

ON PATTERN CLASSIFICATION AND SYSTEM IDENTIFICATION BY PROBABILISTIC NEURAL NETWORKS

LESZEK RUTKOWSKI*, TOMASZ GAŁKOWSKI*

In the paper probabilistic neural networks are discussed in detail. Neural network structures for non-parametric pattern classification and system identification are developed.

1. Introduction

If the process is characterized by the absence of *a priori* information, the most popular methodology for identification and pattern classification is based on non-parametric approach. Such techniques — derived from non-parametric estimates of probability density and regression functions — have been developed by many authors to classify and identify different types of systems (see e.g. Gałkowski and Rutkowski, 1985; 1986; Rutkowski, 1988; 1991; 1993; Rutkowski and Rafajłowicz, 1989). The purpose of this article is to propose neural network structures for implementation of non-parametric algorithms. We shall define a neural network structure as a collection of parallel processors connected together in the form of a directed graph, so organized that the network structure corresponds to the non-parametric problem being considered. We shall use the so-called probabilistic neural networks (Specht, 1990). It should be emphasized that our neural networks do not require learning phase and are asymptotically optimal.

2. Density Estimation

Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be a sequence of independent, identically distributed random variables taking values in \mathbb{R}^d and having a probability density function f . The Parzen-Rosenblatt estimate of f is given by the formula

$$\hat{f}_n(\mathbf{x}) = \frac{1}{nh_n^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{X}_i}{h_n}\right) \quad (1)$$

where K is an appropriately chosen function fulfilling the conditions:

$$\sup_y |K(\mathbf{y})| < \infty \quad (2)$$

$$\int_{\mathbb{R}^d} |K(\mathbf{y})| d\mathbf{y} < \infty \quad (3)$$

* Institute of Electronics and Control Systems, Technical University of Częstochowa, 42–200 Częstochowa, Poland

$$\lim_{\|\mathbf{y}\| \rightarrow \infty} \|\mathbf{y}\|^d |K(\mathbf{y})| = 0 \tag{4}$$

$$\int_{\mathbb{R}^d} K(\mathbf{y}) d\mathbf{y} = 1 \tag{5}$$

In eqn. (4) symbol $\|\cdot\|$ stands for the Euclidean vector norm. Sequence h_n is a function of n and satisfies the conditions

$$\lim_{n \rightarrow \infty} h_n = 0 \quad \text{and} \quad \lim_{n \rightarrow \infty} nh_n^d = \infty \tag{6}$$

We assume that the function K is of the form

$$K(\mathbf{x}) = (2\pi)^{-\frac{1}{2}d} e^{-\frac{1}{2}\|\mathbf{x}\|^2} \tag{7}$$

where $\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x}$. Then we can rewrite estimator (1) as follows:

$$\hat{f}_n(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} nh_n^d} \sum_{i=1}^n \exp\left(-\frac{(\mathbf{x} - \mathbf{X}_i)^T (\mathbf{x} - \mathbf{X}_i)}{2h_n^2}\right) \tag{8}$$

Observe that

$$\begin{aligned} (\mathbf{x} - \mathbf{X}_i)^T (\mathbf{x} - \mathbf{X}_i) &= -2 \left(x^{(1)} X_i^{(1)} + x^{(2)} X_i^{(2)} + \dots + x^{(d)} X_i^{(d)} \right) \\ &+ \left(x^{(1)} \right)^2 + \left(x^{(2)} \right)^2 + \dots + \left(x^{(d)} \right)^2 \\ &+ \left(X_i^{(1)} \right)^2 + \left(X_i^{(2)} \right)^2 + \dots + \left(X_i^{(d)} \right)^2 \end{aligned} \tag{9}$$

Now assuming normalization of the vectors \mathbf{x} and \mathbf{X}_i formula (8) simplifies to

$$\hat{f}_n(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} nh_n^d} \sum_{i=1}^n \exp\left(-\frac{(1 - \mathbf{x}^T \mathbf{X}_i)}{h_n^2}\right) \tag{10}$$

Figure 1 shows a neural realization of algorithm (10). The proposed net has d inputs and two layers. The first layer consists of n neurons and each neuron has d weights. The output layer has a single neuron with the linear activation function. We should emphasize that the proposed network does not require a training procedure (optimal choosing of connection weights). The succeeding coordinates of the observation vectors $\mathbf{X}_i, i = 1, \dots, n$ play the role of the weights.

