} \quad \cdots \quad \frac{dy(t | \boldsymbol{\theta})}{dr_M^T} \right)^T$$

$$\boldsymbol{\theta} = [\mathbf{a}_1^T \quad \mathbf{a}_2^T \quad \cdots \quad \mathbf{a}_M^T \quad b_1 \quad b_2 \quad \cdots \quad b_M \quad r_0 \quad r_1 \quad r_2 \quad \cdots \quad r_M]^T$$

$$= [\theta_1, \dots, \theta_{M*d+2M+1}]$$

$$y(t | \theta)$$

$$= F(\mathbf{x}[t]; \theta) = \sum_{m=1}^M r_m \tanh(\mathbf{a}_m^T \mathbf{x}[t] + b_m) + r_0$$

$$\frac{dy(t | \theta)}{d\mathbf{a}_k}$$

$$d\mathbf{a}_k$$

$$= \frac{d}{d\mathbf{a}_k} \sum_{m=1}^M r_m \tanh(\mathbf{a}_m^T \mathbf{x}[t] + b_m) + r_0$$

$$= r_k (1 - \tanh(\mathbf{a}_k^T \mathbf{x}[t] + b_k))^2 \mathbf{x}[t]$$

$$\begin{aligned} & \frac{dy(t | \boldsymbol{\theta})}{db_k} \\ &= \frac{d}{db_k} \sum_{m=1}^M r_m \tanh(\mathbf{a}_m^T \mathbf{x}[t] + b_m) + r_0 \\ &= r_k (1 - \tanh(\mathbf{a}_k^T \mathbf{x}[t] + b_k)^2) \end{aligned}$$

$$\frac{dy(t | \boldsymbol{\theta})}{dr_k}$$

$$dr_k$$

$$= \frac{d}{dr_k} \sum_{m=1}^M r_m \tanh(\mathbf{a}_m^T \mathbf{x}[t] + b_m) + r_0$$

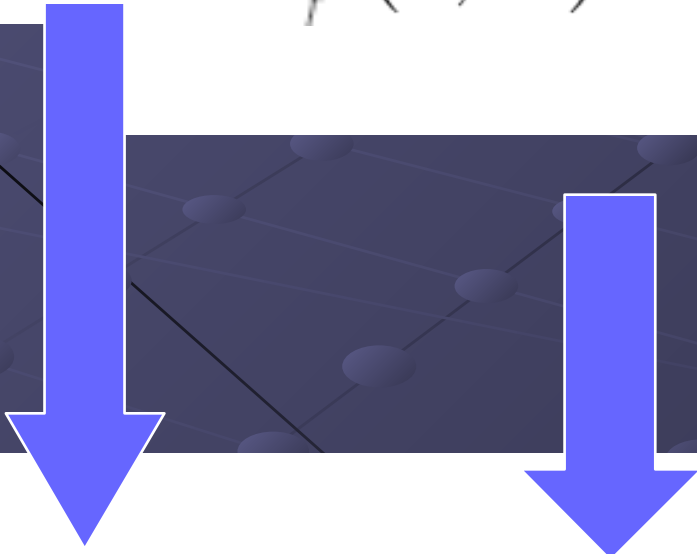
$$= \tanh(\mathbf{a}_k^T \mathbf{x}[t] + b_k)$$

Error

derivative

$$\varepsilon(t, \boldsymbol{\theta}) = y[t] - y(t|\boldsymbol{\theta})$$

$$\psi(t, \boldsymbol{\theta}) = \frac{dy(t|\boldsymbol{\theta})}{d\boldsymbol{\theta}}$$


$$\left. \frac{dE_s(\boldsymbol{\theta})}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_i} = \nabla(\boldsymbol{\theta}_i) = \frac{-1}{N} \sum_{t=1}^N \varepsilon(t, \boldsymbol{\theta}_i) \psi(t, \boldsymbol{\theta}_i)$$

