

Radial-Basis Function Networks

v 1.0 – March 2002



Michel Verleysen

Radial-Basis Function Networks - 1

Radial-Basis Function Networks

- // Origin: Cover's theorem
- // Interpolation problem
- // Regularization theory
- // Generalized RBFN
 - // Universal approximation
 - // Comparison with MLP
 - // RBFN = kernel regression
- // Learning
 - // Centers
 - // Widths
 - // Multiplying factors
 - // Other forms



Michel Verleysen

Radial-Basis Function Networks - 2

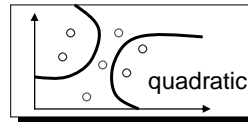
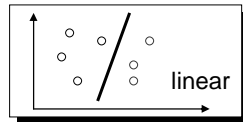
Origin: Covers' theorem

- ⚡ Covers' theorem on separability of patterns (1965)
- ⚡ $\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^P$ assigned to two classes C^1, C^2
- ⚡ φ -separability:

$$\exists \mathbf{w} \left\{ \begin{array}{l} \mathbf{w}^T \varphi(\mathbf{x}) > 0 \quad \mathbf{x} \in C^1 \\ \mathbf{w}^T \varphi(\mathbf{x}) < 0 \quad \mathbf{x} \in C^2 \end{array} \right.$$

- ⚡ Cover's theorem:
 - ⚡ non-linear functions $\varphi(\mathbf{x})$
 - ⚡ dimension hidden space $>$ dimension input space
 - probability of separability closer to 1

- ⚡ Example



Michel Verleysen

Radial-Basis Function Networks - 3

Interpolation problem

- ⚡ Given points (\mathbf{x}^k, t^k) , $\mathbf{x}^k \in \mathfrak{X}^d$, $t^k \in \mathfrak{Y}$, $1 \leq k \leq P$:
- ⚡ Find $F: \mathfrak{X}^d \rightarrow \mathfrak{Y}$ that satisfies

$$F(\mathbf{x}^k) = t^k, \quad k = 1 \dots P$$

- ⚡ RBF technique (Powell, 1988):

$$F(\mathbf{x}) = \sum_{k=1}^P w_k \varphi(\|\mathbf{x} - \mathbf{x}^k\|)$$

- ⚡ $\varphi(\|\mathbf{x} - \mathbf{x}^k\|)$ are arbitrary non-linear functions (RBF)
- ⚡ as many functions as data points
- ⚡ centers fixed at known points \mathbf{x}^k



Michel Verleysen

Radial-Basis Function Networks - 4

Interpolation problem

$$F(\mathbf{x}^k) = t^k \qquad F(\mathbf{x}) = \sum_{k=1}^P w_k \phi(\|\mathbf{x} - \mathbf{x}^k\|)$$

$$\begin{bmatrix} \phi_{11} & \phi_{12} & \cdots & \phi_{1P} \\ \phi_{21} & \phi_{22} & \cdots & \phi_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ \phi_{P1} & \phi_{P2} & \cdots & \phi_{PP} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_P \end{bmatrix} = \begin{bmatrix} t^1 \\ t^2 \\ \vdots \\ t^P \end{bmatrix}$$

where
 $\phi_{kl} = \phi(\|\mathbf{x}^k - \mathbf{x}^l\|)$

⚡ Into matrix form: $\Phi \mathbf{w} = \mathbf{x} \rightarrow \mathbf{w} = \Phi^{-1} \mathbf{x}$

⚡ Vital question: is Φ non-singular ?



