# Efficient Training of RBF Networks Via the BYY Automated Model Selection Learning Algorithms

jwma@math.pku.edu.cn

**Abstract.**

## 1 Introduction

## 2  BYY-AMS Adaptive Gradient Learning Algorithm

$x \in X \subset R^d$

$y \in Y \subset R^m$

$p(x,y) = p(x)p(y|x)$   $q(x,y) = q(x|y)q(y)$

$y \in Y = \{1, 2, \cdots, K\} \subset R$   $D_x = \{x_t\}_{t=1}^{N}$

$p(y,x)$ $p(x)$ $q(x,y)$ $q(y)$

$$H(p,q) = \int p(y,x)p(x) \quad q(x,y)q(y) \, dxdy - z_q$$

$z_q$

$p(y,x)$ $q(x,y)$

$\theta$ YY

YY $q(y=j) = \alpha_j$

$\alpha_j \geq \quad \sum_{j=}^{K} \alpha_j = \quad$ $z_q$ $z_q$

$p(x)$ $p(x) = \dfrac{1}{N}\sum_{t=}^{N} \delta(x-x_t)$

$$p(y=j,x) = p(j,x) = \frac{\alpha_j q(x,\theta_j)}{q(x)\Theta_K}$$

$$U_j(x) = \alpha_j q(x, m_j, \Sigma_j) \quad j = 1 \cdots K, \ J(\Theta_K) \qquad \text{(W)}$$

$$J(\Theta_{K}) = \frac{1}{N} \sum_{t=}^{N} J_t(\Theta_{K}), \ J_t(\Theta_K) = \sum_{j=}^{K} \frac{U_j(x_t)}{\sum_{i=}^{K} U_i(x_t)}, \ U_j(x_t) - $$

$$J(\Theta_K) \qquad \beta_j, \ m_j, \ B_j$$

$$\frac{\partial J_t(\Theta_K)}{\partial \beta_j} = \frac{1}{q(x_t, \Theta_k)} \sum_{i=}^{K} \lambda_i(t), \ \delta_{ij} - \alpha_j U_i(x_t) \qquad (\ )$$

$$\frac{\partial J_t(\Theta_K)}{\partial m_j} = p(j, x_t) \lambda_j(t) \Sigma_j^{-} (x_t - m_j) \qquad (\ )$$

$$vec\left(\frac{\partial J_t(\Theta_K)}{\partial B_j}\right) = \frac{\partial (B_j B_j^T)}{\partial B_j} vec\left(\frac{\partial J_t(\Theta_k)}{\partial \Sigma_j}\right) \qquad (\ )$$

$\delta_{ij}$ ... $vec A$ ...

$$\lambda_i(t) = -\sum_{l=}^{K} p(l, x_t) - \delta_{il} \ \alpha_l q(x_t, m_l, \Sigma_l) \qquad (\ )$$

$$\frac{\partial J_t(\Theta_K)}{\partial \Sigma_j} = - p(j, x_t) \lambda_j(t) \ \Sigma_j^{-} (x_t - m_j)(x_t - m_j)^T \Sigma_j^{-} - \Sigma_j^{-} \qquad (\ )$$

$$\frac{\partial (BB^T)}{\partial B} = I_{d\times d} \otimes B_{d\times d}^T + E_{d\times d} \bullet B_{d\times d}^T \otimes I_{d\times d} \qquad$$

$\otimes$ ...

$$E_{d\times d} = \frac{\partial B^T}{\partial B} = \Gamma_{ij\ d\times d} = \begin{pmatrix} \Gamma & \cdots & \Gamma_d \\ \vdots & \ddots & \vdots \\ \Gamma_d & \cdots & \Gamma_{dd} \end{pmatrix}_{d\times d} \qquad$$

