



Neural Networks

Lecture 8

Radial-basis function networks

Igor Farkaš

2019

2

Radial-Basis-Function neural network

- Inputs x , weights w , outputs y
- Output activation:

$$y_i = \sum_{k=1}^K w_{ik} h_k(x) + w_{i0}$$

- h_k = radial activ. function, e.g.

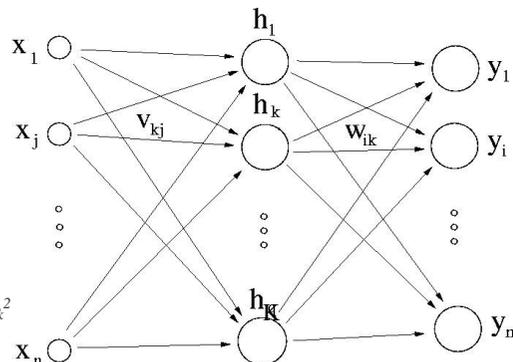
$$h_k(x) = \varphi_k(\|x - v_k\|) = \exp(-\|x - v_k\|^2 / \sigma_k^2)$$

$v_k \sim$ center k , $\sigma_k \sim$ its width

$\varphi(d)$ are (usually) **local** functions because for $d \rightarrow \infty$ $\varphi(d) \rightarrow 0$

σ affects generalization

- v_k used for approximation of unconditional probability density of input data $p(x)$
- RBF as a receptive field (easier than that of an MLP)



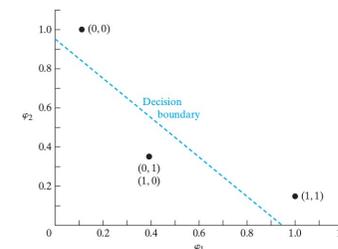
3

Combined NN models

- combination of unsupervised and supervised learning
- independent optimization, can be much faster than gradient descent, with similar results
- unsupervised learning \rightarrow clustering
- more hidden units may be needed (compared to a completely supervised model)
- Examples:
 - \rightarrow counter-propagation networks (Hecht-Nielsen, 1987)
 - \rightarrow learning vector quantization (Kohonen, 1990)
 - \rightarrow radial-basis-function networks (Moody & Darken, 1989)

Separability of patterns

- Data projection into high-dim. space:
A complex pattern classification problem cast in a high-dim. space nonlinearly is more likely to be linearly separable than in a low-dim. space (Cover, 1965).
- Consider binary partitioning (dichotomy) for x_1, x_2, \dots, x_N (classes C_1, C_2). Dichotomy $\{C_1, C_2\}$ is ϕ -separable, where $\phi(x) = [\varphi_1(x), \varphi_2(x), \dots, \varphi_K(x)]$, if $\exists w \in \mathbb{R}^K$ such that for $\forall x \in C_1$: $w^T \cdot \phi(x) > 0$ and for $\forall x \in C_2$: $w^T \cdot \phi(x) < 0$.
- $\{\varphi_k(x)\}$ - **feature** functions (hidden space), $k = 1, 2, \dots, K$
- Sometimes, non-linear transformation can result in linear separability without having to increase data dimension (e.g. XOR problem):



$$\varphi_k(x) = \exp(-\|x - v_k\|^2) \quad v_1 = [0 \ 0], \quad v_2 = [1 \ 1]$$

Input Pattern x	First Hidden Function $\varphi_1(x)$	Second Hidden Function $\varphi_2(x)$
(1,1)	1	0.1353
(0,1)	0.3678	0.3678
(0,0)	0.1353	1
(1,0)	0.3678	0.3678

4

Interpolation problem

- Mapping data into higher dimensions can be useful:
- Then we can deal with multivariate interpolation in high-dim. space (Davis, 1963):
Given the sets $\{\mathbf{h}_i \in \mathbb{R}^K, d_i \in \mathbb{R}\}$, find a function F that satisfies the condition: $F(\mathbf{h}_i) = d_i, i=1,2,\dots,N$. (in strict sense)
- For RBF, we get the set of linear equations: $\mathbf{w}^T \mathbf{h}_i = d_i, i = 1,2,\dots,N$.
- If \mathbf{H}^{-1} exists, the solution is $\mathbf{w} = \mathbf{H}^{-1} \mathbf{d}$
- How can we be sure that **interpolation matrix** \mathbf{H} is nonsingular?
- Theorem: Let $\{\mathbf{x}_i \in \mathbb{R}^n\}$ be a set of distinct points ($i=1,2,\dots,N$). Then \mathbf{H} [$N \times N$] with elements $h_{ij} = \varphi_{ij}(\|\mathbf{x}_i - \mathbf{x}_j\|)$, is nonsingular. (Michelli, 1986)
- a large class of RBFs satisfies this condition