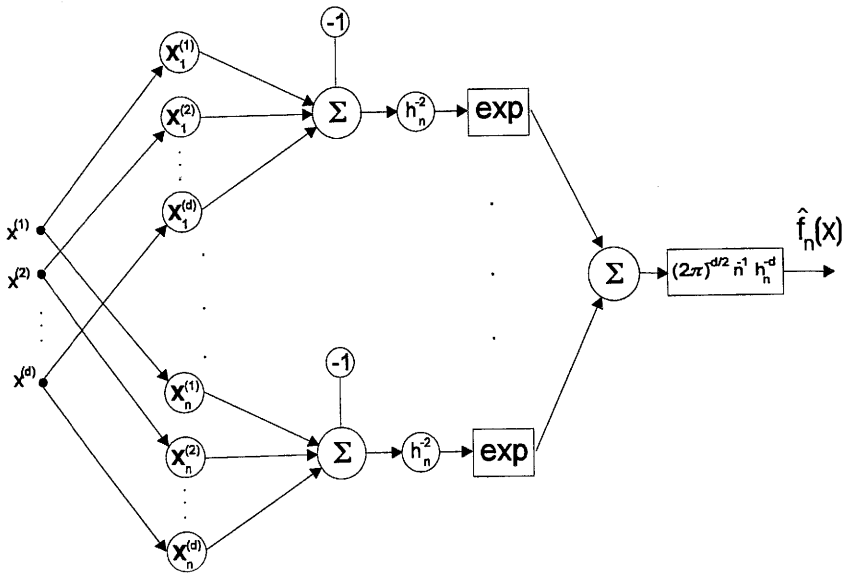


Fig. 1. Neural network density estimation.

The following theorem is a consequence of results presented by Parzen (1962) and Cacoullos (1966).

Theorem 1. *If the number of neurons in the first layer of the probability neural network presented in Fig. 1 is chosen to satisfy conditions (6), then*

$$E \left[\hat{f}_n(\mathbf{x}) - f_n(\mathbf{x}) \right]^2 \xrightarrow{n} 0 \tag{11}$$

in the points at which f is continuous. ■

3. Pattern Classification

Let $(\mathbf{X}, \mathbf{Y}), (\mathbf{X}_1, \mathbf{Y}_1), (\mathbf{X}_2, \mathbf{Y}_2), \dots$, be a sequence of i.i.d. pairs of random variables; \mathbf{Y} takes the values from the set; $S = \{1, \dots, M\}$, whereas \mathbf{X} takes the values in \mathbb{R}^d . The problem is to estimate \mathbf{Y} from \mathbf{X} and \mathbf{V}_n , where $\mathbf{V}_n = (\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_n, \mathbf{Y}_n)$ is a learning sequence. Suppose that p_m and f_m , $m = 1, \dots, M$, are the prior class probabilities and the class conditional densities, respectively. Define

$$T_{im} = \begin{cases} 1 & \text{if } \mathbf{Y}_i = m \\ 0 & \text{if } \mathbf{Y}_i \neq m \end{cases}$$

for $i = 1, 2, \dots, n$ and $m = 1, 2, \dots, M$. The Bayes discriminate function is given by

$$g_m(\mathbf{x}) = f(\mathbf{x}) E[T_{nm} | \mathbf{X}_n = \mathbf{x}] \tag{12}$$

where $f(\mathbf{x}) = \sum_{m=1}^M p_m f_m(\mathbf{x})$.

We consider a procedure of classifying every \mathbf{x} to a class m , $m \in S$, which maximizes $\hat{g}_{nm}(\mathbf{x})$, where $\hat{g}_{nm}(\mathbf{x})$ is the following estimate of the Bayes discriminate function

$$\hat{g}_{nm}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} nh_n^d} \sum_{i=1}^n T_{im} \exp\left(-\frac{1 - \mathbf{x}^T \mathbf{X}_i}{h_n^2}\right) \quad (13)$$

Figure 2 shows the neural network classifying pattern $\mathbf{x} \in \mathbb{R}^d$ to the class m , $m \in S = \{1, 2\}$, for which the expression

$$\tilde{g}_{nm}(\mathbf{x}) = \sum_{i=1}^n T_{im} \exp\left(-\frac{1 - \mathbf{x}^T \mathbf{X}_i}{h_n^2}\right)$$

takes the maximum value.