Michelli's theorem

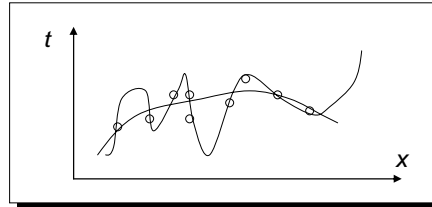
⚡ If points \mathbf{x}^k are distinct, Φ is non-singular (regardless of the dimension of the input space)

⚡ Valid for a large class of RBF functions:

$\phi(\mathbf{x}) = \sqrt{\ \mathbf{x} - \mathbf{c}\ ^2 + k^2} \quad (k > 0)$	non-localized function
$\phi(\mathbf{x}) = \frac{1}{\sqrt{\ \mathbf{x} - \mathbf{c}\ ^2 + k^2}}$	localized functions
$\phi(\mathbf{x}) = \exp\left(-\frac{\ \mathbf{x} - \mathbf{c}\ ^2}{2\sigma^2}\right) \quad (\sigma > 0)$	



Learning: ill-posed problem



⚡ Necessity for *regularization*

⚡ Error criterion:

$$E(F) = \underbrace{\frac{1}{2P} \sum_{k=1}^P (t^k - F(\mathbf{x}^k))}_{\text{MSE}} + \underbrace{\lambda \frac{1}{2} C(\mathbf{w})}_{\text{regularization}}$$



Michel Verleysen

Radial-Basis Function Networks - 7

Solution to the regularization problem

⚡ Poggio & Girosi (1990):

⚡ if $C(\mathbf{w})$ is a (problem-dependent) linear differential operator, the solution to

$$E(F) = \frac{1}{2P} \sum_{k=1}^P (t^k - F(\mathbf{x}^k)) + \lambda \frac{1}{2} C(\mathbf{w})$$

is of the following form:

$$F(\mathbf{x}) = \sum_{k=1}^P w_k G(\mathbf{x}, \mathbf{x}^k)$$

where $G()$ is a Green's function,

$$\mathbf{w} = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{t}$$

$$G_{kl} = G(\mathbf{x}^k, \mathbf{x}^l)$$



Michel Verleysen

Radial-Basis Function Networks - 8

Interpolation - Regularization

⚡ Interpolation

$$F(\mathbf{x}) = \sum_{k=1}^P w_k \varphi(\|\mathbf{x} - \mathbf{x}^k\|)$$

$$\mathbf{w} = \Phi^{-1} \mathbf{t}$$

⚡ Exact interpolator

⚡ Possible RBF:

$$\varphi(\mathbf{x}, \mathbf{x}^k) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}^k\|^2}{2\sigma^2}\right)$$

⚡ Regularization

$$F(\mathbf{x}) = \sum_{k=1}^P w_k G(\mathbf{x}, \mathbf{x}^k)$$

$$\mathbf{w} = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{t}$$

⚡ Exact interpolator

⚡ Equal to the « interpolation » solution iff $\lambda=0$

⚡ Example of Green's function:

$$G(\mathbf{x}, \mathbf{x}^k) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}^k\|^2}{2\sigma^2}\right)$$

One RBF / Green's function for each learning pattern!



Generalized RBFN (GRBFN – RBFN)

⚡ As many radial functions as learning patterns:

⚡ computationally (too) intensive
(inversion of $P \times P$ matrix grows with P^3)

⚡ ill-conditioned matrix

⚡ regularization not easy (problem-specific)

→ *Generalized* RBFN approach!

$$F(\mathbf{x}) = \sum_{i=1}^K w_i \varphi(\|\mathbf{x} - \mathbf{c}_i\|)$$

Typically:

⚡ $K \ll P$

$$\varphi(\|\mathbf{x} - \mathbf{c}_i\|) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$

Parameters:

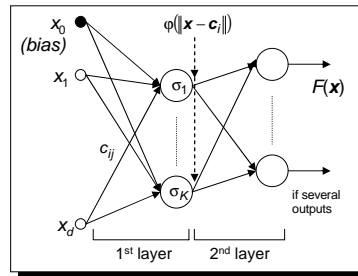
$\mathbf{c}_p, \sigma_p, w_i$



Radial-Basis Function Networks (RBFN)

$$F(\mathbf{x}) = \sum_{i=1}^K w_i \varphi(\|\mathbf{x} - \mathbf{c}_i\|)$$

$$\varphi(\|\mathbf{x} - \mathbf{c}_i\|) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$



Possibilities:

- /// several outputs (common hidden layer)
- /// bias (recommended) (see extensions)



RBFN: universal approximation

/// Park & Sandberg 1991:

/// For any continuous input-output mapping function $f(\mathbf{x})$

$$\exists F(\mathbf{x}) = \sum_{i=1}^K w_i \varphi(\|\mathbf{x} - \mathbf{c}_i\|) \mid L_p(f(\mathbf{x}), F(\mathbf{x})) < \varepsilon \quad (\varepsilon > 0, p \in [1, \infty])$$

- /// The theorem is stronger (radial symmetry not needed)
- /// K not specified
- /// Provides a theoretical basis for *practical* RBFN!