5

Various types of radial-basis functions

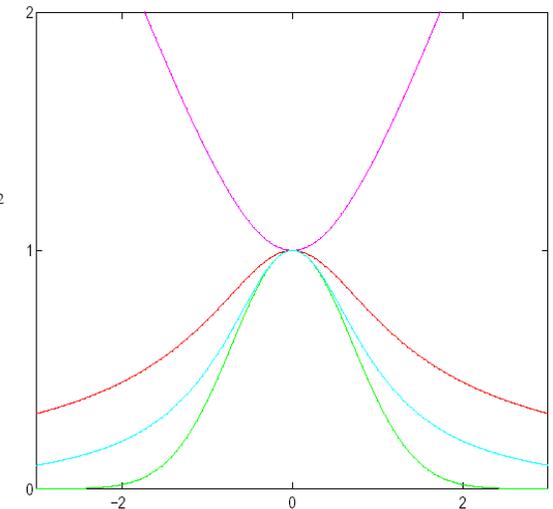
Gaussian: $\varphi(r) = \exp(-r^2/\sigma^2)$

Multiquadrics: $\varphi(r) = (r^2+c^2)^{1/2}$

Inverse multiquadrics: $\varphi(r) = (r^2+c^2)^{-1/2}$

Cauchy: $\varphi(r) = 1/(1+r^2)$
 $r \in \mathbb{R}, c > 0$

Matrix \mathbf{H} for nonlocal multiquadrics is not positive definite. Despite that it can be used to approximate a smooth I/O mapping with greater accuracy than those that yield positive-definite interpolation matrix (Powell, 1988).



6

Training RBF networks

- two-stage process
- **nonlinear** (layer 1) and **linear** (layer 2) optimization strategies are applied to different learning tasks
- **Approaches for layer 1:**
 - fixed centers selected at random
 - self-organized selection of centers
- **Approaches for layer 2**
 - via pseudoinverse \mathbf{H}^+ : then $\mathbf{w} = \mathbf{H}^+ \mathbf{d}$
 - online stochastic optimization (delta rule),
 - online deterministic algorithm (RLS)
- Yet another method: supervised selection of centers and output weight setting (not described here)

7

Fixed centers selected at random

- “sensible” approach if training data are distributed in a representative manner:
 $G(\|\mathbf{x} - \mathbf{v}_j\|^2) = \exp(-K\|\mathbf{x} - \mathbf{v}_j\|^2/d_{\max}^2)$
 K – number of centers, $d_{\max} = \max_{k,l} \{\|\mathbf{v}_k - \mathbf{v}_l\|\}$, $\Rightarrow \sigma = d_{\max}/(2K)^{1/2}$
- RBFs become neither too flat nor too wide
- Alternative: individual widths σ_j , inversely proportional to density $p(\mathbf{x})$ – requires experimentation with data
- relatively insensitive to regularization, for larger data sets

8

Self-organized selection of centers

Self-organization: **K-means** clustering:

Initialization: randomize $\{v_1(0), v_2(0), \dots, v_k(0)\}$

Two steps: (until stopping criterion is met)

1. minimize $J(C) = \min_{\{v_k\}} \sum_{k=1}^K \sum_{C(i)=k} \|\mathbf{x}(i) - \mathbf{v}_k\|^2$ for given encoder C
 - by updating cluster centers: $\{v_k(t)\}$
2. optimize the encoder: $C(i) = \arg \min_k \|\mathbf{x}(i) - \mathbf{v}_k\|^2$
 - by reassigning inputs to clusters

Given a set of N observations, find the encoder C that assigns these observations to the K clusters in such a way that, within each cluster, the average measure of dissimilarity of the assigned observations from the cluster mean is minimized.

- no guarantee for finding an optimum

Recursive Least Squares (RLS)

- RBF centers can be updated recursively
- How to compute optimal output weights, recursively, too?
- RLS algorithm summary: given $\{\phi(p), d(p)\}, p=1,2,\dots,N; p \equiv t$
- *Initialize:* $w(0) = \mathbf{0}, \mathbf{P}(0) = \lambda^{-1} \mathbf{I}$, with $\lambda > 0, \lambda \approx 0$, regularizer $\frac{1}{2} \lambda \|\mathbf{w}\|^2$
- *Repeat:*
 1. $\mathbf{P}(t) = \mathbf{P}(t-1) - \frac{\mathbf{P}(t-1) \Phi(t) \Phi^T(t) \mathbf{P}(t-1)}{1 + \Phi^T(t) \mathbf{P}(t-1) \Phi(t)}$
 2. $\mathbf{g}(t) = \mathbf{P}(t) \cdot \phi(t)$ (gain)
 3. $a(t) = d(t) - \mathbf{w}^T(t-1) \phi(t)$ (prior estimation error)
 4. $\mathbf{w}(t) = \mathbf{w}(t-1) + \mathbf{g}(t) \cdot a(t)$