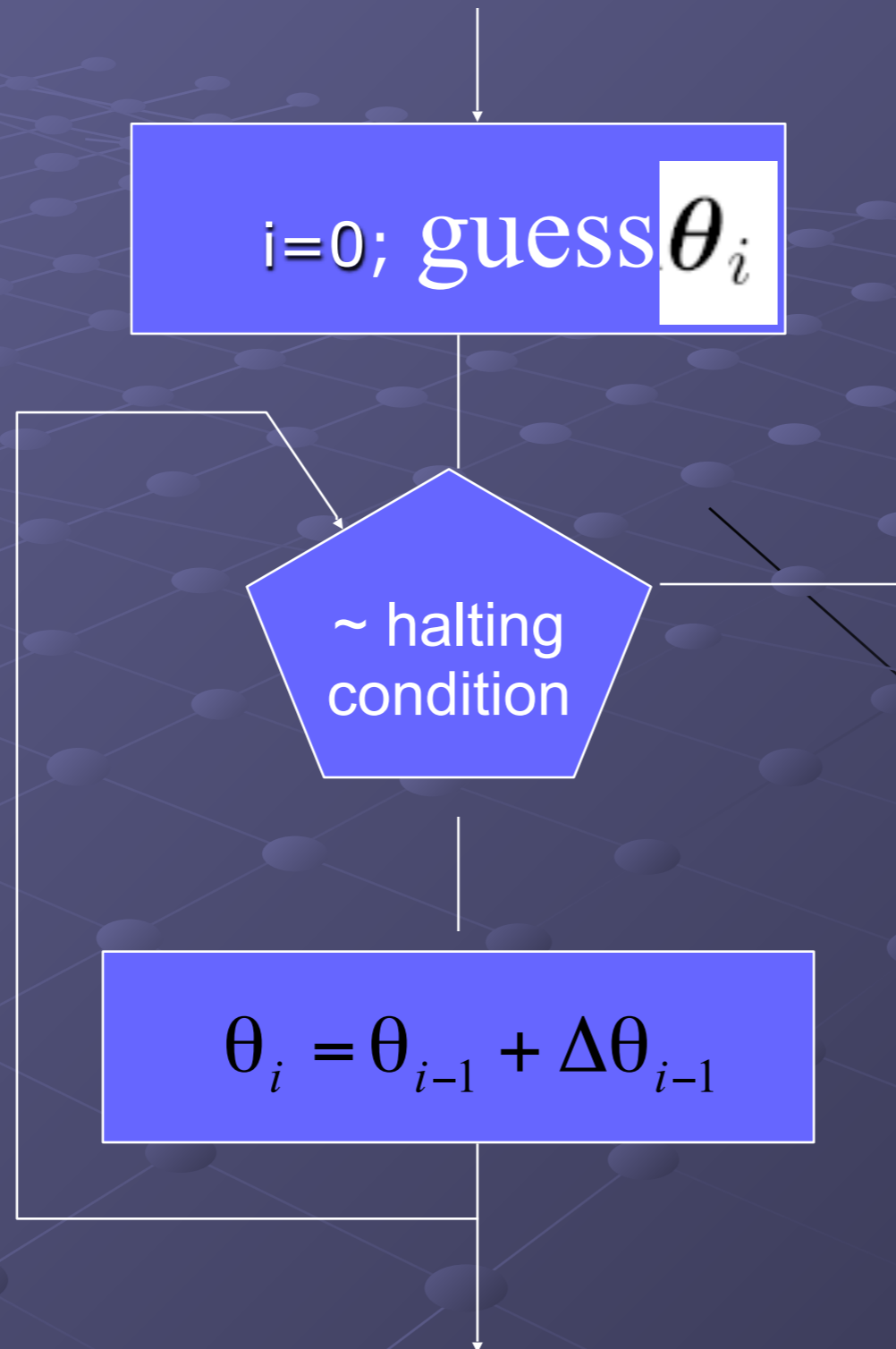
$$\Delta \boldsymbol{\theta}_i \propto - \left. \frac{dE_S(\boldsymbol{\theta})}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_i}$$

$$= -\nabla(\boldsymbol{\theta}_i)$$

# Exercise

- Gradient descent method for learning MLP networks
- Gradient descent method for learning RBF networks

