



INTERNATIONAL ATOMIC ENERGY AGENCY  
UNITED NATIONS EDUCATIONAL, SCIENTIFIC AND CULTURAL ORGANIZATION  
**INTERNATIONAL CENTRE FOR THEORETICAL PHYSICS**  
I.C.T.P., P.O. BOX 586, 34100 TRIESTE, ITALY, CABLE: CENTRATOM TRIESTE



**SMR.853 - 71**

**ANTONIO BORSELLINO COLLEGE ON NEUROPHYSICS**

**(15 May - 9 June 1995)**

---

**"Regularization Theory and Neural Networks Architectures"**

**Federico Girosi**  
**Center for Biological and Computational Learning**  
**Massachusetts Institute of Technology**  
**Cambridge, MA 02139**  
**U.S.A.**

---

**These are preliminary lecture notes, intended only for distribution to participants.**

## Regularization Theory and Neural Networks Architectures

Federico Girosi

Michael Jones

Tomaso Poggio

*Center for Biological and Computational Learning,*

*Department of Brain and Cognitive Sciences and*

*Artificial Intelligence Laboratory,*

*Massachusetts Institute of Technology, Cambridge, MA 02139 USA*

We had previously shown that regularization principles lead to approximation schemes that are equivalent to networks with one layer of hidden units, called *regularization networks*. In particular, standard smoothness functionals lead to a subclass of regularization networks, the well known radial basis functions approximation schemes. This paper shows that regularization networks encompass a much broader range of approximation schemes, including many of the popular general additive models and some of the neural networks. In particular, we introduce new classes of smoothness functionals that lead to different classes of basis functions. Additive splines as well as some tensor product splines can be obtained from appropriate classes of smoothness functionals. Furthermore, the same generalization that extends radial basis functions (RBF) to hyper basis functions (HBF) also leads from additive models to ridge approximation models, containing as special cases Breiman's hinge functions, some forms of projection pursuit regression, and several types of neural networks. We propose to use the term *generalized regularization networks* for this broad class of approximation schemes that follow from an extension of regularization. In the probabilistic interpretation of regularization, the different classes of basis functions correspond to different classes of prior probabilities on the approximating function spaces, and therefore to different types of smoothness assumptions. In summary, different multilayer networks with one hidden layer, which we collectively call generalized regularization networks, correspond to different classes of priors and associated smoothness functionals in a classical regularization principle. Three broad classes are (1) radial basis functions that can be generalized to hyper basis functions, (2) some tensor product splines, and (3) additive splines that can be generalized to schemes of the type of ridge approximation, hinge functions, and several perceptron-like neural networks with one hidden layer.

## 1 Introduction

---

In recent years we and others have argued that the task of learning from examples can be considered in many cases to be equivalent to multivariate function approximation, that is, to the problem of approximating a smooth function from sparse data, the examples. The interpretation of an approximation scheme in terms of networks and vice versa has also been extensively discussed (Barron and Barron 1988; Poggio and Girosi 1989, 1990a,b; Girosi 1992; Broomhead and Lowe 1988; Moody and Darken 1988, 1989; White 1989, 1990; Ripley 1994; Omohundro 1987; Kohonen 1990; Lapedes and Farber 1988; Rumelhart *et al.* 1986; Hertz *et al.* 1991; Kung 1993; Sejnowski and Rosenberg 1987; Hurlbert and Poggio 1988; Poggio 1975).

In a series of papers we have explored a quite general approach to the problem of function approximation. The approach regularizes the ill-posed problem of function approximation from sparse data by assuming an appropriate prior on the class of approximating functions. Regularization techniques (Tikhonov 1963; Tikhonov and Arsenin 1977; Morozov 1984; Bertero 1986; Wahba 1975, 1979, 1990) typically impose smoothness constraints on the approximating set of functions. It can be argued that some form of smoothness is necessary to allow meaningful generalization in approximation type problems (Poggio and Girosi 1989, 1990). A similar argument can also be used (see Section 9.1) in the case of classification where smoothness is a condition on the classification boundaries rather than on the input-output mapping itself. Our use of regularization, which follows the classical technique introduced by Tikhonov, identifies the approximating function as the minimizer of a cost functional that includes an *error term* and a smoothness functional, usually called a *stabilizer*. In the Bayesian interpretation of regularization (see Kimeldorf and Wahba 1971; Wahba 1990; Bertero *et al.* 1988; Marroquin *et al.* 1987; Poggio *et al.* 1985) the stabilizer corresponds to a smoothness prior, and the error term to a model of the noise in the data (usually gaussian and additive).

In Poggio and Girosi (1989, 1990) and Girosi (1992) we showed that regularization principles lead to approximation schemes that are equivalent to networks with one "hidden" layer, which we call *regularization networks* (RN). In particular, we described how a certain class of radial stabilizers—and the associated priors in the equivalent Bayesian formulation—lead to a subclass of regularization networks, the already-known radial basis functions (Powell 1987, 1992; Franke 1982, 1987; Micchelli 1986; Kansa 1990a,b; Madych and Nelson 1990a,b; Dyn 1987, 1991; Hardy 1971, 1990; Buhmann 1990; Lancaster and Salkauskas 1986; Broomhead and Lowe 1988; Moody and Darken 1988, 1989; Poggio and Girosi 1990; Girosi 1992). The regularization networks with radial stabilizers we studied include many classical one-dimensional (Schumaker 1981; de Boor 1978) as well as multidimensional splines and approximation tech-

niques, such as radial and nonradial gaussian, thin-plate splines (Duchon 1977; Meinguet 1979; Grimson 1982; Cox 1984; Eubank 1988) and multi-quadratic functions (Hardy 1971, 1990). In Poggio and Girosi (1990a,b) we extended this class of networks to Hyper Basis Functions (HBF). In this paper we show that an extension of regularization networks, which we propose to call *Generalized Regularization Networks* (GRN), encompasses an even broader range of approximation schemes including, in addition to HBF, tensor product splines, many of the general additive models, and some of the neural networks. As expected, GRN have approximation properties of the same type as already shown for some of the neural networks (Girosi and Poggio 1990a; Cybenko 1989; Hornik *et al.* 1989; White 1990; Irie and Miyake 1988; Funahashi 1989; Barron 1991, 1994; Jones 1992; Mhaskar and Micchelli 1992, 1993; Mhaskar 1993a,b).

The plan of the paper is as follows. We first discuss the solution of the variational problem of regularization. We then introduce three different classes of stabilizers—and the corresponding priors in the equivalent Bayesian interpretation—that lead to different classes of basis functions: the well-known radial stabilizers, tensor-product stabilizers, and the new additive stabilizers that underlie additive splines of different types. It is then possible to show that the same argument that extends radial basis functions to hyper basis functions also leads from additive models to some ridge approximation schemes, defined as

$$f(\mathbf{x}) = \sum_{\mu=1}^K h_{\mu}(\mathbf{w}_{\mu} \cdot \mathbf{x})$$

where  $h_{\mu}$  are appropriate one-dimensional functions.

Special cases of ridge approximation are Breiman's hinge functions (1993), projection pursuit regression (PPR) (Friedman and Stuetzle 1981; Huber 1985; Diaconis and Freedman 1984; Donoho and Johnstone 1989; Moody and Yarvin 1991), and multilayer perceptrons (Lapedes and Farber 1988; Rumelhart *et al.* 1986; Hertz *et al.* 1991; Kung 1993; Sejnowski and Rosenberg 1987). Simple numerical experiments are then described to illustrate the theoretical arguments.

In summary, the chain of our arguments shows that some ridge approximation schemes are approximations of regularization networks with appropriate additive stabilizers. The form of  $h_{\mu}$  depends on the stabilizer, and includes in particular cubic splines (used in typical implementations of PPR) and one-dimensional gaussians. Perceptron-like neural networks with one hidden layer and with a gaussian activation function are included. It seems impossible, however, to directly derive from regularization principles the sigmoidal activation functions typically used in feedforward neural networks. We discuss, however, in a simple example, the close relationship between basis functions of the hinge, the sigmoid and the gaussian type.

The appendices deal with observations related to the main results of the paper and more technical details.

## 2 The Regularization Approach to the Approximation Problem

Suppose that the set  $g = \{(x_i, y_i) \in R^d \times R\}_{i=1}^N$  of data has been obtained by random sampling a function  $f$ , belonging to some space of functions  $X$  defined on  $R^d$ , in the presence of noise, and suppose we are interested in recovering the function  $f$ , or an estimate of it, from the set of data  $g$ . This problem is clearly ill-posed, since it has an infinite number of solutions. To choose one particular solution we need to have some a priori knowledge of the function that has to be reconstructed. The most common form of a priori knowledge consists in assuming that the function is *smooth*, in the sense that two similar inputs correspond to two similar outputs. The main idea underlying regularization theory is that the solution of an ill-posed problem can be obtained from a variational principle, which contains both the data and prior smoothness information. Smoothness is taken into account by defining a *smoothness functional*  $\phi[f]$  in such a way that lower values of the functional correspond to smoother functions. Since we look for a function that is simultaneously close to the data and also smooth, it is natural to choose as a solution of the approximation problem the function that minimizes the following functional:

$$H[f] = \sum_{i=1}^N [f(x_i) - y_i]^2 + \lambda \phi[f] \quad (2.1)$$

where  $\lambda$  is a positive number that is usually called the *regularization parameter*. The first term is enforcing closeness to the data, and the second smoothness, while the regularization parameter controls the trade-off between these two terms, and can be chosen according to cross-validation techniques (Allen 1974; Wahba and Wold 1975; Golub *et al.* 1979; Craven and Wahba 1979; Utreras 1979; Wahba 1985) or to some other principle, such as structural risk minimization (Vapnik 1988).

It can be shown that for a wide class of functionals  $\phi$ , the solutions of the minimization of the functional (2.1) all have the same form. Although a detailed and rigorous derivation of the solution of this problem is out of the scope of this paper, a simple derivation of this general result is presented in Appendix A. In this section we just present a family of smoothness functionals and the corresponding solutions of the variational problem. We refer the reader to the current literature for the mathematical details (Wahba 1990; Madych and Nelson 1990a; Dyn 1987).