RBFN and kernel regression

⚡ non-linear regression model

$$t^k = f(\mathbf{x}^k) + \varepsilon^k = y^k + \varepsilon^k, 1 \leq k \leq P$$

⚡ estimation of $f(\mathbf{x})$: average of t around \mathbf{x} . More precisely:

$$\begin{aligned} f(\mathbf{x}) &= E[y|\mathbf{x}] \\ &= \int_{-\infty}^{\infty} y f_Y(y|\mathbf{x}) dy \\ &= \frac{\int_{-\infty}^{\infty} y f_{\mathbf{X},Y}(\mathbf{x}, y) dy}{f_{\mathbf{X}}(\mathbf{x})} \end{aligned}$$

⚡ Need for estimates of $f_{\mathbf{X},Y}(\mathbf{x}, y)$ and $f_{\mathbf{X}}(\mathbf{x})$

→ Parzen-Rosenblatt density estimator



Parzen-Rosenblatt density estimator

$$\hat{f}_{\mathbf{X}}(\mathbf{x}) = \frac{1}{Ph^d} \sum_{k=1}^P K\left(\frac{\mathbf{x} - \mathbf{x}^k}{h}\right)$$

with $K()$ continuous, bounded, symmetric about the origin, with maximum value at 0, and with unit integral, is consistent (asymptotically unbiased).

⚡ Estimation of

$$\hat{f}_{\mathbf{X},Y}(\mathbf{x}, y) = \frac{1}{Ph^{d+1}} \sum_{k=1}^P K\left(\frac{\mathbf{x} - \mathbf{x}^k}{h}\right) K\left(\frac{y - y^k}{h}\right)$$



RBFN and kernel regression

$$\hat{f}(\mathbf{x}) = \frac{\int_{-\infty}^{\infty} y f_{\mathbf{x},\mathbf{y}}(\mathbf{x}, y) dy}{\hat{f}_{\mathbf{x}}(\mathbf{x})} \quad \leftarrow \quad f(\mathbf{x}) = \frac{\int_{-\infty}^{\infty} y f_{\mathbf{x},\mathbf{y}}(\mathbf{x}, y) dy}{f_{\mathbf{x}}(\mathbf{x})}$$

$$= \frac{\sum_{k=1}^P y^k K\left(\frac{\mathbf{x} - \mathbf{x}^k}{h}\right)}{\sum_{k=1}^P K\left(\frac{\mathbf{x} - \mathbf{x}^k}{h}\right)}$$

- ⚡ Weighted average of y^i
- ⚡ called Nadaraya-Watson estimator (1964)
- ⚡ equivalent to *Normalized RBFN* in the unregularized context



RBFN ↔ MLP

⚡ RBFN

- ⚡ single hidden layer
- ⚡ non-linear hidden layer
linear output layer
- ⚡ argument of hidden units:
Euclidean norm
- ⚡ universal approximation
property
- ⚡ local approximators
- ⚡ splitted learning

⚡ MLP

- ⚡ single or multiple hidden layers
- ⚡ non-linear hidden layer
linear or non-linear output layer
- ⚡ argument of hidden units:
scalar product
- ⚡ universal approximation
property
- ⚡ global approximators
- ⚡ global learning



RBFN: learning strategies

$$F(\mathbf{x}) = \sum_{i=1}^K w_i \varphi(\|\mathbf{x} - \mathbf{c}_i\|) \quad \varphi(\|\mathbf{x} - \mathbf{c}_i\|) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$

