



Neural Networks

Lecture 8

Expansion of hidden-layer dimension

Igor Farkaš

2020

2

Combined NN models

- combination of unsupervised and supervised learning
- independent optimization, can be much faster than gradient descent, with similar results
- unsupervised learning → clustering
- more hidden units may be needed (compared to a completely supervised model)
- Examples:
 - learning vector quantization (Kohonen, 1990)
 - classifier on top of trained SOM
 - radial-basis-function networks (Moody & Darken, 1989)

Changes in data dimensionality

- Neural networks process data by nonlinearly transforming them over layers
- **Dimensionality reduction** has many advantages:
 - allows to extract features
 - leads to abstraction(s)
 - allows robustness against noise
- **Dimensionality expansion** leads to what?
 - allows linear separability of inputs
 - hence, their better separability

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.

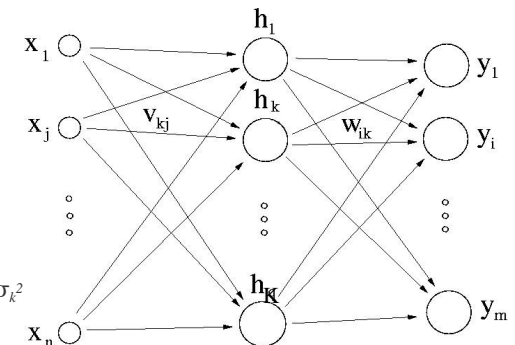
$$h_k(x) = \varphi_k(\|x - v_k\|) = \exp(-\|x - v_k\|^2 / \sigma_k^2)$$

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

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

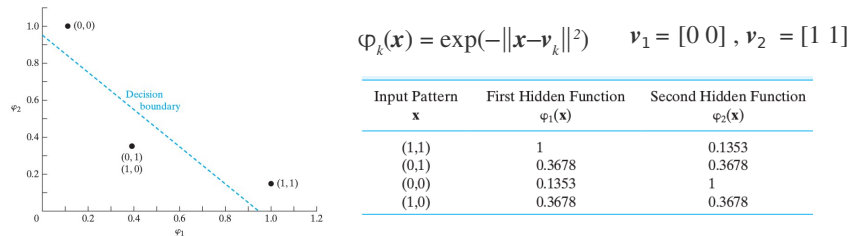
σ affects generalization

- 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)



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 $\forall x \in C_1: w^T \Phi(x) > 0$ and for $\forall x \in C_2: w^T \Phi(x) < 0$.
- $\{\varphi_k(x)\}$ – **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):



5

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 $\{h_i \in \mathbb{R}^K, d_i \in \mathbb{R}\}$, find a function F that satisfies the condition: $F(h_i) = d_i, i=1, 2, \dots, N$. (in strict sense)
- For RBF, we get the set of linear equations: $w^T h_i = d_i, i = 1, 2, \dots, N$.
- If H^{-1} exists, the solution is $w = H^{-1} d$
- How can we be sure that **interpolation matrix** H is nonsingular?
- Theorem: Let $\{x_i \in \mathbb{R}^n\}$ be a set of distinct points ($i=1, 2, \dots, N$). Then H $[N \times N]$ with elements $h_{ij} = \varphi_{ij}(\|x_i - x_j\|)$, is nonsingular. (Michelli, 1986)
- a large class of RBFs satisfies this condition

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 H^+ : then $w = H^+ d$
 - 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:
 $G(\|x - v_j\|^2) = \exp(-K\|x - v_j\|^2/d_{\max}^2)$
 K – number of centers, $d_{\max} = \max_{k,l} \{\|v_k - v_l\|\}$, $\Rightarrow \sigma = d_{\max}/(2K)^{1/2}$
- RBFs become neither too flat nor too wide
- Alternative: individual widths σ_j , inversely proportional to density $p(x)$ – 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 $\{\mathbf{v}_1(0), \mathbf{v}_2(0), \dots, \mathbf{v}_K(0)\}$

Two steps: (until stopping criterion is met)

1. minimize $J(C) = \min_{\mathbf{v}_k} \sum_{k=1}^K \sum_{C(i)=k} \|\mathbf{x}(i) - \mathbf{v}_k\|^2$ for given encoder C
 - by updating cluster centers: $\{\mathbf{v}_k(t)\}$
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

9

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:* $\mathbf{w}(0) = \mathbf{0}, \mathbf{P}(0) = \lambda^{-1} \mathbf{I}$, 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. $\mathbf{g}(t) = \mathbf{P}(t)\Phi(t)$ (gain)
 3. $a(t) = d(t) - \mathbf{w}^T(t-1)\Phi(t)$ (prior estimation error)
 4. $\mathbf{w}(t) = \mathbf{w}(t-1) + \mathbf{g}(t)a(t)$