We first need to give a more precise definition of what we mean by smoothness and define a class of suitable smoothness functionals. We refer to smoothness as a measure of the "oscillatory" behavior of a function. Therefore, within a class of differentiable functions, one function will be said to be smoother than another one if it oscillates less. If we look at the functions in the frequency domain, we may say that a function is smoother than another one if it has less energy at high frequency (smaller bandwidth). The high frequency content of a function can be

measured by first high-pass filtering the function, and then measuring the power, that is the  $L_2$  norm, of the result. In formulas, this suggests defining smoothness functionals of the form

$$\phi[f] = \int_{R^d} ds \frac{|\tilde{f}(s)|^2}{\tilde{G}(s)} \quad (2.2)$$

where the tilde indicates the Fourier transform,  $\tilde{G}$  is some positive function that tends to zero as  $\|s\| \rightarrow \infty$  (so that  $1/\tilde{G}$  is an high-pass filter) and for which the class of functions such that this expression is well defined is not empty. For a well defined class of functions  $G$  (Madych and Nelson 1990a; Dyn 1991; Dyn *et al.* 1989) this functional is a seminorm, with a finite dimensional null space  $\mathcal{N}$ . The next section will be devoted to giving examples of the possible choices for the stabilizer  $\phi$ . For the moment we just assume that it can be written as in equation 2.2, and make the additional assumption that  $\tilde{G}$  is symmetric, so that its Fourier transform  $G$  is real and symmetric. In this case it is possible to show (see Appendix A for a sketch of the proof) that the function that minimizes the functional (2.1) has the form

$$f(\mathbf{x}) = \sum_{i=1}^N c_i G(\mathbf{x} - \mathbf{x}_i) + \sum_{\alpha=1}^k d_{\alpha} \psi_{\alpha}(\mathbf{x}) \quad (2.3)$$

where  $\{\psi_{\alpha}\}_{\alpha=1}^k$  is a basis in the  $k$ -dimensional null space  $\mathcal{N}$  of the functional  $\phi$ , that in most cases is a set of polynomials, and therefore will be referred to as the "polynomial term" in equation 2.3. The coefficients  $d_{\alpha}$  and  $c_i$  depend on the data, and satisfy the following linear system:

$$(G + \lambda I)\mathbf{c} + \Psi^T \mathbf{d} = \mathbf{y} \quad (2.4)$$

$$\Psi \mathbf{c} = 0 \quad (2.5)$$

where  $I$  is the identity matrix, and we have defined

$$\begin{aligned} (\mathbf{y})_i &= y_i, & (\mathbf{c})_i &= c_i, & (\mathbf{d})_i &= d_i \\ (G)_{ij} &= G(\mathbf{x}_i - \mathbf{x}_j), & (\Psi)_{\alpha i} &= \psi_{\alpha}(\mathbf{x}_i) \end{aligned}$$

Notice that if the data term in equation 2.1 is replaced by  $\sum_{i=1}^N V[f(\mathbf{x}_i) - y_i]$  where  $V$  is any differentiable function, the solution of the variational principle has still the form 2.3, but the coefficients cannot be found any more by solving a linear system of equations (Girosi 1991; Girosi *et al.* 1991).

The existence of a solution to the linear system shown above is guaranteed by the existence of the solution of the variational problem. The case of  $\lambda = 0$  corresponds to pure interpolation. In this case the existence of an exact solution of the linear system of equations depends on the properties of the basis function  $G$  (Micchelli 1986).

The approximation scheme of equation 2.3 has a simple interpretation in terms of a network with one layer of hidden units, which we call a *Regularization Network* (RN). Appendix B describes the extension to the vector output scheme.

In summary, the argument of this section shows that using a regularization network of the form 2.3, for a certain class of basis functions  $G$ , is equivalent to minimizing the functional 2.1. In particular, the choice of  $G$  is equivalent to the corresponding choice of the smoothness functional 2.2.

**2.1 Dual Representation of Regularization Networks.** Consider an approximating function of the form 2.3, neglecting the “polynomial term” for simplicity. A compact notation for this expression is

$$f(\mathbf{x}) = \mathbf{c} \cdot \mathbf{g}(\mathbf{x}) \quad (2.6)$$

where  $\mathbf{g}(\mathbf{x})$  is the vector of functions such that  $[\mathbf{g}(\mathbf{x})]_i = G(\mathbf{x} - \mathbf{x}_i)$ . Since the coefficients  $\mathbf{c}$  satisfy the linear system 2.4, solution 2.6 becomes

$$f(\mathbf{x}) = (G + \lambda I)^{-1} \mathbf{y} \cdot \mathbf{g}(\mathbf{x})$$

We can rewrite this expression as

$$f(\mathbf{x}) = \sum_{i=1}^N y_i b_i(\mathbf{x}) = \mathbf{y} \cdot \mathbf{b}(\mathbf{x}) \quad (2.7)$$

in which the vector  $\mathbf{b}(\mathbf{x})$  of basis functions is defined

$$\mathbf{b}(\mathbf{x}) = (G + \lambda I)^{-1} \mathbf{g}(\mathbf{x}) \quad (2.8)$$

and now depends on all the data points and on the regularization parameter  $\lambda$ . The representation 2.7 of the solution of the approximation problem is known as the *dual* of equation 2.6, and the basis functions  $b_i(\mathbf{x})$  are called the *equivalent kernels*, because of the similarity between equation 2.7 and the kernel smoothing technique that we will define in Section 2.2 (Silverman 1984; Härdle 1990; Hastie and Tibshirani 1990). While in equation 2.6 the “difficult” part is the computation of the vector of coefficients  $\mathbf{c}$ , the set of basis functions  $\mathbf{g}(\mathbf{x})$  being easily built, in equation 2.7 the “difficult” part is the computation of the basis functions  $\mathbf{b}(\mathbf{x})$ , the coefficients of the expansion being explicitly given by the  $y_i$ . Notice that  $\mathbf{b}(\mathbf{x})$  depends on the distribution of the data in the input space and that the kernels  $b_i(\mathbf{x})$ , unlike the kernels  $G(\mathbf{x} - \mathbf{x}_i)$ , are not translated replicas of the same kernel. Notice also that, as shown in Appendix B, a dual representation of the form 2.7 exists for all the approximation schemes that consists of linear superpositions of arbitrary numbers of basis functions, as long as the error criterion that is used to determine the parameters of the approximation is quadratic.

The dual representation provides an intuitive way of looking at the approximation scheme 2.3: the value of the approximating function at an

evaluation point  $x$  is explicitly expressed as a weighted sum of the values  $y_i$  of the function at the examples  $x_i$ . This concept is not new in approximation theory, and has been used, for example, in the theory of quasi-interpolation. The case in which the data points  $\{x_i\}$  coincide with the multi-integers  $Z^d$ , where  $Z$  is the set of integers number, has been extensively studied in the literature, and it is also known as *Schoenberg's approximation*, (Schoenberg 1946a, 1969; Rabut 1991, 1992; Madych and Nelson 1990a; Jackson 1988; de Boor 1990; Buhmann 1990, 1991; Dyn *et al.* 1989). In this case, an approximation  $f^*$  to a function  $f$  is sought of the form

$$f^*(x) = \sum_{j \in Z^d} f(j) \psi(x - j) \quad (2.9)$$

where  $\psi$  is some fast-decaying function that is a linear combination of radial basis functions. The approximation scheme 2.9 is therefore a linear superposition of radial basis functions in which the functions  $\psi(x - j)$  play the role of equivalent kernels. Quasi-interpolation is interesting because it could provide good approximation without the need of solving complex minimization problems or solving large linear systems. For a discussion of such noniterative training algorithms see Mhaskar (1993b) and references therein.

Although difficult to prove rigorously, we can expect the kernels  $b_i(x)$  to decrease with the distance of the data points  $x_i$  from the evaluation point, so that only the neighboring points affect the estimate of the function at  $x$ , providing therefore a "local" approximation scheme. Even if the original basis function  $G$  is not "local," like the multiquadric  $G(x) = \sqrt{1 + \|x\|^2}$ , the basis functions  $b_i(x)$  are bell shaped, local functions, whose locality will depend on the choice of the basis function  $G$ , on the density of data points, and on the regularization parameter  $\lambda$ . This shows that apparently "global" approximation schemes can be regarded as local, *memory-based* techniques (see equation 2.7) (Mhaskar 1993b). It should be noted however, that these techniques do not have the highest possible degree of locality, since the parameter that controls the locality is the regularization parameter  $\lambda$ , that is the same for all the kernels. It is possible to devise even more local techniques, in which each kernel has a parameter that controls its locality (Bottou and Vapnik 1992; Vapnik, personal communication).

When the data are equally spaced on an infinite grid, we expect the basis functions  $b_i(x)$  to become translation invariant, and therefore the dual representation 2.7 becomes a convolution filter. For a study of the properties of these filters in the case of one-dimensional cubic splines see the work of Silverman (1984), who gives explicit results for the shape of the equivalent kernel.

Let us consider some simple experiments that show the shape of the equivalent kernels in specific situations. We first considered a data set composed of 36 equally spaced points on the domain  $[0, 1] \times [0, 1]$ , at the nodes of a regular grid with spacing equal to 0.2. We use the multiquadric

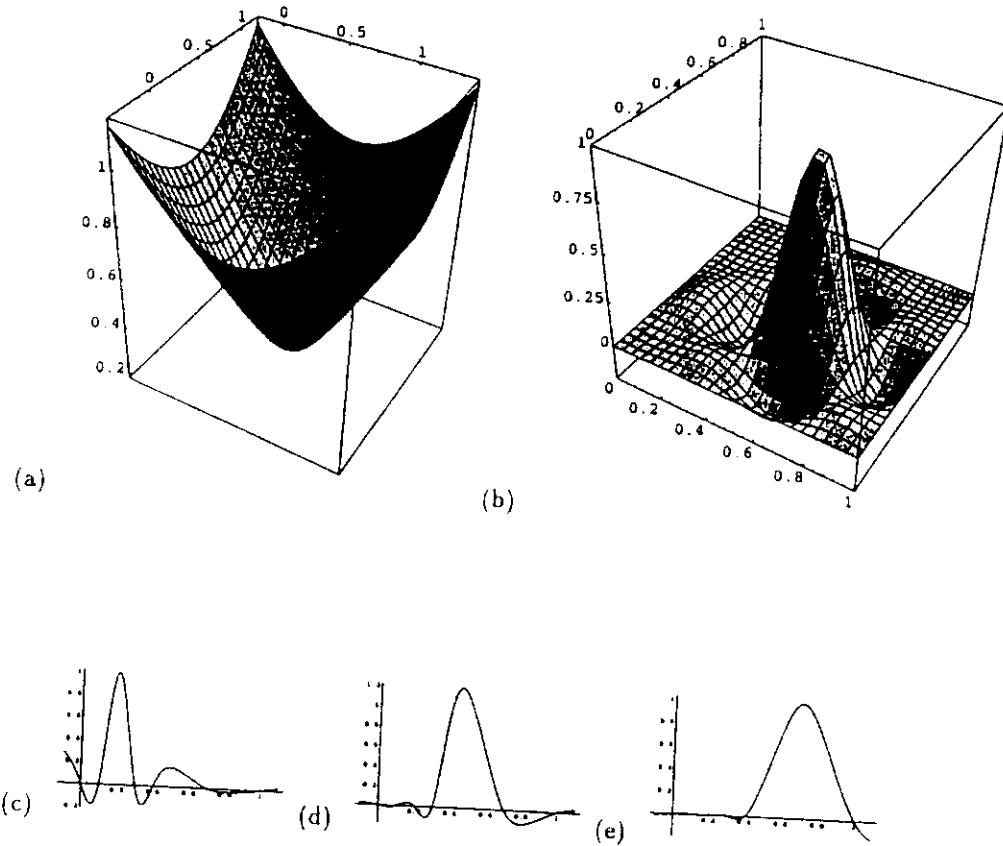


Figure 1: (a) The multiquadric function. (b) An equivalent kernel for the multiquadric basis function in the cases of two-dimensional equally spaced data. (c,d,e) The equivalent kernels  $b_3$ ,  $b_5$ , and  $b_6$ , for nonuniform one-dimensional multiquadric interpolation (see text for explanation).

basis functions  $G(\mathbf{x}) = \sqrt{\sigma^2 + \|\mathbf{x}\|^2}$ , where  $\sigma$  has been set to 0.2. Figure 1a shows the original multiquadric function, and Figure 1b the equivalent kernel  $b_{16}$ , in the case of  $\lambda = 0$ , where, according to definition 2.8

$$b_i(\mathbf{x}) = \sum_{j=1}^{36} (G^{-1})_{ij} G(\mathbf{x} - \mathbf{x}_j)$$

All the other kernels, except those close to the border, are very similar, since the data are equally spaced, and translation invariance holds approximately.