- ⚡ Parameters to be determined: \mathbf{c}_i , σ_i , w_i
- ⚡ Traditional learning strategy: splitted computation
 1. centers \mathbf{c}_i
 2. widths σ_i
 3. weights w_i



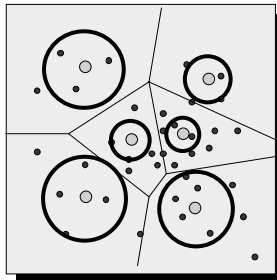
RBFN: computation of centers

- ⚡ Idea: centers \mathbf{c}_i must have the (density) properties of learning points \mathbf{x}^k
 - vector quantization
 - ⚡ selected at random (in learning set)
 - ⚡ competitive learning
 - ⚡ frequency-sensitive learning
 - ⚡ Kohonen maps
- ⚡ This phase only uses the \mathbf{x}^k information, not the t^k



RBFN: computation of widths

- ⚡ Universal approximation property: valid with identical widths
- ⚡ In practice (limited learning set): variable widths σ_i
- ⚡ Idea: RBFN use *local* clusters



- ⚡ choose σ_i according to standard deviation of clusters



RBFN: computation of weights

$$F(\mathbf{x}) = \sum_{i=1}^K w_i \varphi(\|\mathbf{x} - \mathbf{c}_i\|) \quad \varphi(\|\mathbf{x} - \mathbf{c}_i\|) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)$$

constants !

- ⚡ Problem becomes linear !
- ⚡ Solution of least square criterion $E(F) = \frac{1}{2P} \sum_{k=1}^P (t^k - F(\mathbf{x}^k))^2$ leads to

$$\mathbf{w} = \Phi^+ \mathbf{t} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

where

$$\Phi \equiv \varphi_{ki} = \varphi(\|\mathbf{x}^k - \mathbf{c}_i\|)$$

- ⚡ In practise: use SVD !



RBFN: gradient descent

$$F(\mathbf{x}) = \sum_{i=1}^K w_i \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right) \quad 1$$

3-steps method: supervised unsupervised

- ⚡ Once \mathbf{c}_i , σ_i , w_i have been set by the previous method, possibility of gradient descent on *all* parameters
- ⚡ Some improvement, but
 - ⚡ learning speed
 - ⚡ local minima
 - ⚡ risk of non-local basis functions
 - ⚡ etc.



More elaborated models

- ⚡ Add constant and linear terms

$$F(\mathbf{x}) = \sum_{i=1}^K w_i \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right) + \sum_{i=1}^d w'_i x_i + w'_0$$

good idea (very difficult to approximate a constant with kernels...)

- ⚡ Use normalized RBFN

$$F(\mathbf{x}) = \sum_{i=1}^K w_i \frac{\exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right)}{\sum_{j=1}^K \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|^2}{2\sigma_j^2}\right)}$$

basis functions are bounded [0,1] → can be interpreted as probability values (classification)



Back to the widths...

∥ choose σ_i according to standard deviation of clusters

∥ In the literature:

∥ $\sigma = d_{\max} / \sqrt{2K}$ where d_{\max} = maximum distance between centroids [1]

∥ $\sigma_i = \frac{1}{p} \sqrt{\sum_{j=1}^p \|\mathbf{c}_i - \mathbf{c}_j\|^2}$ where index j scans the p nearest centroids to \mathbf{c}_i [2]

∥ $\sigma_i = r \min_j (\|\mathbf{c}_i - \mathbf{c}_j\|)$ where r is an overlap constant [3]

∥

[1] S. Haykin, "Neural Networks a Comprehensive Foundation", Prentice-Hall Inc, second edition, 1999.

[2] J. Moody and C. J. Darken, "Fast learning in networks of locally-tuned processing units", Neural Computation 1, pp. 281-294, 1989.

[3] A. Saha and J. D. Keeler, "Algorithms for Better Representation and Faster Learning in Radial Basis Function Networks", Advances in Neural Information Processing Systems 2, Edited by David S. Touretzky, pp. 482-489, 1989.