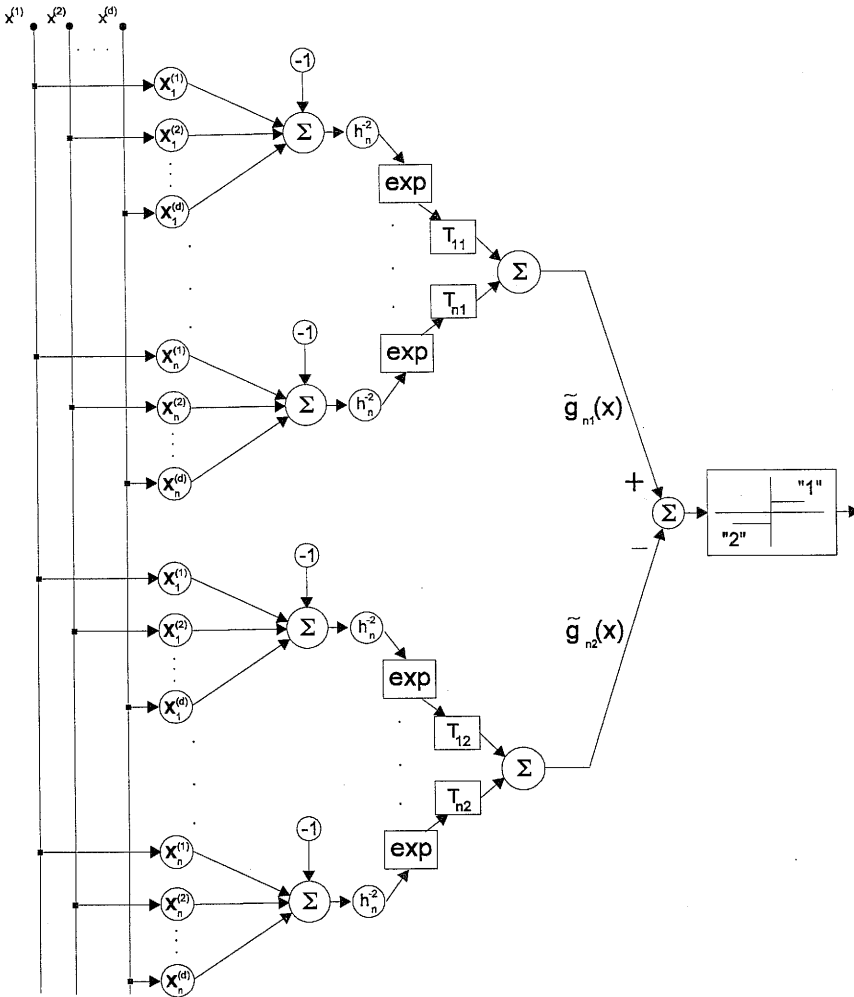


Fig. 2. Neural network pattern classification.

Let \hat{Y}_n be a decision obtained according to the classification procedure defined above. Let $D_n = P(\hat{Y}_n \neq Y | V_n)$. The procedure is said to be weakly (strongly) Bayes risk consistent if $D_n \xrightarrow{n} D_0$ in probability (with probability one), where D_0 is the Bayes probability error. The result below follows from the theorem presented by Greblicki *et al.* (1984).

Theorem 2. *If the number of neurons of the probabilistic neural network presented in Fig. 2 satisfies conditions (6), then $D_n \xrightarrow{n} D_0$ in probability. ■*

4. Neural Realization of the Recurrent Non-parametric Estimation Algorithms

The following recurrent estimate of the probability density function was first proposed by Wolverton and Wagner (1969)

$$\hat{f}_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_i^d} K\left(\frac{\mathbf{x} - \mathbf{X}_i}{h_i}\right) \tag{14}$$

Expression (14) can be rewritten as

$$\hat{f}_n(\mathbf{x}) = \frac{n-1}{n} f_{n-1}(\mathbf{x}) + \frac{1}{nh_n^d} K\left(\frac{\mathbf{x} - \mathbf{X}_n}{h_n}\right) \tag{15}$$