Consider now a one-dimensional example with a multiquadric basis function:

$$G(x) = \sqrt{\sigma^2 + x^2}$$

The data set was chosen to be a nonuniform sampling of the interval  $[0, 1]$ , that is the set

$$\{0.0, 0.1, 0.2, 0.3, 0.4, 0.7, 1.0\}$$

In Figure 1c, d, and e we have drawn, respectively, the equivalent kernels  $b_3$ ,  $b_5$ , and  $b_6$ , under the same definitions. Notice that all of them are bell-shaped, although the original basis function is an increasing, cup-shaped function. Notice, moreover, that the shape of the equivalent kernels changes from  $b_3$  to  $b_6$ , becoming broader in moving from a high to low sample density region. This phenomenon has been shown by Silverman (1984) for cubic splines, but we expect it to appear in much more general cases.

The connection between regularization theory and the dual representation 2.7 becomes clear in the special case of "continuous" data, for which the regularization functional has the form

$$H[f] = \int dx [f(x) - y(x)]^2 + \lambda \phi[f] \quad (2.10)$$

where  $y(x)$  is the function to be approximated. This functional can be intuitively seen as the limit of the functional 2.1 when the number of data points goes to infinity and their spacing is uniform. It is easily seen that, when the stabilizer  $\phi[f]$  is of the form 2.2, the solution of the regularization functional 2.10 is

$$f(x) = y(x) * B(x) \quad (2.11)$$

where  $B(x)$  is the Fourier transform of

$$\tilde{B}(s) = \frac{\tilde{G}(s)}{\lambda + \tilde{G}(s)}$$

[see Poggio *et al.* (1988) for some examples of  $B(x)$ ]. The solution 2.11 is therefore a filtered version of the original function  $y(x)$  and, consistently with the results of Silverman (1984), has the form 2.7, where the equivalent kernels are translates of the function  $B(x)$  defined above. Notice the effect of the regularization parameter: for  $\lambda = 0$  the equivalent kernel  $B(x)$  is a Dirac delta function, and  $f(x) = y(x)$  (no noise), while for  $\lambda \rightarrow \infty$  we have  $B(x) = G(x)/\lambda$  and  $f = G/\lambda * y$  (a low-pass filter).

The dual representation is illuminating and especially interesting for the case of a multi-output network—approximating a vector field—that is discussed in Appendix B.

**2.2 Normalized Kernels.** An approximation technique very similar to radial basis functions is the so-called *normalized Radial Basis Functions*

(Moody and Darken 1988, 1989). A normalized radial basis functions expansion is a function of the form

$$f(\mathbf{x}) = \frac{\sum_{\alpha=1}^n c_{\alpha} G(\mathbf{x} - \mathbf{t}_{\alpha})}{\sum_{\alpha=1}^n G(\mathbf{x} - \mathbf{t}_{\alpha})} \quad (2.12)$$

The only difference between equation 2.12 and radial basis functions is the normalization factor in the denominator, which is an estimate of the probability distribution of the data. A discussion about the relation between normalized gaussian basis function networks, gaussian mixtures, and gaussian mixture classifiers can be found in the work of Tresp *et al.* (1993). In the rest of this section we show that a particular version of this approximation scheme has again a tight connection to regularization theory.

Let  $P(\mathbf{x}, y)$  be the joint probability of inputs and outputs of the network, and let us assume that we have a sample of  $N$  pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$  randomly drawn according to  $P$ . Our goal is to build an estimator (a network)  $f$  that minimizes the *expected risk*:

$$I[f] = \int d\mathbf{x} dy P(\mathbf{x}, y) [y - f(\mathbf{x})]^2 \quad (2.13)$$

This cannot be done, since the probability  $P$  is unknown, and usually the *empirical risk*

$$I_{\text{emp}}[f] = \frac{1}{N} \sum_{i=1}^N [y_i - f(\mathbf{x}_i)]^2 \quad (2.14)$$

is minimized instead. An alternative consists in obtaining an approximation of the probability  $P(\mathbf{x}, y)$  first, and then in minimizing the expected risk. If this option is chosen, one could use the regularization approach to probability estimation (Vapnik and Stefanyuk 1978; Aidun and Vapnik 1989; Vapnik 1982) that leads to the well-known technique of *Parzen windows*. A Parzen window estimator  $P^*$  for the probability distribution of a set of data  $\{\mathbf{z}_i\}_{i=1}^N$  has the form

$$P^*(\mathbf{z}) = \frac{1}{Nh} \sum_{i=1}^N \Phi\left(\frac{\mathbf{z} - \mathbf{z}_i}{h}\right) \quad (2.15)$$

where  $\Phi$  is an appropriate kernel, for example a gaussian, whose  $L_1$  norm is 1, and where  $h$  is a positive parameter, that, for simplicity, we set to 1 from now on. If the joint probability  $P(\mathbf{x}, y)$  in the expected risk 2.13 is approximated with a Parzen window estimator  $P^*$ , we obtain an approximated expression for the expected risk,  $I^*[f]$ , that can be explicitly minimized. In order to show how this can be done, we notice that we need to approximate the probability distribution  $P(\mathbf{x}, y)$ , and therefore

the random variable  $z$  of equation 2.15 is  $z = (x, y)$ . Hence, we choose a kernel of the following form:<sup>1</sup>

$$\Phi(z) = K(\|x\|)K(y)$$

where  $K$  is a standard one-dimensional, symmetric kernel, like the gaussian. The Parzen window estimator to  $P(x, y)$  is therefore

$$P^*(x, y) = \frac{1}{N} \sum_{i=1}^N K(\|x - x_i\|) K(y - y_i) \quad (2.16)$$

An approximation to the expected risk is therefore obtained as

$$I^*[f] = \frac{1}{N} \sum_{i=1}^N \int dx dy K(\|x - x_i\|) K(y - y_i) [y - f(x)]^2$$

In order to find an analytical expression for the minimum of  $I^*[f]$  we impose the stationarity constraint:

$$\frac{\delta I^*[f]}{\delta f(s)} = 0$$

that leads to the following equation:

$$\sum_{i=1}^N \int dx dy K(\|x - x_i\|) K(y - y_i) [y - f(x)] \delta(x - s) = 0$$

Performing the integral over  $x$ , and using the fact that  $\|K\|_{L_1} = 1$  we obtain

$$f(x) = \frac{\sum_{i=1}^N K(\|x - x_i\|) \int dy K(y - y_i) y}{\sum_{i=1}^N K(\|x - x_i\|)}$$

Performing a change of variable in the integral of the previous expression and using the fact that the kernel  $K$  is symmetric, we finally conclude that the function that minimizes the approximated expected risk is

$$f(x) = \frac{\sum_{i=1}^N y_i K(\|x - x_i\|)}{\sum_{i=1}^N K(\|x - x_i\|)} \quad (2.17)$$

The right-hand side of the equation converges to  $f$  when the number of examples goes to infinity, provided that the scale factor  $h$  tends to zero at an appropriate rate. This form of approximation is known as *kernel regression*, or *Nadaraya-Watson estimator*, and it has been the subject of extensive study in the statistics community (Nadaraya 1964; Watson 1964; Rosenblatt 1971; Priestley and Chao 1972; Gasser and Müller 1985; Devroye and Wagner 1980). A similar derivation of equation 2.17 has been given by Specht (1991), but we should remark that this equation

<sup>1</sup>Any kernel of the form  $\Phi(z) = K(x, y)$  in which the function  $K$  is even in each of the variables  $x$  and  $y$  would lead to the same conclusions that we obtain for this choice.

is usually derived in a different way, within the framework of locally weighted regression, assuming a locally constant model (Hardle 1990) with a local weight function  $K$ .

Notice that this equation has the form of equation 2.12, in which the centers coincide with the examples, and the coefficients  $c_i$  are simply the values  $y_i$  of the function at the data points  $\mathbf{x}_i$ . On the other hand, the equation is an estimate of  $f$ , which is linear in the observations  $y_i$  and has therefore also the general form of equation 2.7.

The Parzen window estimator, and therefore expression 2.17, can be derived in the framework of regularization theory (Vapnik and Stefanyuk 1978; Aidun and Vapnik 1989; Vapnik 1982) under a smoothness assumption on the probability distribution that has to be estimated. This means that in order to derive equation 2.17, a smoothness assumption has to be made on the joint probability distribution  $P(\mathbf{x}, y)$ , rather than on the regression function as in 2.2.

### 3 Classes of Stabilizers

In the previous section we considered the class of stabilizers of the form

$$\phi[f] = \int_{R^d} d\mathbf{s} \frac{|\tilde{f}(\mathbf{s})|^2}{G(\mathbf{s})} \quad (3.1)$$

and we have seen that the solution of the minimization problem always has the same form. In this section we discuss three different types of stabilizers belonging to the class 3.1, corresponding to different properties of the basis functions  $G$ . Each of them corresponds to different a priori assumptions on the smoothness of the function that must be approximated.

**3.1 Radial Stabilizers.** Most of the commonly used stabilizers have radial symmetry, that is, they satisfy the following equation:

$$\phi[f(\mathbf{x})] = \phi[f(R\mathbf{x})]$$

for any rotation matrix  $R$ . This choice reflects the a priori assumption that all the variables have the same relevance, and that there are no privileged directions. Rotation invariant stabilizers correspond to radial basis function  $G(\|\mathbf{x}\|)$ . Much attention has been dedicated to this case, and the corresponding approximation technique is known as radial basis functions (Powell 1987, 1990; Franke 1982, 1987; Micchelli 1986; Kansa, 1990a,b; Madych and Nelson 1990a; Dyn 1987, 1991; Hardy 1971, 1990; Buhmann 1990; Lancaster and Salkauskas 1986; Broomhead and Lowe 1988; Moody and Darken 1988, 1989; Poggio and Girosi 1990; Girosi 1992). The class of admissible radial basis functions is the class of conditionally positive definite functions (Micchelli 1986) of any order, since it has been shown

(Madych and Nelson 1990a; Dyn 1991) that in this case the functional of equation 3.1 is a seminorm, and the associated variational problem is well defined. All the radial basis functions can therefore be derived in this framework. We explicitly give two important examples.