# Flow Chart





# Calculate the change

$$\begin{aligned}\Delta \boldsymbol{\theta}_i &\propto -\left. \frac{dE_S(\boldsymbol{\theta})}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_i} \\ &= \frac{-1}{N} \sum_{t=1}^N \varepsilon(t, \boldsymbol{\theta}_i) \psi(t, \boldsymbol{\theta}_i)\end{aligned}$$

- Prepare  $a, b, r$  from given  $\theta$

- Substitute  $a, b, r$  to

$$\varphi(t, \boldsymbol{\theta}) = \frac{dy(t|\boldsymbol{\theta})}{d\boldsymbol{\theta}}$$

- Substitute  $a, b, r$  to

$$\varepsilon(t, \boldsymbol{\theta}) = y[t] - y(t|\boldsymbol{\theta})$$

for all  $t$

Substitute  $a, b, r$  to

$$\varphi(t, \boldsymbol{\theta}) = \frac{dy(t|\boldsymbol{\theta})}{d\boldsymbol{\theta}}$$

● Substitute current a,b,r to

$$\frac{dy(t | \boldsymbol{\theta})}{d\boldsymbol{\theta}}$$
$$= \left( \frac{dy(t | \boldsymbol{\theta})}{d\mathbf{a}_1^T} \quad \frac{dy(t | \boldsymbol{\theta})}{d\mathbf{a}_2^T} \quad \dots \quad \frac{dy(t | \boldsymbol{\theta})}{db_1^T} \quad \dots \quad \frac{dy(t | \boldsymbol{\theta})}{dr_1^T} \quad \dots \quad \frac{dy(t | \boldsymbol{\theta})}{dr_M^T} \right)^T$$

Substitute  $a, b, r$  to

$$\frac{dy(t | \boldsymbol{\theta})}{d\mathbf{a}_k} = r_k (1 - \tanh(\mathbf{a}_k^T \mathbf{x}[t] + b_k))^2 \mathbf{x}[t]$$

$$\frac{dy(t | \boldsymbol{\theta})}{db_k} = r_k (1 - \tanh(\mathbf{a}_k^T \mathbf{x}[t] + b_k))^2$$

$$\frac{dy(t | \boldsymbol{\theta})}{dr_k} = \tanh(\mathbf{a}_k^T \mathbf{x}[t] + b_k)$$

for all  $k$

# Newton-Gauss Method

Second order expansion

# Linear expansion of errors

Error : nonlinear function of  $\boldsymbol{\theta}$

$$\varepsilon(t, \boldsymbol{\theta}) = y[t] - y(t|\boldsymbol{\theta})$$

Linear expansion at  $\boldsymbol{\theta} = \boldsymbol{\theta}_i$

$$\begin{aligned}\tilde{\varepsilon}_i(t, \boldsymbol{\theta}) &= \varepsilon(t, \boldsymbol{\theta}_i) + (\boldsymbol{\theta} - \boldsymbol{\theta}_i)^T \frac{d\varepsilon(t, \boldsymbol{\theta})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_i} \\ &= \varepsilon(t, \boldsymbol{\theta}_i) - (\boldsymbol{\theta} - \boldsymbol{\theta}_i)^T \frac{dy(t|\boldsymbol{\theta})}{d\boldsymbol{\theta}} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_i}\end{aligned}$$

# Linear expansion of error

Let

$$\psi(t, \boldsymbol{\theta}) = \frac{dy(t|\boldsymbol{\theta})}{d\boldsymbol{\theta}}$$

$$\tilde{\varepsilon}_i(t, \boldsymbol{\theta}) = \varepsilon(t, \boldsymbol{\theta}_i) - (\boldsymbol{\theta} - \boldsymbol{\theta}_i)^T \psi(t, \boldsymbol{\theta}_i)$$

is a linear function of  $\boldsymbol{\theta}$



# Criteria for Newton-Gauss

Mean square  
error

$$E_S(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{t=1}^N \varepsilon^2(t, \boldsymbol{\theta})$$