For the kernel K of form (7), using the same arguments as those in Section 2, one gets

$$\hat{f}_n(\mathbf{x}) = \frac{n-1}{n} f_{n-1}(\mathbf{x}) + \frac{1}{(2\pi)^{d/2} nh_n^d} \exp\left(-\frac{1 - \mathbf{x}^T \mathbf{X}_n}{h_n^2}\right) \tag{16}$$

Figure 3 shows the neural network performing algorithm (16) in a fixed point $\mathbf{x} \in \mathbb{R}^d$. The net consists of one neuron in the first layer having d inputs — coordinates of the vector \mathbf{X}_n , $n = 1, 2, \dots$. Let us notice that the role of weights is played by the coordinates of the vector \mathbf{x} . The second layer also consists of only one neuron with the feedback typical for recurrent neural networks. When the density is to be estimated at several points $\mathbf{x}_1, \dots, \mathbf{x}_L$, then the proposed structure should be copied L times. We shall obtain the neural network of L neurons processing the input observations in the parallel way.

Observe that the recurrent version of algorithm (13) takes the form

$$\hat{g}_{n,m}(\mathbf{x}) = \frac{n-1}{n} \hat{g}_{n-1,m}(\mathbf{x}) + \frac{1}{(2\pi)^{d/2} nh_n^d} T_{n,m} \exp\left(-\frac{1 - \mathbf{x}^T \mathbf{X}_n}{h_n^2}\right) \tag{17}$$

The corresponding recurrent neural network for pattern classification is shown in Fig. 4.

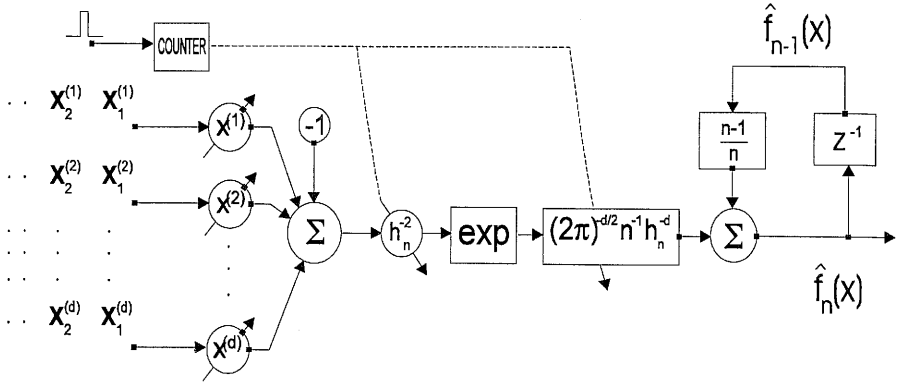


Fig. 3. Neural network density estimation — recurrent algorithm.

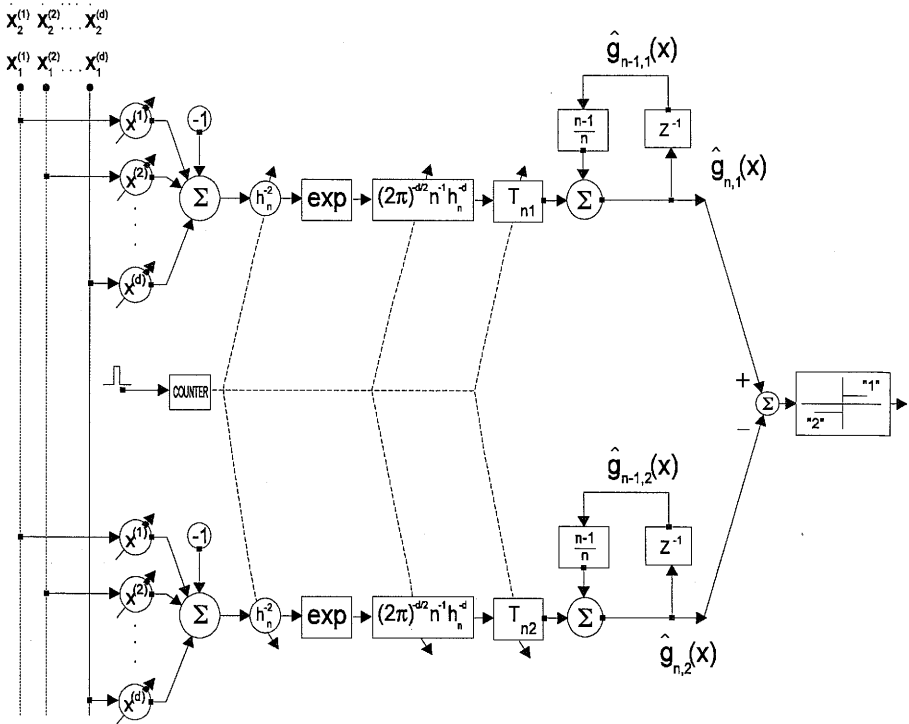


Fig. 4. Neural network pattern classification — recurrent algorithm.