**3.1.1 Duchon Multidimensional Splines.** Duchon (1977) considered measures of smoothness of the form

$$\phi[f] = \int_{\mathbb{R}^d} ds \|s\|^{2m} |\tilde{f}(s)|^2$$

In this case  $\tilde{G}(s) = 1/\|s\|^{2m}$  and the corresponding basis function is therefore

$$G(x) = \begin{cases} \|x\|^{2m-d} \ln \|x\| & \text{if } 2m > d \text{ and } d \text{ is even} \\ \|x\|^{2m-d} & \text{otherwise} \end{cases} \quad (3.2)$$

In this case the null space of  $\phi[f]$  is the vector space of polynomials of degree at most  $m$  in  $d$  variables, whose dimension is

$$k = \binom{d+m-1}{d}$$

These basis functions are radial and conditionally positive definite, so that they represent just particular instances of the well known radial basis functions technique (Micchelli 1986; Wahba 1990). In two dimensions, for  $m = 2$ , equation 3.2 yields the so-called "thin plate" basis function  $G(x) = \|x\|^2 \ln \|x\|$  (Harder and Desmarais 1972; Grimson 1982).

**3.1.2 The Gaussian.** A stabilizer of the form

$$\phi[f] = \int_{\mathbb{R}^d} ds e^{\frac{\|s\|^2}{\beta}} |\tilde{f}(s)|^2$$

where  $\beta$  is a fixed positive parameter, has  $\tilde{G}(s) = e^{-\|s\|^2/\beta}$  and as basis function the gaussian function (Poggio and Girosi 1989; Yuille and Grzywacz 1988). The gaussian function is positive definite, and it is well known from the theory of reproducing kernels (Aronszajn 1950) that positive definite functions (Stewart 1976) can be used to define *norms* of the type 3.1. Since  $\phi[f]$  is a norm, its null space contains only the zero element, and the additional null space terms of equation 2.3 are not needed, unlike in Duchon splines. A disadvantage of the gaussian is the appearance of the scaling parameter  $\beta$ , while Duchon splines, being homogeneous functions, do not depend on any scaling parameter. However, it is possible to devise good heuristics that furnish suboptimal, but still good, values of  $\beta$ , or good starting points for cross-validation procedures.

**3.1.3 Other Basis Functions.** Here we give a list of other functions that can be used as basis functions in the radial basis functions technique, and that are therefore associated with the minimization of some functional. In the following, we indicate as "p.d." the positive definite functions, which do not need any polynomial term in the solution, and as "c.p.d.  $k$ " the conditionally positive definite functions of order  $k$ , which need a polynomial of degree  $k$  in the solution. It is a well known fact that positive definite functions tend to zero at infinity whereas conditionally positive functions tend to infinity.

$G(r) = e^{-\beta r^2}$	Gaussian, p.d.
$G(r) = \sqrt{r^2 + c^2}$	multiquadric, c.p.d. 1
$G(r) = \frac{1}{\sqrt{c^2 + r^2}}$	inverse multiquadric, p.d.
$G(r) = r^{2n+1}$	thin plate splines, c.p.d. $n$
$G(r) = r^{2n} \ln r$	thin plate splines, c.p.d. $n$

**3.2 Tensor Product Stabilizers.** An alternative to choosing a radial function  $\tilde{G}(s)$  in the stabilizer 3.1 is a *tensor product* type of basis function, that is a function of the form

$$\tilde{G}(s) = \prod_{j=1}^d \tilde{g}(s_j) \quad (3.3)$$

where  $s_j$  is the  $j$ th coordinate of the vector  $s$ , and  $\tilde{g}$  is an appropriate one-dimensional function. When  $g$  is positive definite the functional  $\phi\{f\}$  is clearly a norm and its null space is empty. In the case of a conditionally positive definite function the structure of the null space can be more complicated and we do not consider it here. Stabilizers with  $\tilde{G}(s)$  as in equation 3.3 have the form

$$\phi\{f\} = \int_{R^d} ds \frac{|\tilde{f}(s)|^2}{\prod_{j=1}^d \tilde{g}(s_j)}$$

which leads to a *tensor product* basis function

$$G(x) = \prod_{j=1}^d g(x_j)$$

where  $x_j$  is the  $j$ th coordinate of the vector  $x$  and  $g(x)$  is the Fourier transform of  $\tilde{g}(s)$ . An interesting example is the one corresponding to the choice

$$\tilde{g}(s) = \frac{1}{1 + s^2}$$

which leads to the basis function

$$G(\mathbf{x}) = \prod_{j=1}^d e^{-|x_j|} = e^{-\sum_{j=1}^d |x_j|} = e^{-\|\mathbf{x}\|_{L_1}}$$

This basis function is interesting from the point of view of VLSI implementations, because it requires the computation of the  $L_1$  norm of the input vector  $\mathbf{x}$ , which is usually easier to compute than the Euclidean norm  $L_2$ . However, this basis function is not very smooth, and its performance in practical cases should first be tested experimentally. Notice that if the approximation is needed for computing derivatives smoothness of an appropriate degree is clearly a necessary requirement (see Poggio *et al.* 1988). We notice that the choice

$$\tilde{g}(s) = e^{-s^2}$$

leads again to the gaussian basis function  $G(\mathbf{x}) = e^{-\|\mathbf{x}\|^2}$ .

**3.3 Additive Stabilizers.** We have seen in the previous section how some tensor product approximation schemes can be derived in the framework of regularization theory. We now will see that it is also possible to derive the class of *additive approximation* schemes in the same framework, where by additive approximation we mean an approximation of the form

$$f(\mathbf{x}) = \sum_{\mu=1}^d f_{\mu}(x^{\mu}) \quad (3.4)$$

where  $x^{\mu}$  is the  $\mu$ th component of the input vector  $\mathbf{x}$  and the  $f_{\mu}$  are one-dimensional functions that will be defined as the *additive components* of  $f$  (from now on Greek letter indices will be used in association with components of the input vectors). Additive models are well known in statistics (Hastie and Tibshirani 1986, 1987, 1990; Stone 1985; Wahba 1990; Buja *et al.* 1989) and can be considered as a generalization of linear models. They are appealing because, being essentially a superposition of one-dimensional functions, they have a low complexity, and they share with linear models the feature that the effects of the different variables can be examined separately.

The simplest way to obtain such an approximation scheme is to choose, if possible, a stabilizer that corresponds to an additive basis function:

$$G(\mathbf{x}) = \sum_{\mu=1}^d \theta_{\mu} g(x^{\mu}) \quad (3.5)$$

where  $\theta_{\mu}$  are certain fixed parameters and  $g$  is a one-dimensional basis function. Such a choice would lead to an approximation scheme of the form 3.4 in which the additive components  $f_{\mu}$  have the form

$$f_{\mu}(x^{\mu}) = \theta_{\mu} \sum_{i=1}^N c_i g(x^{\mu} - x_i^{\mu}) \quad (3.6)$$

Notice that the additive components are not independent at this stage, since there is only one set of coefficients  $c_i$ . We postpone the discussion of this point to Section 4.2.

We would like then to write stabilizers corresponding to the basis function 3.5 in the form 3.1, where  $\tilde{G}(s)$  is the Fourier transform of  $G(x)$ . We notice that the Fourier transform of an additive function like the one in equation 3.5 exists only in the generalized sense (Gelfand and Shilov 1964), involving the  $\delta$  distribution. For example, in two dimensions we obtain

$$\tilde{G}(s) = \theta_x \tilde{g}(s_x) \delta(s_y) + \theta_y \tilde{g}(s_y) \delta(s_x) \quad (3.7)$$

and the interpretation of the reciprocal of this expression is delicate. However, *almost* additive basis functions can be obtained if we approximate the delta functions in equation 3.7 with gaussians of very small variance. Consider, for example in two dimensions, the stabilizer

$$\phi[f] = \int_{R^d} ds \, \epsilon \frac{|\tilde{f}(s)|^2}{\theta_x \tilde{g}(s_x) e^{-(s_y/\epsilon)^2} + \theta_y \tilde{g}(s_y) e^{-(s_x/\epsilon)^2}} \quad (3.8)$$

This corresponds to a basis function of the form

$$G(x, y) = \theta_x g(x) e^{-\epsilon^2 y^2} + \theta_y g(y) e^{-\epsilon^2 x^2} \quad (3.9)$$

In the limit of  $\epsilon$  going to zero the denominator in expression 3.8 approaches equation 3.7, and the basis function 3.9 approaches a basis function that is the sum of one-dimensional basis functions. In this paper we do not discuss this limit process in a rigorous way. Instead we outline another way to obtain additive approximations in the framework of regularization theory.

Let us assume that we know a priori that the function  $f$  that we want to approximate is additive, that is

$$f(x) = \sum_{\mu=1}^d f_{\mu}(x^{\mu})$$

We then apply the regularization approach and impose a smoothness constraint, not on the function  $f$  as a whole, but on each single additive component, through a regularization functional of the form (Wahba 1990; Hastie and Tibshirani 1990):

$$H[f] = \sum_{i=1}^N \left[ y_i - \sum_{\mu=1}^d f_{\mu}(x_i^{\mu}) \right]^2 + \lambda \sum_{\mu=1}^d \frac{1}{\theta_{\mu}} \int_R ds \, \frac{|\tilde{f}_{\mu}(s)|^2}{\tilde{g}(s)}$$

where  $\theta_{\mu}$  are given positive parameters that allow us to impose different degrees of smoothness on the different additive components. The min-

imizer of this functional is found with the same technique described in Appendix A, and, skipping null space terms, it has the usual form

$$f(\mathbf{x}) = \sum_{i=1}^N c_i G(\mathbf{x} - \mathbf{x}_i) \quad (3.10)$$

where

$$G(\mathbf{x} - \mathbf{x}_i) = \sum_{\mu=1}^d \theta_{\mu} g(x^{\mu} - x_i^{\mu})$$

as in equation 3.5.

We notice that the additive component of equation 3.10 can be written as

$$f_{\mu}(x^{\mu}) = \sum_{i=1}^N c_i^{\mu} g(x^{\mu} - x_i^{\mu})$$

where we have defined

$$c_i^{\mu} = c_i \theta_{\mu}$$

The additive components are therefore not independent because the parameters  $\theta_{\mu}$  are fixed. If the  $\theta_{\mu}$  were free parameters, the coefficients  $c_i^{\mu}$  would be independent, as well as the additive components.

Notice that the two ways we have outlined for deriving additive approximation from regularization theory are equivalent. They both start from a priori assumptions of additivity and smoothness of the class of functions to be approximated. In the first technique the two assumptions are woven together in the choice of the stabilizer (equation 3.8); in the second they are made explicit and exploited sequentially.

#### 4 Extensions: From Regularization Networks to Generalized Regularization Networks

---

In this section we will first review some extensions of regularization networks, and then will apply them to radial basis functions and to additive splines.