Quadratic form

$$L_i(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{t=1}^N \tilde{\varepsilon}_i^2(t, \boldsymbol{\theta})$$

# Quadratic form

Expansion at  $\boldsymbol{\theta} = \boldsymbol{\theta}_i$

$$L_i(\boldsymbol{\theta}) = L_i(\boldsymbol{\theta}_i) + \nabla(\boldsymbol{\theta}_i)^T (\boldsymbol{\theta} - \boldsymbol{\theta}_i) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_i)^T R(\boldsymbol{\theta}_i) (\boldsymbol{\theta} - \boldsymbol{\theta}_i),$$

gradient  $\nabla(\boldsymbol{\theta}_i) = \frac{-1}{N} \sum_{t=1}^N \varepsilon(t, \boldsymbol{\theta}_i) \psi(t, \boldsymbol{\theta}_i)$

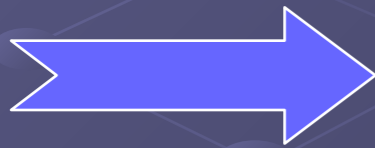
Gauss-Newton Hessian

$$R(\boldsymbol{\theta}_i) = \frac{1}{N} \sum_{t=1}^N \psi(t, \boldsymbol{\theta}_i) \psi^T(t, \boldsymbol{\theta}_i)$$

# Newton-Gauss method

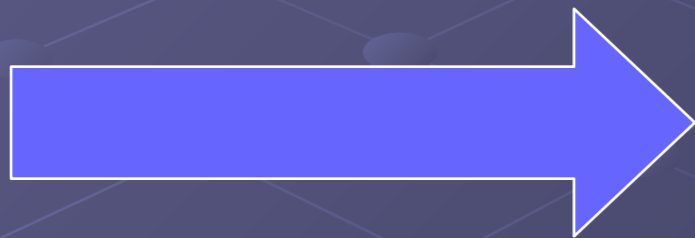
$$\frac{dL_i(\boldsymbol{\theta})}{d\boldsymbol{\theta}} = \nabla(\boldsymbol{\theta}_i) + R(\boldsymbol{\theta}_i)(\boldsymbol{\theta} - \boldsymbol{\theta}_i) = 0$$

$$\Delta\boldsymbol{\theta}_i = \boldsymbol{\theta} - \boldsymbol{\theta}_i$$



$$R(\boldsymbol{\theta}_i)\Delta\boldsymbol{\theta}_i = -\nabla(\boldsymbol{\theta}_i)$$

$$R(\boldsymbol{\theta}_i) \Delta \boldsymbol{\theta}_i = -\nabla(\boldsymbol{\theta}_i)$$



$$\Delta \boldsymbol{\theta}_i = -(R(\boldsymbol{\theta}_i))^{-1} \nabla(\boldsymbol{\theta}_i)$$

# Levenberg-Marquardt method

$\Delta\theta_i$  is determined by

$$(R(\boldsymbol{\theta}_i) + \lambda\mathbf{I})\Delta\boldsymbol{\theta}_i = -\nabla(\boldsymbol{\theta}_i)$$

$$\lambda = 0$$

$$(R(\boldsymbol{\theta}_i) + \lambda \mathbf{I}) \Delta \boldsymbol{\theta}_i = -\nabla(\boldsymbol{\theta}_i)$$



$$R(\boldsymbol{\theta}_i) \Delta \boldsymbol{\theta}_i = -\nabla(\boldsymbol{\theta}_i)$$