5. Identification by Regression Function Estimation

5.1. Stochastic Input Signal

Let (X, Y) be a pair of random variables. X takes values in \mathbb{R}^d , whereas Y takes values in \mathbb{R} . Let f be the marginal Lebesgue density of X . Based on the sample

$(\mathbf{X}_1, \mathbf{Y}_1), \dots, (\mathbf{X}_n, \mathbf{Y}_n)$ of independent observations of (\mathbf{X}, \mathbf{Y}) , we wish to estimate the regression \mathbb{R} of \mathbf{Y} on \mathbf{X} , i.e. $\mathbb{R}(\mathbf{x}) = E[\mathbf{Y}|\mathbf{X} = \mathbf{x}]$. The probabilistic neural network that provides estimates of function \mathbb{R} and converges to the underlying (linear or non-linear) regression surface is presented in Fig. 4. In system identification the independent variable \mathbf{X} is the system input and the dependent variable \mathbf{Y} is the system output. Assume that $k(\mathbf{x}, y)$ represents the joint probability density function of a vector random variable \mathbf{X} and a scalar random variable \mathbf{Y} . The conditional mean of \mathbf{Y} given \mathbf{x} (particular measured value of random variable \mathbf{X}) is given by

$$\mathbb{R}(\mathbf{x}) = E[\mathbf{Y}|\mathbf{X} = \mathbf{x}] = \frac{1}{f(\mathbf{x})} \int_{-\infty}^{+\infty} yk(\mathbf{x}, y) dy \tag{18}$$

If the probability density functions f and k are unknown, then they should be estimated from a sample of observations of \mathbf{X} and \mathbf{Y} . The joint probability density estimator derived from the Gaussian kernel takes the form:

$$\hat{k}_n(\mathbf{x}, y) = \frac{1}{(2\pi)^{(d+1)/2} nh_n^{d+1}} \sum_{i=1}^n \exp\left(-\frac{(1 - \mathbf{x}^T \mathbf{X}_i)}{h_n^2}\right) \exp\left(-\frac{(y - \mathbf{Y}_i)^2}{2h_n^2}\right) \tag{19}$$

Replacing $k(\mathbf{x}, y)$ by $\hat{k}_n(\mathbf{x}, y)$ and $f(\mathbf{x})$ by $\hat{f}_n(\mathbf{x})$ in (18) and performing the indicated integration we get the following estimate of the regression function:

$$\hat{\mathbb{R}}_n(\mathbf{x}) = \frac{\sum_{i=1}^n \mathbf{Y}_i \exp\left(\frac{1 - \mathbf{x}^T \mathbf{X}_i}{h_n^2}\right)}{\sum_{i=1}^n \exp\left(\frac{1 - \mathbf{x}^T \mathbf{X}_i}{h_n^2}\right)} \tag{20}$$

Figure 5 shows a neural network realization of algorithm (20).

One may easily derive the neural network structure for recurrent regression function estimation. The corresponding net is shown in Fig. 6.

Theorem 3. *If $E|\mathbf{Y}| < \infty$ and the number of neurons of the probabilistic neural network presented in Fig. 5 satisfies conditions (6), then $\hat{\mathbb{R}}_n \xrightarrow{n} \mathbb{R}$ in probability for almost all $\mathbf{x} \in \mathbb{R}^d$.* ■

This result follows the theorem given by Greblicki *et al.* (1984).

5.2. Deterministic Input Signal

Many engineering problems are concerned with systems described by the following equation:

$$y_i = \mathbb{R}(x_i) + Z_i, \quad i = 1, \dots, n$$

relating the input x_i and output y_i , and the measurement noise Z_i .