A fundamental problem in almost all practical applications in learning and pattern recognition is the choice of the relevant input variables. It may happen that some of the variables are more relevant than others, that some variables are just totally irrelevant, or that the relevant variables are linear combinations of the original ones. It can therefore be useful to work not with the original set of variables  $\mathbf{x}$ , but with a linear transformation of them,  $\mathbf{W}\mathbf{x}$ , where  $\mathbf{W}$  is a possibly rectangular matrix. In the framework of regularization theory, this can be taken into account by making the assumption that the approximating function  $f$  has the form  $f(\mathbf{x}) = F(\mathbf{W}\mathbf{x})$  for some smooth function  $F$ . The smoothness assumption is now made

directly on  $F$ , through a smoothness functional  $\phi[F]$  of the form 3.1. The regularization functional is expressed in terms of  $F$  as

$$H[F] = \sum_{i=1}^N [y_i - F(\mathbf{z}_i)]^2 + \lambda \phi[F]$$

where  $\mathbf{z}_i = \mathbf{W}\mathbf{x}_i$ . The function that minimizes this functional is clearly, accordingly to the results of Section 2, of the form

$$F(\mathbf{z}) = \sum_{i=1}^N c_i G(\mathbf{z} - \mathbf{z}_i)$$

(plus eventually a polynomial in  $\mathbf{z}$ ). Therefore the solution for  $f$  is

$$f(\mathbf{x}) = F(\mathbf{W}\mathbf{x}) = \sum_{i=1}^N c_i G(\mathbf{W}\mathbf{x} - \mathbf{W}\mathbf{x}_i) \quad (4.1)$$

This argument is rigorous for given and known  $\mathbf{W}$ , as in the case of classical radial basis functions. Usually the matrix  $\mathbf{W}$  is unknown, and it must be estimated from the examples. Estimating both the coefficients  $c_i$  and the matrix  $\mathbf{W}$  by least squares is usually not a good idea, since we would end up trying to estimate a number of parameters that is larger than the number of data points (though one may use regularized least squares). Therefore, it has been proposed (Moody and Darken 1988, 1989; Broomhead and Lowe 1988; Poggio and Girosi 1989, 1990a) that the approximation scheme of equation 4.1 be replaced with a similar one, in which the basic shape of the approximation scheme is retained, but the number of basis functions is decreased. The resulting approximating function that we call the *Generalized Regularization Network* (GRN) is

$$f(\mathbf{x}) = \sum_{\alpha=1}^n c_{\alpha} G(\mathbf{W}\mathbf{x} - \mathbf{W}\mathbf{t}_{\alpha}) \quad (4.2)$$

where  $n < N$  and the *centers*  $\mathbf{t}_{\alpha}$  are chosen according to some heuristic, or are considered as free parameters (Moody and Darken 1988, 1989; Poggio and Girosi 1989, 1990a). The coefficients  $c_{\alpha}$ , the elements of the matrix  $\mathbf{W}$ , and eventually the centers  $\mathbf{t}_{\alpha}$ , are estimated according to a least squares criterion. The elements of the matrix  $\mathbf{W}$  could also be estimated through cross-validation (Allen 1974; Wahba and Wold 1975; Golub *et al.* 1979; Craven and Wahba 1979; Utreras 1979; Wahba 1985), which may be a formally more appropriate technique.

In the special case in which the matrix  $\mathbf{W}$  and the centers are kept fixed, the resulting technique is one originally proposed by Broomhead and Lowe (1988), and the coefficients satisfy the following linear equation

$$\mathbf{G}^T \mathbf{G} \mathbf{c} = \mathbf{G}^T \mathbf{y}$$

where we have defined the following vectors and matrices:

$$(y)_i = y_i, \quad (c)_\alpha = c_\alpha, \quad (G)_{ii} = G(\mathbf{W}\mathbf{x}_i - \mathbf{W}\mathbf{t}_\alpha)$$

This technique, which has become quite common in the neural network community, has the advantage of retaining the form of the regularization solution, while being less complex to compute. A complete theoretical analysis has not yet been given, but some results, in the case in which the matrix  $\mathbf{W}$  is set to identity, are already available (Sivakumar and Ward 1991; Poggio and Girosi 1989).

The next sections discuss approximation schemes of the form 4.2 in the cases of radial and additive basis functions.

**4.1 Extensions of Radial Basis Functions.** In the case in which the basis function is radial, the approximation scheme of equation 4.2 becomes

$$f(\mathbf{x}) = \sum_{\alpha=1}^n c_\alpha G(\|\mathbf{x} - \mathbf{t}_\alpha\|_w)$$

where we have defined the weighted norm

$$\|\mathbf{x}\|_w \equiv \mathbf{x}\mathbf{W}^T\mathbf{W}\mathbf{x} \quad (4.3)$$

The basis functions of equation 4.2 are not radial any more, or, more precisely, they are radial in the metric defined by equation 4.3. This means that the level curves of the basis functions are not circles, but ellipses, whose axis does not need to be aligned with the coordinate axis. Notice that in this case what is important is not the matrix  $\mathbf{W}$  itself, but rather the symmetric matrix  $\mathbf{W}^T\mathbf{W}$ . Therefore, by the Cholesky decomposition, it is sufficient to consider  $\mathbf{W}$  to be upper triangular. The optimal center locations  $\mathbf{t}_\alpha$  satisfy the following set of nonlinear equations (Poggio and Girosi 1990a,b):

$$\mathbf{t}_\alpha = \frac{\sum_i P_i^\alpha \mathbf{x}_i}{\sum_i P_i^\alpha} \quad \alpha = 1, \dots, n \quad (4.4)$$

where  $P_i^\alpha$  are coefficients that depend on all the parameters of the network and are not necessarily positive. The optimal centers are then a weighted sum of the example points. Thus in some cases it may be more efficient to "move" the coefficients  $P_i^\alpha$  rather than the components of  $\mathbf{t}_\alpha$  (for instance when the dimensionality of the inputs is high relative to the number of data points).

The approximation scheme defined by equation 4.2 has been discussed in detail in Poggio and Girosi (1990a) and Girosi (1992), so we will not discuss it further. In the next section we will consider its analogue in the case of additive basis functions.

**4.2 Extensions of Additive Splines.** In the previous sections we have seen an extension of the classical regularization technique. In this section we derive the form that this extension takes when applied to additive splines. The resulting scheme is very similar to projection pursuit regression (Friedman and Stuetzle 1981; Huber 1985; Diaconis and Freedman 1984; Donoho and Johnstone 1989; Moody and Yarvin 1991).

We start from the “classical” additive spline, derived from regularization in Section 3.3:

$$f(\mathbf{x}) = \sum_{i=1}^N c_i \sum_{\mu=1}^d \theta_{\mu} g(x^{\mu} - x_i^{\mu}) \quad (4.5)$$

In this scheme the smoothing parameters  $\theta_{\mu}$  should be known, or can be estimated by cross-validation. An alternative to cross-validation is to consider the parameters  $\theta_{\mu}$  as *free parameters*, and estimate them with a least squares technique together with the coefficients  $c_i$ . If the parameters  $\theta_{\mu}$  are free, the approximation scheme of equation 4.5 becomes the following:

$$f(\mathbf{x}) = \sum_{i=1}^N \sum_{\mu=1}^d c_i^{\mu} g(x^{\mu} - x_i^{\mu})$$

where the coefficients  $c_i^{\mu}$  are now independent. Of course, now we must estimate  $N \times d$  coefficients instead of just  $N$ , and we are likely to encounter an overfitting problem. We then adopt the same idea presented in Section 4, and consider an approximation scheme of the form

$$f(\mathbf{x}) = \sum_{\alpha=1}^n \sum_{\mu=1}^d c_{\alpha}^{\mu} g(x^{\mu} - t_{\alpha}^{\mu}) \quad (4.6)$$

in which the number of centers is smaller than the number of examples, reducing the number of coefficients that must be estimated. We notice that equation 4.6 can be written as

$$f(\mathbf{x}) = \sum_{\mu=1}^d f_{\mu}(x^{\mu})$$

where each additive component has the form

$$f_{\mu}(x^{\mu}) = \sum_{\alpha=1}^n c_{\alpha}^{\mu} g(x^{\mu} - t_{\alpha}^{\mu})$$

Therefore another advantage of this technique is that the *additive components are now independent*, each of them being one-dimensional radial basis functions.

We can now use the same argument from Section 4 to introduce a linear transformation of the inputs  $\mathbf{x} \rightarrow \mathbf{W}\mathbf{x}$ , where  $\mathbf{W}$  is a  $d' \times d$  matrix.

Calling  $\mathbf{w}_\mu$  the  $\mu$ th row of  $W$ , and performing the substitution  $\mathbf{x} \rightarrow W\mathbf{x}$  in equation 4.6, we obtain

$$f(\mathbf{x}) = \sum_{\alpha=1}^n \sum_{\mu=1}^{d'} c_\alpha^\mu g(\mathbf{w}_\mu \cdot \mathbf{x} - t_\alpha^\mu) \quad (4.7)$$

We now define the following one-dimensional function:

$$h_\mu(y) = \sum_{\alpha=1}^n c_\alpha^\mu g(y - t_\alpha^\mu)$$

and rewrite the approximation scheme of equation 4.7 as

$$f(\mathbf{x}) = \sum_{\mu=1}^{d'} h_\mu(\mathbf{w}_\mu \cdot \mathbf{x}) \quad (4.8)$$

Notice the similarity between equation 4.8 and the projection pursuit regression technique: in both schemes the unknown function is approximated by a linear superposition of one-dimensional variables, which are projections of the original variables on certain vectors that have been estimated. In projection pursuit regression the choice of the functions  $h_\mu(y)$  is left to the user. In our case the  $h_\mu$  are one-dimensional radial basis functions, for example, cubic splines, or gaussians. The choice depends, strictly speaking, on the specific prior, that is, on the specific smoothness assumptions made by the user. Interestingly, in many applications of projection pursuit regression the functions  $h_\mu$  have been indeed chosen to be cubic splines but other choices are flexible Fourier series, rational approximations, and orthogonal polynomials (see Moody and Yarvin 1991).

Let us briefly review the steps that bring us from the classical additive approximation scheme of equation 3.6 to a projection pursuit regression-like type of approximation:

1. The regularization parameters  $\theta_\mu$  of the classical approximation scheme 3.6 are considered as free parameters.
2. The number of centers is chosen to be smaller than the number of data points.
3. The true relevant variables are assumed to be some unknown linear combination of the original variables.

We notice that in the extreme case in which each additive component has just one center ( $n = 1$ ), the approximation scheme of equation 4.7 becomes

$$f(\mathbf{x}) = \sum_{\mu=1}^{d'} c^\mu g(\mathbf{w}_\mu \cdot \mathbf{x} - t^\mu) \quad (4.9)$$

When the basis function  $g$  is a gaussian we call—somewhat improperly—a network of this type a *gaussian multilayer perceptron (MLP) network*, because if  $g$  were a threshold function sigmoidal function this would be a multilayer perceptron with one layer of hidden units. The sigmoidal function, typically used instead of the threshold, cannot be derived directly from regularization theory because it is not symmetric, but we will see in Section 6 the relationship between a sigmoidal function and the absolute value function, which is a basis function that can be derived from regularization.

There are a number of computational issues related to how to find the parameters of an approximation scheme like the one of equation 4.7, but we do not discuss them here. We present instead, in Section 7, some experimental results, and will describe the algorithm used to obtain them.

## 5 The Bayesian Interpretation of Generalized Regularization Networks

It is well known that a variational principle such as equation 2.1 can be derived not only in the context of functional analysis (Tikhonov and Arsenin 1977), but also in a probabilistic framework (Kimeldorf and Wahba 1971; Wahba 1980, 1990; Poggio *et al.* 1985; Marroquin *et al.* 1987; Bertero *et al.* 1988). In this section we illustrate this connection informally, without addressing the related mathematical issues.