Newton-Gauss method

Sufficiently large  $\lambda$

$$(R(\boldsymbol{\theta}_i) + \lambda \mathbf{I}) \Delta \boldsymbol{\theta}_i = -\nabla(\boldsymbol{\theta}_i)$$



$$\Delta \boldsymbol{\theta}_i \propto -\nabla(\boldsymbol{\theta}_i)$$

Gradient method

Control

$\lambda$

Current parameter

Next parameter

$$\alpha_i = \frac{E_S(\boldsymbol{\theta}_i) - E_S(\boldsymbol{\theta}_i + \Delta\boldsymbol{\theta}_i)}{E_S(\boldsymbol{\theta}_i) - L_i(\boldsymbol{\theta}_i + \Delta\boldsymbol{\theta}_i)}$$

Actual cost reduction

---

Predicted cost reduction



High

$\alpha_i$

Reduce

$\lambda$

improve efficiency



Force to  
Newton-Gauss  
method

LOW

$\alpha_i$

Increase

$\lambda$

Improve reliability



Force to  
gradient  
method

# Heuristic adaption

(a) If  $\alpha_i > 0.75$ ,  $\lambda \leftarrow 0.5\lambda$ .

(b) If  $\alpha_i < 0.25$ ,  $\lambda \leftarrow 2\lambda$ .

# LM method

1. Initialize  $\theta_i, i=0$  and set  $\lambda$
2. Calculate  $\nabla R(\theta_i)$  and  $R(\theta_i)$  and  $\Delta\theta_i$
3. Update network parameters

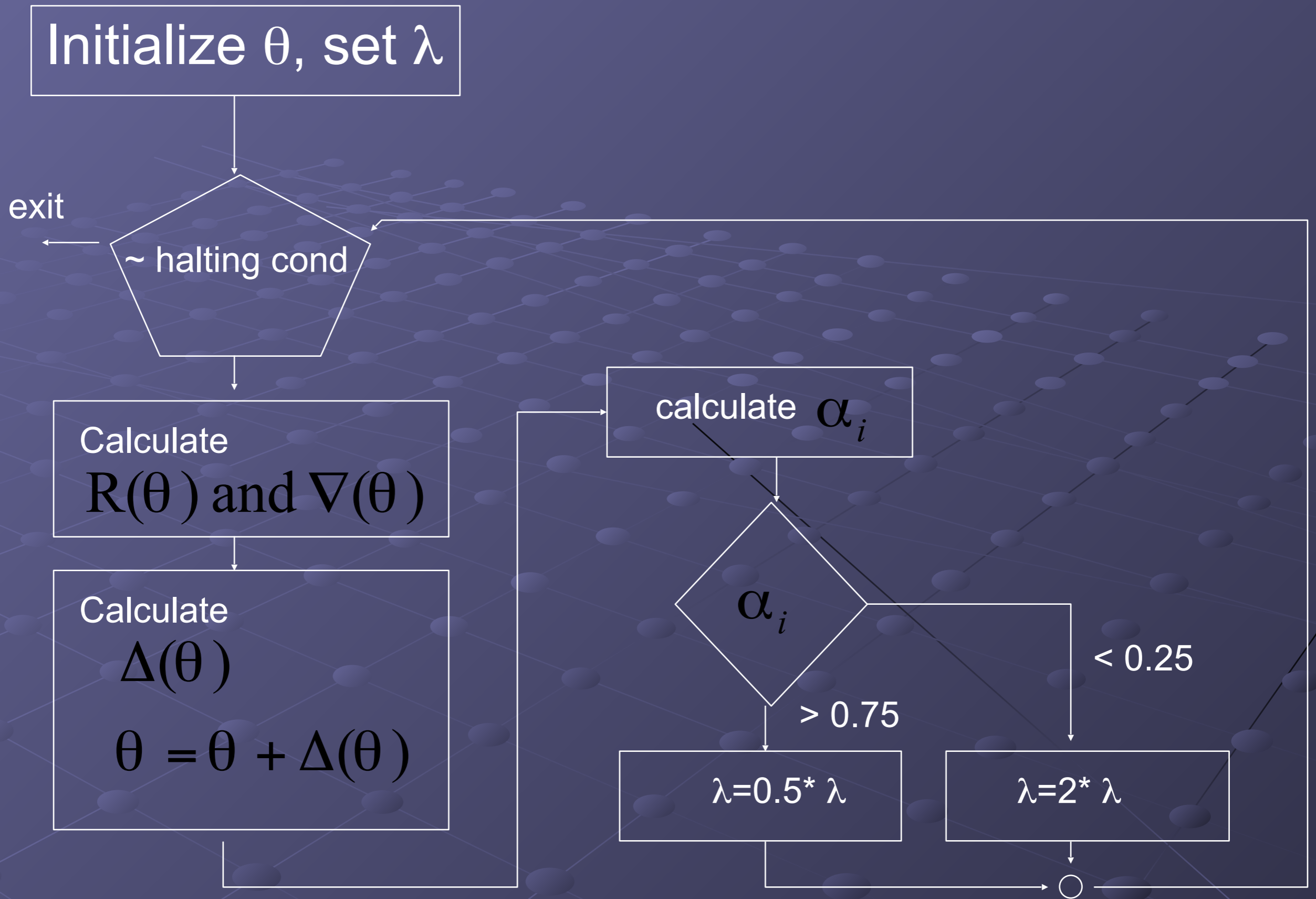
$$\theta_{i+1} = \theta_i + \Delta\theta_i$$