Consider the d -dimensional unit cube $Q = [0, 1]^d$. Let $n^{1/d}$ be an integer, and $i_j = 1, \dots, n^{1/d}$, $j = 1, \dots, d$. Partition the interval $[0, 1]$ on each axis into $n^{1/d}$ subsets $\Delta x_{j,i_j}$. Define the following Cartesian product

$$\Delta x_{1,i_1} \times \Delta x_{2,i_2} \times \dots \times \Delta x_{d,i_d} = Q_{d,i}$$

Let $Q_{d,i} \cap Q_{d,j} = \emptyset$, for $i \neq j$ and $\cup Q_{d,i} = Q_d$. The inputs x_i are selected so that $x_i \in Q_{d,i}$.

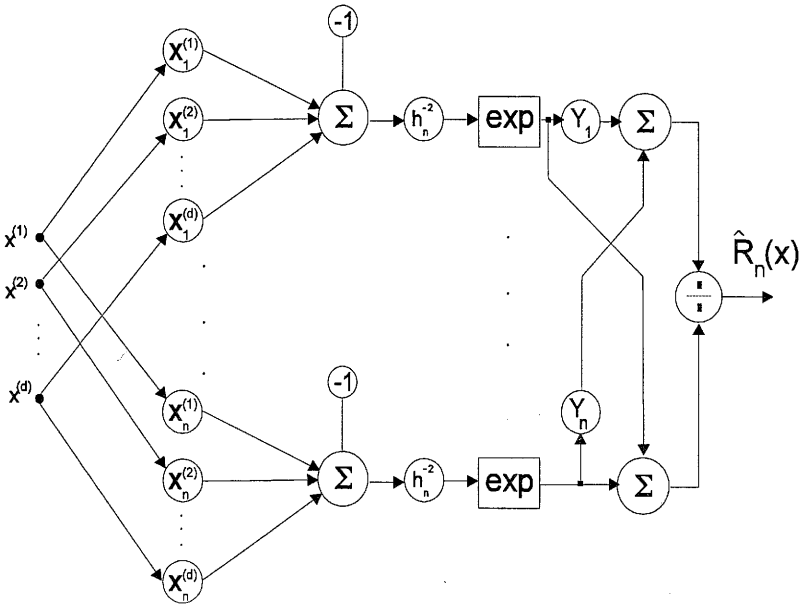


Fig. 5. Neural regression function estimation — stochastic input signal.

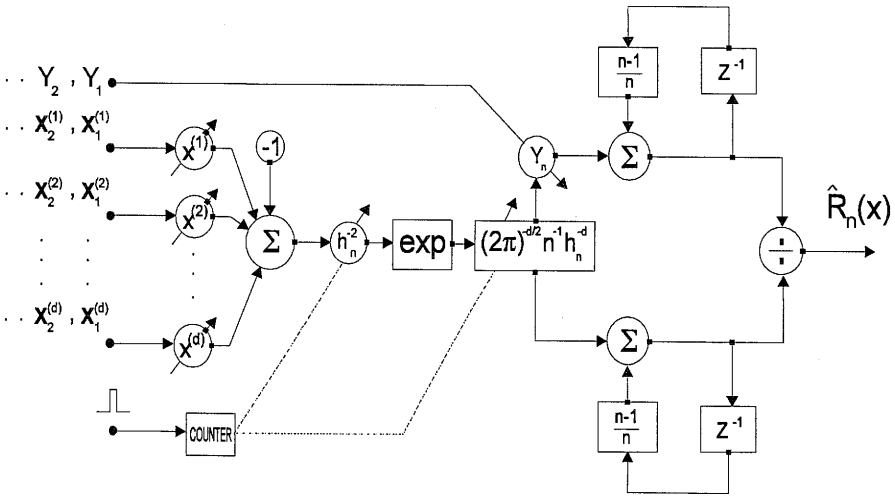


Fig. 6. Recurrent neural regression estimation.

We propose the following algorithm

$$\hat{R}_n(x) = \frac{1}{(2\pi)^{d/2} n h_n^d} \sum_{i=1}^n Y_i \exp\left(-\frac{(1 - x^T X_i)}{h_n^2}\right) \quad (21)$$

The appropriate net is shown in Fig. 7.