Suppose that the set  $g = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}\}_{i=1}^N$  of data has been obtained by random sampling a function  $f$ , defined on  $\mathbb{R}^d$ , in the presence of noise, that is

$$f(\mathbf{x}_i) = y_i + \epsilon_i, \quad i = 1, \dots, N \quad (5.1)$$

where  $\epsilon_i$  are random independent variables with a given distribution. We are interested in recovering the function  $f$ , or an estimate of it, from the set of data  $g$ . We take a probabilistic approach, and regard the function  $f$  as the realization of a random field with a known prior probability distribution. Let us define

- $\mathcal{P}[f | g]$  as the conditional probability of the function  $f$  given the examples  $g$ .
- $\mathcal{P}[g | f]$  as the conditional probability of  $g$  given  $f$ . If the function underlying the data is  $f$ , this is the probability that by random sampling the function  $f$  at the sites  $\{\mathbf{x}_i\}_{i=1}^N$  the set of measurement  $\{y_i\}_{i=1}^N$  is obtained. This is therefore a model of the noise.
- $\mathcal{P}[f]$ : is the a priori probability of the random field  $f$ . This embodies our a priori knowledge of the function, and can be used to impose constraints on the model, assigning significant probability only to those functions that satisfy those constraints.

Assuming that the probability distributions  $\mathcal{P}[g | f]$  and  $\mathcal{P}[f]$  are known, the posterior distribution  $\mathcal{P}[f | g]$  can now be computed by applying the Bayes rule:

$$\mathcal{P}[f | g] \propto \mathcal{P}[g | f] \mathcal{P}[f] \quad (5.2)$$

We now make the assumption that the noise variables in equation 5.1 are normally distributed, with variance  $\sigma$ . Therefore the probability  $\mathcal{P}[g | f]$  can be written as

$$\mathcal{P}[g | f] \propto e^{-(1/2\sigma^2) \sum_{i=1}^N [y_i - f(x_i)]^2}$$

where  $\sigma$  is the variance of the noise.

The model for the prior probability distribution  $\mathcal{P}[f]$  is chosen in analogy with the discrete case (when the function  $f$  is defined on a finite subset of a  $n$ -dimensional lattice) for which the problem can be formalized (see for instance Marroquin *et al.* 1987). The prior probability  $\mathcal{P}[f]$  is written as

$$\mathcal{P}[f] \propto e^{-\alpha \phi[f]} \quad (5.3)$$

where  $\phi[f]$  is a smoothness functional of the type described in Section 3 and  $\alpha$  a positive real number. This form of probability distribution gives high probability only to those functions for which the term  $\phi[f]$  is small, and embodies the a priori knowledge that one has about the system.

Following the Bayes rule (5.2) the a posteriori probability of  $f$  is written as

$$\mathcal{P}[f | g] \propto e^{-(1/2\sigma^2) \sum_{i=1}^N [y_i - f(x_i)]^2 + 2\alpha\sigma^2 \phi[f]} \quad (5.4)$$

One simple estimate of the function  $f$  from the probability distribution 5.4 is the so-called *maximum a posteriori* (MAP) estimate, that considers the function that maximizes the a posteriori probability  $\mathcal{P}[f | g]$ , and therefore minimizes the exponent in equation 5.4. The MAP estimate of  $f$  is therefore the minimizer of the following functional:

$$H[f] = \sum_{i=1}^N [y_i - f(x_i)]^2 + \lambda \phi[f]$$

where  $\lambda = 2\sigma^2\alpha$ . This functional is the same as that of equation 2.1, and from here it is clear that the parameter  $\lambda$ , that is usually called the "regularization parameter" determines the trade-off between the level of the noise and the strength of the a priori assumptions about the solution, therefore controlling the compromise between the degree of smoothness of the solution and its closeness to the data. Notice that functionals of the type 5.3 are common in statistical physics (Parisi 1988), where  $\phi[f]$  plays the role of an energy functional. It is interesting to notice that, in that case, the correlation function of the physical system described by  $\phi[f]$  is the basis function  $G(x)$ .

As we have pointed out (Poggio and Girosi 1989; Rivest, personal communication), prior probabilities can also be seen as a measure of complexity, assigning high complexity to the functions with small probability. It has been proposed by Rissanen (1978) to measure the complexity of a hypothesis in terms of the bit length needed to encode it. It turns out that the MAP estimate mentioned above is closely related to the minimum description length principle: the hypothesis  $f$ , which for given  $g$  can be described in the most compact way, is chosen as the "best" hypothesis. Similar ideas have been explored by others (see for instance Solomonoff 1978). They connect data compression and coding with Bayesian inference, regularization, function approximation, and learning.

## 6 Additive Splines, Hinge Functions, Sigmoidal Neural Nets

In the previous sections we have shown how to extend RN to schemes that we have called GRN, which include ridge approximation schemes of the PPR type, that is

$$f(\mathbf{x}) = \sum_{\mu=1}^{d'} h_{\mu}(\mathbf{w}_{\mu} \cdot \mathbf{x})$$

where

$$h_{\mu}(y) = \sum_{\alpha=1}^n c_{\alpha}^{\mu} g(y - t_{\alpha}^{\mu})$$

The form of the basis function  $g$  depends on the stabilizer, and a list of "admissible"  $G$  has been given in Section 3. These include the absolute value  $g(x) = |x|$ —corresponding to piecewise linear splines, and the function  $g(x) = |x|^3$ —corresponding to cubic splines (used in typical implementations of PPR), as well as gaussian functions. Though it may seem natural to think that sigmoidal multilayer perceptrons may be included in this framework, it is actually impossible to derive *directly* from regularization principles the sigmoidal activation functions typically used in multilayer perceptrons. In the following section we show, however, that there is a close relationship between basis functions of the hinge, the sigmoid and the gaussian type.

**6.1 From Additive Splines to Ramp and Hinge Functions.** We will consider here the one-dimensional case, since multidimensional additive approximations consist of one-dimensional terms. We consider the approximation with the lowest possible degree of smoothness: piecewise linear. The associated basis function  $g(x) = |x|$  is shown in Figure 2a, and the associated stabilizer is given by

$$\phi[f] = \int_{-\infty}^{\infty} ds \, s^2 |\tilde{f}(s)|^2$$

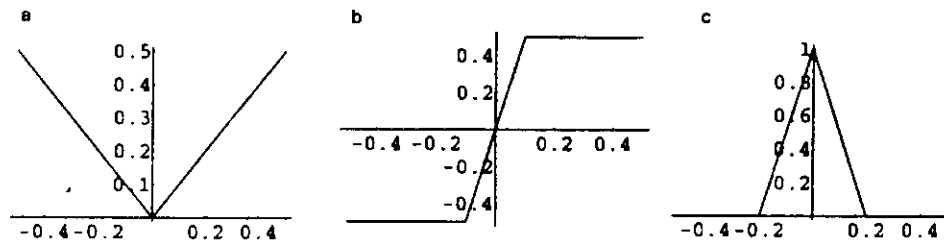


Figure 2: (a) Absolute value basis function,  $|x|$ . (b) Sigmoidal-like basis function  $\sigma_L(x)$ . (c) Gaussian-like basis function  $g_L(x)$ .

This assumption thus leads to approximating a one-dimensional function as the linear combination with appropriate coefficients of translates of  $|x|$ . It is easy to see that a linear combination of two translates of  $|x|$  with appropriate coefficients (positive and negative and equal in absolute value) yields the piecewise linear threshold function  $\sigma_L(x)$  also shown in Figure 2b. Linear combinations of translates of such functions can be used to approximate one-dimensional functions. A similar derivative-like, linear combination of two translates of  $\sigma_L(x)$  functions with appropriate coefficients yields the gaussian-like function  $g_L(x)$  also shown in Figure 2c. Linear combinations of translates of this function can also be used for approximation of a function. Thus any given approximation in terms of  $g_L(x)$  can be rewritten in terms of  $\sigma_L(x)$  and the latter can be in turn expressed in terms of the basis function  $|x|$ .

Notice that the basis functions  $|x|$  underlie the “hinge” technique proposed by Breiman (1993), whereas the basis functions  $\sigma_L(x)$  are sigmoidal-like and the  $g_L(x)$  are gaussian-like. The arguments above show the close relations between all of them, despite the fact that only  $|x|$  is strictly a “legal” basis function from the point of view of regularization [ $g_L(x)$  is not, though the very similar but smoother gaussian is]. Notice also that  $|x|$  can be expressed in terms of “ramp” functions, that is  $|x| = x_+ + x_-$ . Thus a one-hidden-layer perceptron using the activation function  $\sigma_L(x)$  can be rewritten in terms of a generalized regularization network with basis function  $|x|$ . The equivalent kernel is effectively local only if there exist a sufficient number of centers for each dimension ( $\mathbf{w}_\mu \cdot \mathbf{x}$ ). This is the case for projection pursuit regression but not for usual one-hidden-layer perceptrons.

These relationships imply that it may be interesting to compare how well each of these basis functions is able to approximate some simple function. To do this we used the model  $f(x) = \sum_a^n c_a g(w_a x - t_a)$  to approximate the function  $h(x) = \sin(2\pi x)$  on  $[0, 1]$ , where  $g(x)$  is one of the basis functions of Figure 2. Fifty training points and 10,000 test points

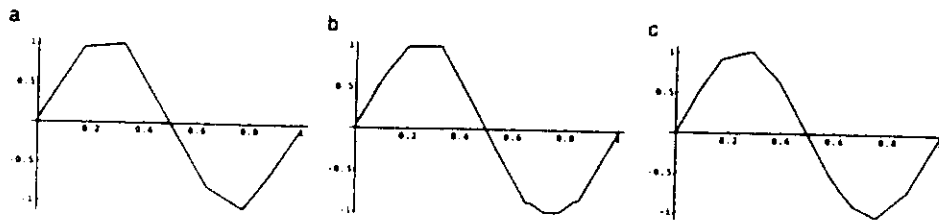


Figure 3: Approximation of  $\sin(2\pi x)$  using 8 basis functions of the (a) absolute value type, (b) sigmoidal-like type, and (c) gaussian-like type.

were chosen uniformly on  $[0, 1]$ . The parameters were learned using the iterative backfitting algorithm (Friedman and Stuetzle 1981; Hastie and Tibshirani 1990; Breiman 1993) that will be described in Section 7. We looked at the function learned after fitting 1, 2, 4, 8, and 16 basis functions. Some of the resulting approximations are plotted in Figure 3.

The results show that the performance of all three basis functions is fairly close as the number of basis functions increases. All models did a good job of approximating  $\sin(2\pi x)$ . The absolute value function did slightly worse and the "gaussian" function did slightly better. It is interesting that the approximation using two absolute value functions is almost identical to the approximation using one "sigmoidal" function, which again shows that two absolute value basis functions can sum to equal one "sigmoidal" piecewise linear function.