4. Calculate  $\alpha_i$
5. Update  $\lambda$

(a) If  $\alpha_i > 0.75$ ,  $\lambda \leftarrow 0.5\lambda$ .

(b) If  $\alpha_i < 0.25$ ,  $\lambda \leftarrow 2\lambda$ .

6. If halting condition hold, exit otherwise go to step 2



# Project

Learning RBF by the LM method

1. (30 pt) Derivation

$$\frac{dy(t|\theta)}{d\theta} = ?$$

2. (140 pt) Matlab codes

3. (50 pt) Testing

# Derivative

$$\begin{aligned}\theta &= [\mathbf{u}_1^T \ \mathbf{u}_2^T \ \cdots \ \mathbf{u}_M^T \ \sigma_1 \ \sigma_2 \ \cdots \ \sigma_M \ w_0 \ w_1 \ w_2 \ \cdots \ w_M]^T \\ &= [\theta_1, \dots, \theta_{M*d+2M+1}]\end{aligned}$$

$$\frac{dy(t|\theta)}{d\theta} = \left[ \frac{dy(t|\theta)}{d\theta_1}, \dots, \frac{dy(t|\theta)}{d\theta_{M*d+2M+1}} \right]^T$$

# Derivative

$$\frac{dy(t|\theta)}{d\mathbf{u}_m} = ?$$

$$\frac{dy(t|\theta)}{d\sigma_m} = ?$$

$$\frac{dy(t|\theta)}{dw_m} = ?$$



# Matlab coding

- (30 pt) Main program
- Matlab Functions
  - a. (10 pt) Calculate  $E_S(\theta_i)$
  - b. (10 pt) Calculate  $L_i(\theta_i)$

# Matlab coding

- Matlab functions

- (10 pt) Calculate

$$\frac{dy(t|\theta)}{d\mathbf{u}_m} \Big|_{\theta=\theta_i}$$

- (10 pt) Calculate

$$\frac{dy(t|\theta)}{d\sigma_m} \Big|_{\theta=\theta_i}$$

- (10 pt) Calculate

$$\frac{dy(t|\theta)}{dw_m} \Big|_{\theta=\theta_i}$$

# Matlab coding

- Matlab functions

- (20 pt) Calculate  $\nabla(\theta_i)$

- (20 pt) Calculate  $R(\theta_i)$

- (20 pt) Calculate  $G(\mathbf{x}|\theta)$

# Test

- Give two examples to test your matlab codes for learning RBF networks by Levenberg-Marquardt method