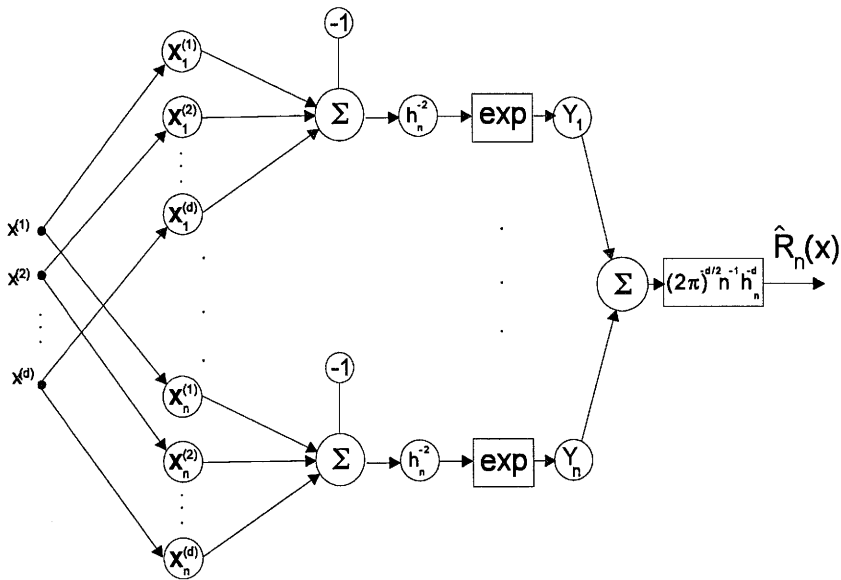


Fig. 7. Neural regression function estimation — deterministic input signal.

Theorem 4. *If the number of neurons of the probabilistic neural network presented in Fig. 7 satisfies conditions (6) and $\sup_{1 \leq i \leq n} \sup_{x, y \in Q_i} \|x - y\| = O(n^{-1/d})$, then $\widehat{\mathbb{R}}_n \xrightarrow{n} \mathbb{R}$ in probability at all continuity points x of \mathbb{R} .* ■

The above theorem follows from theorems presented by Gałkowski and Rutkowski (1985) and Georgiev (1990).

References

Cacoullos T. (1966): *Estimation of multivariate density.* — Annals of Institute of Statistical Mathematics, v.18, pp.179–189.

Gałkowski T. and Rutkowski L. (1985): *Non-parametric recovery of multivariate functions with applications to system identification.* — Proc. IEEE, v.73, pp.942–943.

Gałkowski T. and Rutkowski L. (1986): *Non-parametric fitting of multivariate functions.* — IEEE Trans. Automat. Contr. v.AC-31, No.8, pp.785–787.

Gałkowski T. and Rutkowski L. (1994): *Realizacja algorytmów rozpoznawania i nieparametrycznej estymacji za pomocą sieci neuronowych.* — Proc. 1st Nat. Conf. Neural Networks and Their Applications, Kule – Częstochowa, (Poland), pp.237–242, (in Polish).

Georgiev A.A. (1990): *Non-parametric multiple functions fitting.* — Statistics and Probability Letters, No.3.

- Greblicki W., Krzyżak A. and Pawlak M. (1984): *Distribution-free pointwise consistency of kernel regression estimate*. — The Annals of Statistics, v.12, No.4, pp.1570–1575.
- Parzen E. (1962): *On estimation of a probability density function and mode*. — Annals of Math. Statistics, v.33, pp.1065–1076.
- Rosenblatt M. (1956): *Remarks on some estimates of a density function*. — Annals of Math. Statistics, v.27, pp.823–837.
- Rutkowski L. (1988): *Sequential pattern recognition procedures derived from multiple Fourier series*. — Pattern Recognition Letters, No.8, pp.213–216.
- Rutkowski L. and Rafajłowicz E. (1989): *On global rate of convergence of some non-parametric identification procedures*. — IEEE Trans. Automat. Contr. v.AC-34, No.10, pp.1089–1091.
- Rutkowski L. (1991): *Identification of MISO non-linear regressions in the presence of a wide class of disturbances*. — IEEE Trans. Information Theory, v.IT-37, pp.214–216.
- Rutkowski L. (1993): *Multiple Fourier series procedures of non-linear regressions from noisy data*. — IEEE Trans. Signal Processing, v.41, No.10, pp.3062–3065.
- Specht D.F. (1990): *Probabilistic neural networks*. — Neural Networks, v.3, pp.109–118.
- Wolverton C.T. and Wagner T.J. (1969): *Recursive estimates of probability densities*. — IEEE Trans. SSC, v.SSC-5, No.3