10

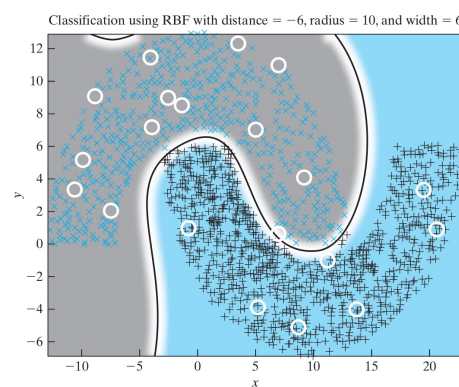
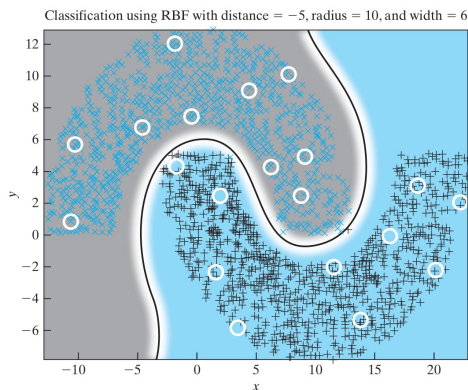
Example using an RBF network

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

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

$\sigma = 2.6$

$\sigma = 2.4$



11

Approximation properties of RBF networks

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

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

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

Then for any continuous function $f(\mathbf{x})$ there exists an RBF network with a set of centers $\mathbf{v}_k \in \mathcal{R}^K$ 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 $G: \mathcal{R}^K \rightarrow \mathcal{R}$.

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

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

12

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

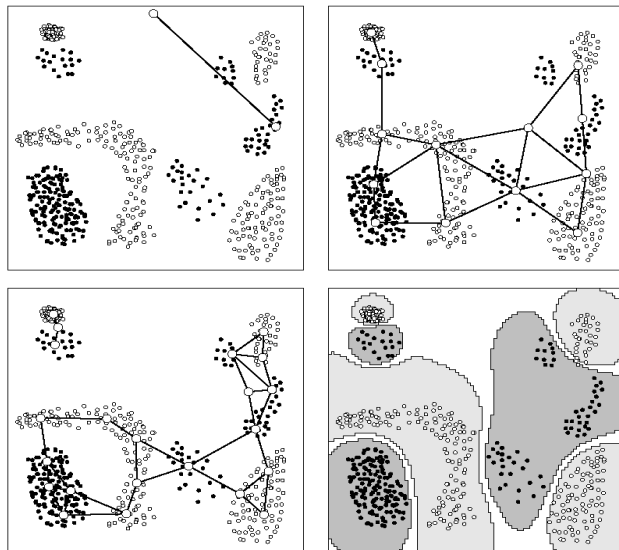
13

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)

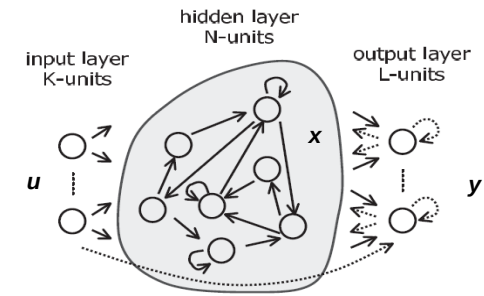
15

Reservoir computing

- A relatively new framework for computation derived from a RNN that maps input signals into **higher dimensional** computational spaces through the dynamics of a fixed, non-linear system called a reservoir (Schrauwen et al, 2007).
- After the input signal is fed into the reservoir, which is treated as a "black box," a **simple readout mechanism** is trained to read the state of the reservoir and map it to the desired output.
- This has **two benefits**: (1) training is performed only at the readout stage, (2) computational efficiency, with very good accuracy in many tasks.
- Best known models are **echo state network** (with classical neurons) and **liquid state machines** (with spiking neurons).

16

Echo-state network



(Jaeger, 2001)

ESN can have an SRN architecture, but also **additional connections** are possible (useful for some tasks).

Reservoir units: usually nonlinear (tanh), can also be linear.

System equations:

$$\mathbf{x}(t) = f(\mathbf{W} \mathbf{x}(t-1) + \mathbf{W}^{\text{inp}} \mathbf{u}(t) + \mathbf{W}^{\text{fb}} \mathbf{y}(t))$$

$$\mathbf{y}(t) = f^{\text{out}}(\mathbf{W}^{\text{out}} \mathbf{z}(t))$$

$$\mathbf{z}(t) = [\mathbf{x}(t); \mathbf{u}(t)]$$

$$\mathbf{W}^{\text{out}} \sim L \times (N+K)$$

Note: these pathways (dotted lines in figure) will not be considered.