$\Gamma_{ij}$ ... $j, i^{th}$ ... $\frac{\partial (BB^T)}{\partial B}$

$$vec \frac{\partial J \, \Theta_{K'}}{\partial B_j} = - p \, j \, x_{t'} \lambda_{j} \, t' \, I_{d\times d} \otimes B_{d\times d}^T + E_{d \times d} \bullet B_{d\times d}^T \otimes I_{d\times d'}$$

$$\times vec \, \Sigma_j^- \, x_t - m_{j'} \, x_t - m_{j'}^T \Sigma_j^- - \Sigma_j^-$$

$$\Delta\beta_j = \frac{\eta}{q \, x_t \, \Theta_{k'}} \sum_{i=}^{K} \lambda_i \, t' \, \delta_{ij} - \alpha_{j'} U_i \, x_{t'}$$

$$\Delta m_j = \eta \, p \, j \, x_{t'} \lambda_{j} \, t' \Sigma_j^- \, x_t - m_{j'}$$

$$\Delta vecB_j = \frac{\eta}{} p \, j \, x_{t'} \lambda_{j} \, t' \, I_{d\times d} \otimes B_{d\times d}^T + E_{d \times d} \bullet B_{d\times d}^T \otimes I_{d\times d'}$$

$$\times vec \, \Sigma_j^- \, x_t - m_{j'} \, x_t - m_{j'}^T \Sigma_j^- - \Sigma_j^-$$

$$\lim_{n\to\infty} \eta \, n' = \quad \sum_{n=}^{\infty} \eta \, n' = \infty$$

$$\eta \, n' = \eta \quad n$$

$$D_x$$

## 3  Training of the RBF Network

$$y_j \, x' = \sum_{j=}^{n} w_{ij} \phi_i \, x'$$

$$\phi_j(x) = \phi_j(\|x - m_j\|) = \exp\left(-\frac{\|x - m_j\|}{\sigma_j}\right)$$

$$y(x) = \sum_{j=1}^{n} \lambda_j \phi_j(x) = \sum_{j=1}^{n} \lambda_j \exp\left(-\frac{\|x - m_j\|}{\sigma_j}\right)$$

$$E = \frac{1}{2}\sum_{i=1}^{N} (y_i - f(x_i))^2 = \frac{1}{2}\sum_{i=1}^{N} \left(y_i - \sum_{j=1}^{n} \lambda_j \phi_j(x_i)\right)^2$$

$$= \frac{1}{2}\sum_{i=1}^{N} \left(y_i - \sum_{j=1}^{n} \lambda_j \exp\left(-\frac{\|x_i - m_j\|}{\sigma_j}\right)\right)^2$$

$$\begin{cases} \Delta\lambda_j = \eta_\lambda \sum_{i=1}^{N} \left(y_i - \sum_{l=1}^{n} \lambda_l \phi_l(x_i)\right) \phi_j(x_i) \\ \Delta m_j = \eta_m \sum_{i=1}^{N} \left(y_i - \sum_{l=1}^{n} \lambda_l \phi_l(x_i)\right) \phi_j(x_i)(x_i - m_j)\lambda_j / \sigma_j \\ \Delta\sigma_j = \eta_\sigma \sum_{i=1}^{N} \left(y_i - \sum_{l=1}^{n} \lambda_l \phi_l(x_i)\right) \phi_j(x_i)(x_i - m_j)^T(x_i - m_j)\lambda_j / \sigma_j \end{cases}$$

$$\sigma_j = \frac{1}{N_j} \sum_{x_t \in C_j} (x_t - m_j)^T (x_t - m_j)$$

## 4  Experiment Results

### 4.1  On the Noisy XOR Problem

**Fig. 1.** ... X

**Fig. 2.** ...

## 4.2 On the Mackey-Glass Time Series Prediction

$$x(t+1) = -bx(t) + \frac{ax(t-\tau)}{1+x(t-\tau)},$$

$a =$ $b =$ $\tau =$

$x(t-\ ) x(t-\ ) x(t-\ ) x(t) x(t+\ )$ $\le t \le$

$y_i = f(x_i)$ $x_i = [x(t-\ ) x(t-\ ) x(t-\ ) x(t)]^T$ $y_i = x(t+\ )$

$i = \cdots N$

**Fig. 3.**



**Fig. 4.**

## 5   Conclusions

# References