## 7 Numerical Illustrations

**7.1 Comparing Additive and Nonadditive Models.** To illustrate some of the ideas presented in this paper and to provide some practical intuition about the various models, we present numerical experiments comparing the performance of additive and nonadditive networks on two-dimensional problems. In a model consisting of a sum of two-dimensional gaussians, the model can be changed from a nonadditive radial basis function network to an additive network by "elongating" the gaussians along the two coordinate axes  $x$  and  $y$ . This allows us to measure the performance of a network as it changes from a nonadditive scheme to an additive one.

Five different models were tested. The first three differ only in the variances of the gaussian along the two coordinate axes. The ratio of the

$x$  variance to the  $y$  variance determines the elongation of the gaussian. These models all have the same form and can be written as

$$f(\mathbf{x}) = \sum_{i=1}^N c_i [G_1(\mathbf{x} - \mathbf{x}_i) + G_2(\mathbf{x} - \mathbf{x}_i)]$$

where

$$G_1 = e^{-[(x^2/\sigma_1) + (y^2/\sigma_2)]}, \quad G_2 = e^{-[(x^2/\sigma_2) + (y^2/\sigma_1)]}$$

The models differ only in the values of  $\sigma_1$  and  $\sigma_2$ . For the first model,  $\sigma_1 = 0.5$  and  $\sigma_2 = 0.5$  (RBF), for the second model  $\sigma_1 = 10$  and  $\sigma_2 = 0.5$  (elliptical gaussian), and for the third model,  $\sigma_1 = \infty$  and  $\sigma_2 = 0.5$  (additive). These models correspond to placing two gaussians at each data point  $\mathbf{x}_i$ , with one gaussian elongated in the  $x$  direction and one elongated in the  $y$  direction. In the first case (RBF) there is no elongation, in the second case (elliptical gaussian) there is moderate elongation, and in the last case (additive) there is infinite elongation.

The fourth model is a generalized regularization network model, of the form 4.9, that uses a gaussian basis function:

$$f(\mathbf{x}) = \sum_{\alpha=1}^n c_{\alpha} e^{-(\mathbf{w}_{\alpha} \cdot \mathbf{x} - t_{\alpha})^2}$$

In this model, to which we referred earlier as a gaussian MLP network (equation 4.9), the weight vectors, centers, and coefficients are all learned.

In order to see how sensitive were the performances to the choice of basis function, we also repeated the experiments for model 4 with a sigmoid (that is *not* a basis function that can be derived from regularization theory) replacing the gaussian basis function. In our experiments we used the standard sigmoid function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Models 1 to 5 are summarized in Table 1: notice that only model 5 is a multilayer perceptron in the standard sense.

In the first three models, the centers were fixed in the learning algorithm and equal to the training examples. The only parameters that were learned were the coefficients  $c_i$ , that were computed by solving the linear system of equations 2.4. The fourth and the fifth models were trained by fitting one basis function at a time according to the following recursive algorithm with backfitting (Friedman and Stuezle 1981; Hastie and Tibshirani 1990; Breiman 1993)

- Add a new basis function;
- Optimize the parameters  $\mathbf{w}_{\alpha}$ ,  $t_{\alpha}$ , and  $c_{\alpha}$  using the "random step" algorithm (Caprile and Girosi 1990) described below;

Table 1: The Five Models Tested in our Numerical Experiments.

Model 1		
	$f(x, y) = \sum_{i=1}^{20} c_i \left[ e^{-\left(\frac{(x-x_i)^2}{\sigma_1} + \frac{(y-y_i)^2}{\sigma_2}\right)} + e^{-\left(\frac{(x-x_i)^2}{\sigma_2} + \frac{(y-y_i)^2}{\sigma_1}\right)} \right]$	$\sigma_1 = \sigma_2 = 0.5$
Model 2		
	$f(x, y) = \sum_{i=1}^{20} c_i \left[ e^{-\left(\frac{(x-x_i)^2}{\sigma_1} + \frac{(y-y_i)^2}{\sigma_2}\right)} + e^{-\left(\frac{(x-x_i)^2}{\sigma_2} + \frac{(y-y_i)^2}{\sigma_1}\right)} \right]$	$\sigma_1 = 10, \sigma_2 = 0.5$
Model 3		
	$f(x, y) = \sum_{i=1}^{20} c_i \left[ e^{-\frac{(x-x_i)^2}{\sigma}} + e^{-\frac{(y-y_i)^2}{\sigma}} \right]$	$\sigma = 0.5$
Model 4		
	$f(x, y) = \sum_{\alpha=1}^n c_{\alpha} e^{-(\mathbf{w}_{\alpha} \cdot \mathbf{x} - t_{\alpha})^2}$	—
Model 5		
	$f(x, y) = \sum_{\alpha=1}^n c_{\alpha} \sigma(\mathbf{w}_{\alpha} \cdot \mathbf{x} - t_{\alpha})$	—

- Backfitting: for each basis function  $\alpha$  added so far:
  - hold the parameters of all other functions fixed;
  - reoptimize the parameters of function  $\alpha$ ;
- Repeat the backfitting stage until there is no significant decrease in  $L_2$  error.

The “random step” (Caprile and Girosi 1990) is a stochastic optimization algorithm that is very simple to implement and that usually finds good local minima. The algorithm works as follows: pick random changes to each parameter such that each random change lies within some interval  $[a, b]$ . Add the random changes to each parameter and then calculate the new error between the output of the network and the target values. If the error decreases, then keep the changes and double the length of the interval for picking random changes. If the error increases, then throw out the changes and halve the size of the interval. If the length of the interval becomes less than some threshold, then reset the length of the interval to some larger value.

The five models were each tested on two different functions: a two-dimensional additive function

$$h_{\text{add}}(x, y) = \sin(2\pi x) + 4(y - 0.5)^2$$

and the two-dimensional Gabor function

$$g_{\text{Gabor}}(x, y) = e^{-\|\mathbf{x}\|^2} \cos [0.75\pi(x + y)]$$

Table 2: A Summary of the Results of Our Numerical Experiments.<sup>a</sup>

	Model 1	Model 2	Model 3	Model 4	Model 5
$h_{\text{add}}(x, y)$					
Training	0.000036	0.000067	0.000001	0.000170	0.000743
Test	0.011717	0.001598	0.000007	0.001422	0.026699
$g_{\text{Gabor}}(x, y)$					
Training	0.000000	0.000000	0.000000	0.000001	0.000044
Test	0.003818	0.344881	67.95237	0.033964	0.191055

<sup>a</sup>Each table entry contains the  $L_2$  errors for both the training set and the test set.

The training data for the functions  $h_{\text{add}}$  and  $g_{\text{Gabor}}$  consisted of 20 points picked from a uniform distribution on  $[0, 1] \times [0, 1]$  and  $[-1, 1] \times [-1, 1]$ , respectively. Another 10,000 points were randomly chosen to serve as test data. The results are summarized in Table 2 (see Girosi *et al.* 1993 for a more extensive description of the results).

As expected, the results show that the additive model 3 was able to approximate the additive function,  $h_{\text{add}}(x, y)$  better than both the RBF model 1 and the elliptical gaussian model 2, and that there seems to be a smooth degradation of performance as the model changes from the additive to the radial basis function. Just the opposite results are seen in approximating the nonadditive Gabor function,  $g_{\text{Gabor}}(x, y)$ , shown in Figure 4a. The RBF model 1 did very well, while the additive model 3 did a very poor job, as shown in Figure 4b. However, Figure 4c shows that the GRN scheme (model 4) gives a fairly good approximation, because the learning algorithm finds better directions for projecting the data than the  $x$  and  $y$  axis as in the pure additive model.

Notice that the first three models we considered had a number of parameters equal to the number of data points, and were supposed to exactly interpolate the data, so that one may wonder why the training errors are not exactly zero. The reason is the ill-conditioning of the associated linear system, which is a typical problem of radial basis functions (Dyn *et al.* 1986).

## 8 Hardware and Biological Implementation of Network Architectures

We have seen that different network architectures can be derived from regularization by making somewhat different assumptions on the classes of functions used for approximation. Given the basic common roots, one is tempted to argue—and numerical experiments support the claim—that there will be small differences in average performance of the various architectures (see also Lippmann 1989; Lippmann and Lee 1991).

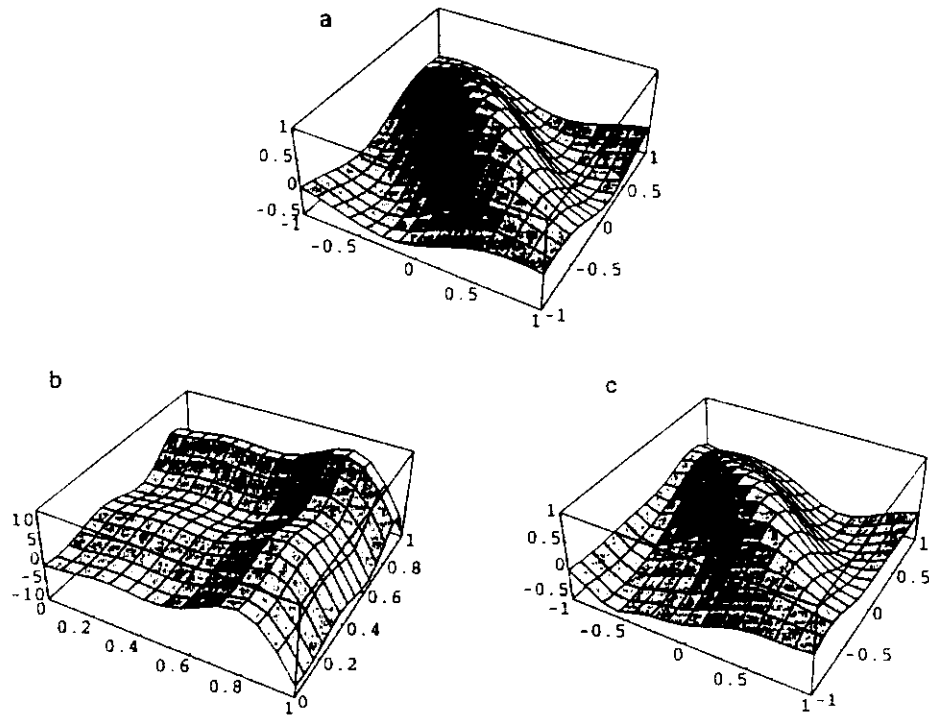


Figure 4: (a) The function to be approximated  $g(x,y)$ . (b) Additive gaussian model approximation of  $g(x,y)$  (model 3). (c) GRN approximation of  $g(x,y)$  (model 4).

It therefore becomes interesting to ask which architectures are easier to implement in hardware.

All the schemes that use the same number of centers as examples—such as RBF and additive splines—are expensive in terms of memory requirements (if there are many examples) but have a simple learning stage. More interesting are the schemes that use fewer centers than examples (and use the linear transformation  $W$ ). There are at least two perspectives for our discussion: we can consider implementation of radial vs. additive schemes and we can consider different activation functions. Let us first discuss radial vs. nonradial functions such as a gaussian RBF vs. a gaussian MLP network. For VLSI implementations, the main difference is in computing a scalar product rather than an  $L_2$  distance, which is usually more expensive both for digital and analog VLSI. The  $L_2$  distance, however, might be replaced with the  $L_1$  distance, that is a sum of absolute values, which can be computed efficiently. Notice that a radial basis functions scheme that uses the  $L_1$  norm has been derived in Section 3.2 from a tensor-product stabilizer.