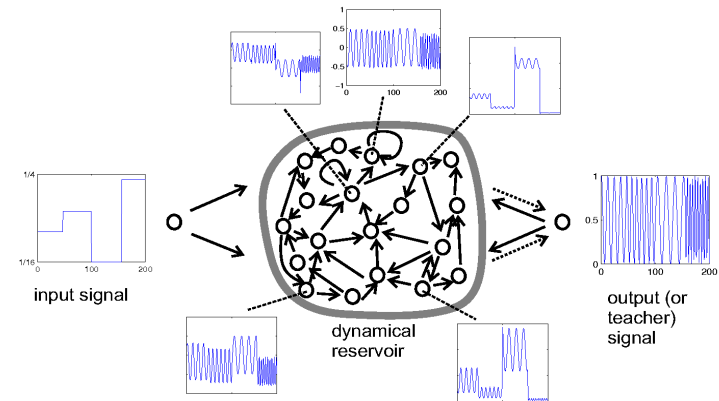
17

ESN training

- **Initialize the ESN**
 - create the reservoir with **echo-state property** (asymptotic properties of reservoir dynamics are given by driving signal): (Jaeger, 2001)
 - Network $F: X \times U \rightarrow X$ (with compactness condition) has the **echo state property** w.r.t. U , if for any left infinite input sequence $u^{-\infty} \in U^{-\infty}$ and any two state vector sequences $x^{-\infty}, y^{-\infty} \in X^{-\infty}$ compatible with $u^{-\infty}$, it holds that $x_0 = y_0$.
 - small random input weights (with uniform or gaussian distribution)
- **Collect reservoir states**
 - feed the input sequence into the network (recursively apply the state equation)
- **Compute output weights**
 - Supervised learning, via pseudoinverse of \mathbf{X} , or RLS
- ESN reservoir has a Markov property (in symbolic dynamics)

19

Echo State Network (ctd)



- studied issues: memory capacity, information transfer, ...
- edge of stability = interesting regime (may be optimal w.r.t. info processing)

18

ESN properties

Echo-state property (ESP): depends on spectral properties of \mathbf{W} = (typically) random sparse matrix, measures:

- spectral radius: $\rho(\mathbf{W}) = |\lambda_{\max}|$, i.e. largest absolute eigenvalue,
- spectral norm: $s_{\max}(\mathbf{W})$ = largest singular value, relation: $0 \leq \rho(\mathbf{W}) \leq s_{\max}(\mathbf{W})$
- Criteria for ESP: $s_{\max}(\mathbf{W}) < 1 \rightarrow$ too strict, $\rho(\mathbf{W}) < 1$ not sufficient
- New recipe (Yildiz & Jaeger, 2012): (i) random $w_{ij} \geq 0$, (ii) scale \mathbf{W} so that $\rho(\mathbf{W}) < 1$, (iii) change the signs of a desired number of entries to get some $w_{ij} < 0$ as well.
- $\rho(\mathbf{W}) \approx 1$ tends to be a “turning point” in behavior (e.g. memory capacity)

Memory capacity (MC): reflects the ability to retrieve input data from the reservoir

- scalar i.i.d. inputs assumed, MC depends on \mathbf{W} , \mathbf{W}^{inp} , reservoir size N , sparsity,...

$$\text{MC} = \sum_{k=1}^{k_{\max}} \text{MC}_k = \sum_{k=1}^{k_{\max}} \frac{\text{cov}^2(u(t-k), y_k(t))}{\text{var}(u(t)) \cdot \text{var}(y_k(t))} \quad y_k(t) = \mathbf{w}_k^{\text{out}} \mathbf{x}(t) = \tilde{u}(t-k)$$

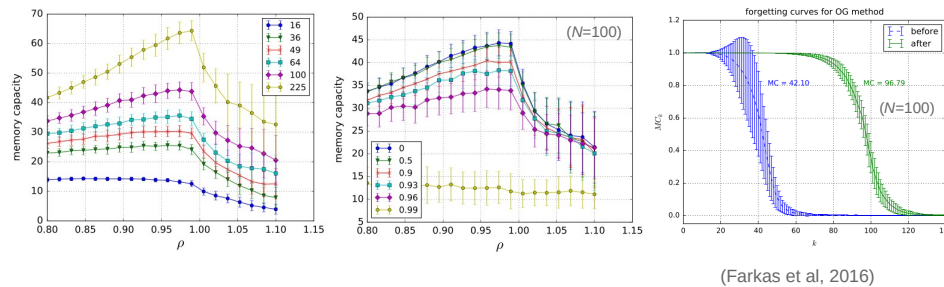
$$k_{\max} = L$$

Reservoir stability – measured by characteristic Lyapunov exponent (Sprott, 2003), that quantifies the average divergence of state space trajectories under perturbations.

20

Memory capacity – calculated

- MC depends on spectral radius ρ and grows with reservoir size N (left)
 - for $\rho > 1$ the dynamics may become unstable
- MC degrades very gracefully for sparse reservoirs (middle)
- MC can be increased by (iterative) reservoir orthogonalization (right)
 - reaching the theoretical limit (N)



21

Summary

- **RBF** – hybrid feedforward NN model
 - hidden layer unsupervised (high-dim. projection), output layer supervised (linear readout)
 - various training algorithms for setting RBF centers
 - RLS for computing output weights, or pseudoinverse
 - universal approximator (like MLP)
 - applicable for function approximation and classification
- **ESN** – fast recurrent NN, only linear readout trained
 - reservoir = high-dim. spatio-temporal embedding
 - good for time series prediction and memory tasks with Markov properties

22