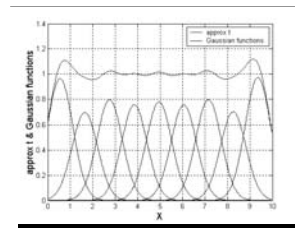


Basic example

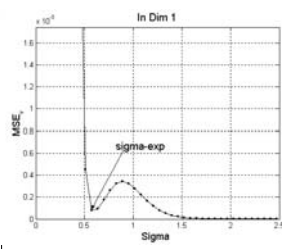
∥ Approximation of $f(\mathbf{x}) = 1$ with a d -dimensional RBFN

∥ In theory: identical w_i

∥ Experimentally: side effects
→ only middle taken into account



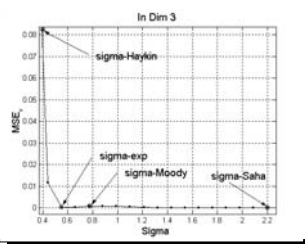
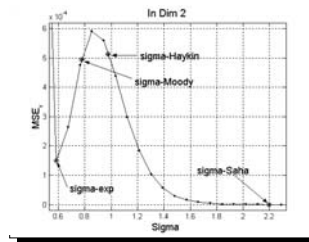
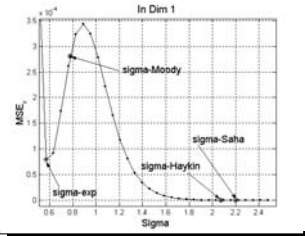
∥



← Error versus width



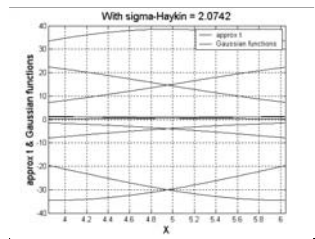
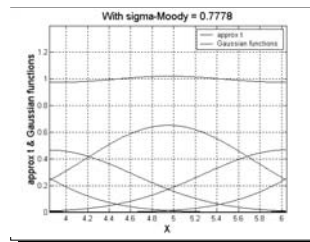
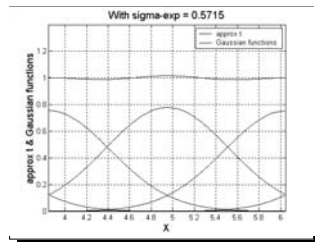
Basic example: erros vs space dimension



Michel Verleysen

Radial-Basis Function Networks - 25

Basic example: local decomposition?

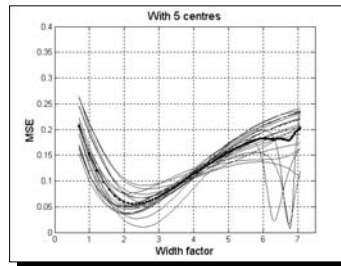


Michel Verleysen

Radial-Basis Function Networks - 26

Multiple local minima in error curve

- ⚡ Choose the first minimum to preserve the locality of clusters



- ⚡ The first local minimum is usually less sensitive to variability



Some concluding comments

- ⚡ RBFN: easy learning (compared to MLP)
 - ⚡ in a cross-validation scheme: important!
- ⚡ Many RBFN models
- ⚡ Even more RBFN learning schemes...
- ⚡ Results not very sensitive to unsupervised part of learning (\mathbf{c}_i, σ_i)
- ⚡ Open work for a priori (proble-dependent) choice of widths σ_i



Sources and references

- ⚡ Most of the basic concepts developed in these slides come from the excellent book:
 - ⚡ Neural networks – a comprehensive foundation, S. Haykin, Macmillan College Publishing Company, 1994.
- ⚡ Some supplementary comments come from the tutorial on RBF:
 - ⚡ An overview of Radial Basis Function Networks, J. Ghosh & A. Nag, in: Radial Basis Function Networks 2, R.J. Howlett & L.C. Jain eds., Physica-Verlag, 2001.
- ⚡ The results on the basic exemple were generated by my colleague N. Benoudjit, and are submitted for publication.