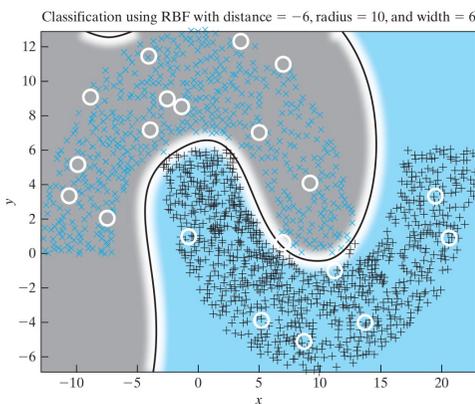
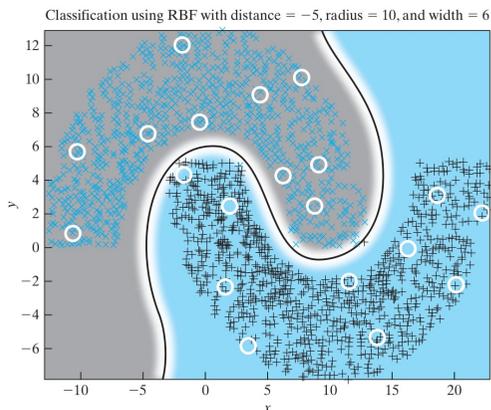
Example using an RBF network

Two-moons classification task: 20 Gaussian units, 1000 points used for training, 2000 for testing. Different widths (σ) used.

$$\sigma = \frac{d_{\max}}{\sqrt{2K}}$$

$\sigma = 2.6$

$\sigma = 2.4$



Approximation properties of RBF networks

Theorem: (Park & Sandberg, 1991) Let $G: \mathfrak{R}^K \rightarrow \mathfrak{R}$ be an integrable bounded function such that G is continuous and $\int_{\mathfrak{R}^K} G(\mathbf{x}) d\mathbf{x} \neq 0$. The family of RBF networks consists of functions $F: \mathfrak{R}^K \rightarrow \mathfrak{R}$:

$$F(\mathbf{x}) = \sum_{k=1}^K w_k G((\mathbf{x} - \mathbf{v}_k)/\sigma)$$

where $\sigma > 0, w_k \in \mathfrak{R}$ and $\mathbf{v}_k \in \mathfrak{R}^K$.

Then for any continuous function $f(\mathbf{x})$ there exists an RBF network with a set of centers $\mathbf{v}_k \in \mathfrak{R}^q$ and a common width $\sigma > 0$ such that $F(\mathbf{x})$ realized by RBF network is close to $f(\mathbf{x})$ in L_p norm, $p \in [1, \infty]$.

Note: Theorem does not require radial symmetry for kernel $G: \mathfrak{R}^K \rightarrow \mathfrak{R}$.

- Useful constraint in RBF design: $K < N$ (number of patterns)
- Gaussian centers as kernels: $\int_{\mathfrak{R}^K} G(\mathbf{x}) d\mathbf{x} = 1$

Kernel $G(\mathbf{x}) =$ continuous, bounded, and real function of \mathbf{x} , symmetric about the origin, where it attains its maximum value.

Comparison of RBF and MLP

- both are nonlinear layered feedforward networks
- both are **universal approximators**, using parametrized compositions of functions of single variables.
- localized vs. distributed representations on hidden layer =>
 - convergence of RBF may be faster
 - MLPs are global, RBF are local => MLP need **fewer** parameters
- different designs of a supervised network:
 - MLP = **stochastic** approximation problem
 - RBF = **hypersurface-fitting** problem in a high-dim. space
- one-stage (MLP) vs. two-stage (RBF) training scheme

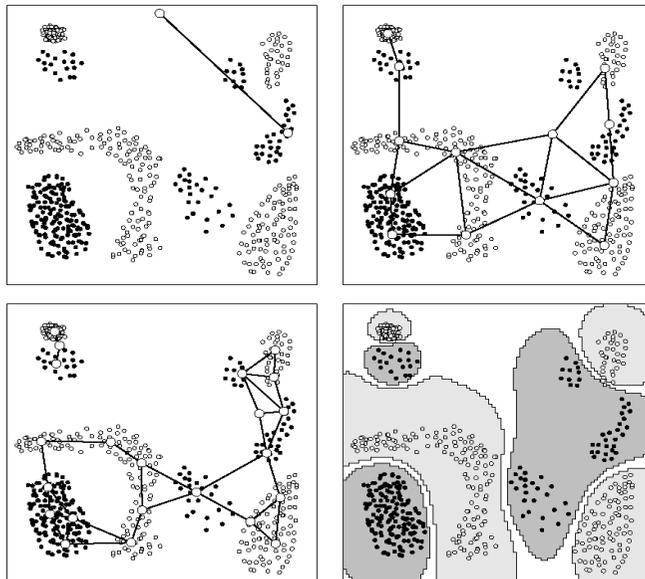
13

Alternative self-organizing modules for center allocation

- Can be useful for input data
 - with varying dimensionality across input domain (e.g. Topology Representing Network)
 - with non-stationary distributions – dynamic networks (Dynamic Cell Structures, Growing CS)
- to be coupled with dynamic linear part
- all based on competitive learning

14

Example: binary classification with a growing RBF net



(Fritzke, 1994)

Summary

- RBF – hybrid feedforward NN model
 - first layer unsupervised, second layer supervised
- we deal with interpolation problem – as a hypersurface reconstruction problem
- various training algorithms for RBF centers
- RLS algorithm for computing output weights
- RBF as universal approximator
- applicable for function approximation and for classification

15

16