Let us consider now different activation functions. Activation functions such as gaussian, sigmoid, or absolute values are equally easy to compute, especially if look-up table approaches are used. In analog hardware it is somewhat simpler to generate a sigmoid than a gaussian, although gaussian-like shapes can be synthesized with fewer than 10 transistors (J. Harris, personal communication).

In practical implementations other issues, such as trade-offs between memory and computation and on-chip learning, are likely to be much more relevant than the specific chosen architecture. In other words, a general conclusion about ease of implementation is not possible: none of the architectures we have considered holds a clear edge.

From the point of view of *biological implementations* the situation is somewhat different. The hidden unit in MLP networks with sigmoidal-like activation functions is a plausible, albeit much oversimplified, model of real neurons. The sigmoidal transformation of a scalar product seems much easier to implement in terms of known biophysical mechanisms than the gaussian of a multidimensional Euclidean distance. On the other hand, it is intriguing to observe that HBF centers and tuned cortical neurons behave alike (Poggio and Hurlbert 1994). In particular, a gaussian HBF unit is maximally excited when each component of the input exactly matches each component of the center. Thus the unit is optimally tuned to the stimulus value specified by its center. Units with multidimensional centers are tuned to complex features, made of the conjunction of simpler features. This description is very like the customary description of cortical cells optimally tuned to some more or less complex stimulus. So-called place coding is the simplest and most universal example of tuning: cells with roughly bell-shaped receptive fields have peak sensitivities for given locations in the input space, and by overlapping, cover all of that space. Thus tuned cortical neurons seem to behave more like gaussian HBF units than like the sigmoidal units of MLP networks: the tuned response function of cortical neurons mostly resembles  $\exp(-\|x - t\|^2)$  more than it does  $\sigma(x \cdot w)$ . When the stimulus to a cortical neuron is changed from its optimal value in any direction, the neuron's response typically decreases. The activity of a gaussian HBF unit would also decline with any change in the stimulus away from its optimal value  $t$ . For the sigmoid unit, though, certain changes away from the optimal stimulus will not decrease its activity, for example, when the input  $x$  is multiplied by a constant  $\alpha > 1$ .

How might, then, multidimensional gaussian receptive fields be synthesized from known receptive fields and biophysical mechanisms?

The simplest answer is that cells tuned to complex features may be constructed from a hierarchy of simpler cells tuned to incrementally larger conjunctions of elementary features. This idea—popular among physiologists—can immediately be formalized in terms of gaussian radial basis functions, since a multidimensional gaussian function can be decomposed into the product of lower dimensional gaussians (Ballard

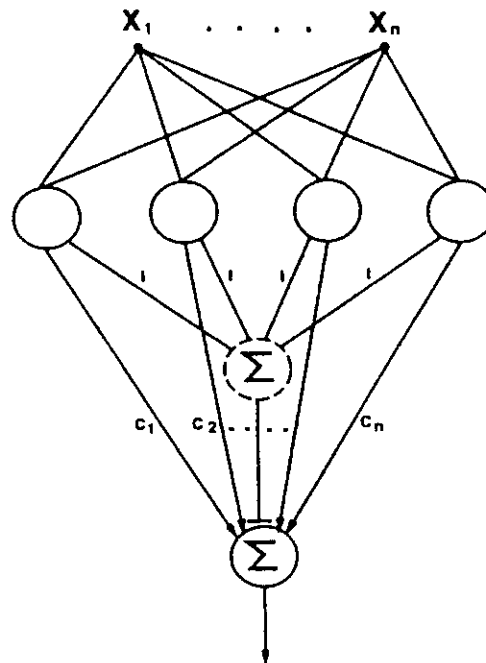


Figure 5: An implementation of the normalized radial basis function scheme. A "pool" cell (dotted circle) summates the activities of the hidden units and then divides the output of the network. The division may be approximated in a physiological implementation by shunting inhibition.

1986; Mel 1988, 1990, 1992; Poggio and Girosi 1990a). There are several biophysically plausible ways to implement gaussian RBF-like units (see Poggio and Girosi 1989; Poggio 1990), but none is particularly simple. Ironically one of the plausible implementations of a RBF unit may exploit circuits based on sigmoidal nonlinearities (see Poggio and Hurlbert 1994). In general, the circuits required for the various schemes described in this paper are reasonable from a biological point of view (Poggio and Girosi 1989; Poggio 1990). For example, the normalized basis function scheme of Section 2.2 could be implemented as outlined in Figure 5 where a "pool" cell summates the activities of all hidden units and shunts the output unit with a shunting inhibition approximating the required division operation.

## 9 Summary and Remarks

A large number of approximation techniques can be written as multilayer networks with one hidden layer. In past papers (Poggio and

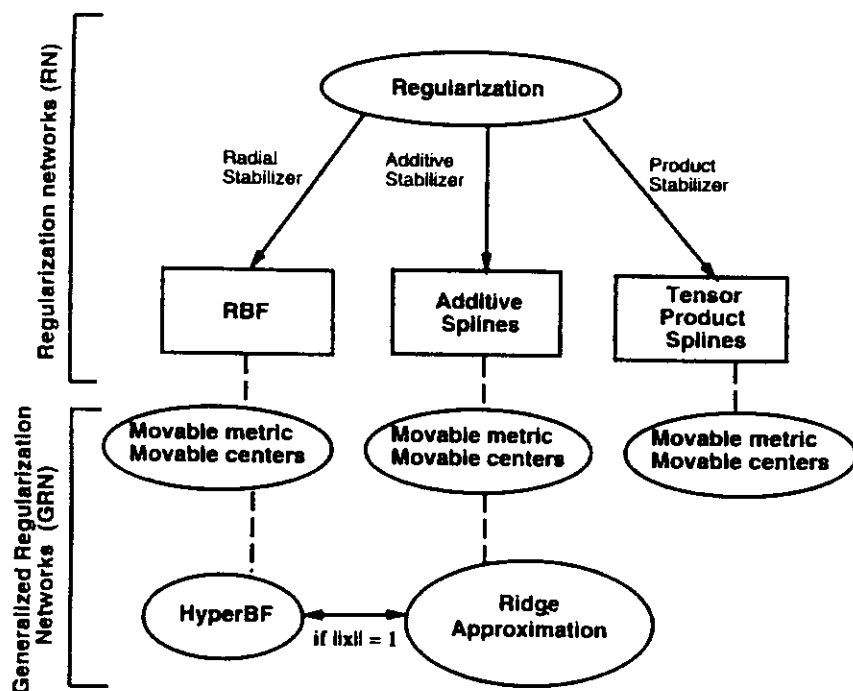


Figure 6: Several classes of approximation schemes and corresponding network architectures can be derived from regularization with the appropriate choice of smoothness priors and associated stabilizers and basis functions, showing the common Bayesian roots.

Girosi 1989, 1990; Girosi 1992) we showed how to derive radial basis functions, hyper basis functions, and several types of multidimensional splines from regularization principles. We had not used regularization to yield approximation schemes of the additive type (Wahba 1990; Hastie and Tibshirani 1990), such as additive splines, ridge approximation of the projection pursuit regression type, and hinge functions. In this paper, we show that appropriate stabilizers can be defined to justify such additive schemes, and that the same extensions that lead from RBF to HBF lead from additive splines to ridge function approximation schemes of the projection pursuit regression type. Our generalized regularization networks include, depending on the stabilizer (that is on the prior knowledge on the functions we want to approximate), HBF networks, ridge approximation, tensor products splines, and perceptron-like networks with one hidden layer and appropriate activation functions (such as the gaussian). Figure 6 shows a diagram of the relationships. Notice that HBF networks and ridge approximation networks are directly related in the special case of normalized inputs (Maruyama *et al.* 1992).

We now feel that a common theoretical framework justifies a large spectrum of approximation schemes in terms of different smoothness constraints imposed within the same regularization functional to solve the ill-posed problem of function approximation from sparse data. The claim is that many different networks and corresponding approximation schemes can be derived from the variational principle

$$H[f] = \sum_{i=1}^N [f(\mathbf{x}_i) - y_i]^2 + \lambda \phi[f] \quad (9.1)$$

They differ because of different choices of stabilizers  $\phi$ , which correspond to different assumptions of smoothness. In this context, we believe that the Bayesian interpretation is one of the main advantages of regularization: it makes clear that different network architectures correspond to different prior assumptions of smoothness of the functions to be approximated.

The common framework we have derived suggests that differences between the various network architectures are relatively minor, corresponding to different smoothness assumptions. One would expect that each architecture will work best for the class of function defined by the associated prior (that is stabilizer), an expectation that is consistent with numerical results in this paper (see also Donoho and Johnstone 1989).

**9.1 Classification and Smoothness.** From the point of view of regularization, the task of classification—instead of regression—may seem to represent a problem since the role of smoothness is less obvious. Consider for simplicity binary classification, in which the output  $y$  is either 0 or 1 and let  $P(\mathbf{x}, y) = P(\mathbf{x})P(y | \mathbf{x})$  be the joint probability of the input-output pairs  $(\mathbf{x}, y)$ . The average cost associated to an estimator  $f(\mathbf{x})$  is the *expected risk* (see Section 2.2)

$$I[f] = \int d\mathbf{x} dy P(\mathbf{x}, y) [y - f(\mathbf{x})]^2$$

The problem of learning is now equivalent to minimizing the expected risk based on  $N$  samples of the joint probability distribution  $P(\mathbf{x}, y)$ , and it is usually solved by minimizing the empirical risk (2.14). Here we discuss two possible approaches to the problem of finding the best estimator:

- If we look for an estimator in the class of real valued functions, it is well known that the minimizer  $f_0$  of  $Q[f]$  is the so-called *regression function*, that is

$$f_0(x) = \int dy y P(y | \mathbf{x}) = P(1 | \mathbf{x}) \quad (9.2)$$

Therefore, a real valued network  $f$  trained on the empirical risk (2.14) will approximate, under certain conditions of consistency

(Vapnik 1982; Vapnik and Chervonenkis 1991), the conditional probability distribution of class 1,  $P(1 | \mathbf{x})$ . In this case our final estimator  $f$  is real valued, and in order to obtain a binary estimator we have to apply a threshold function to it, so that our final solution turns out to be

$$f^*(\mathbf{x}) = \theta[f(\mathbf{x})]$$

where  $\theta$  is the Heaviside function.

- We could look for an estimator with range  $\{0, 1\}$ , for example of the form  $f(\mathbf{x}) = \theta[g(\mathbf{x})]$ . In this case the expected risk becomes the average number of misclassified vectors. The function that minimizes the expected risk is not the regression function any more, but a binary approximation to it.

We argue that in both cases it makes sense to assume that  $f$  (and  $g$ ) is a smooth real-valued function and therefore to use regularization networks to approximate it. The argument is that a natural prior constraint for classification is smoothness of classification boundaries, since otherwise it would be impossible to effectively generalize the correct classification from a set of examples. Furthermore, a condition that usually provides smooth classification boundaries is smoothness of the underlying regressor: a smooth function usually has "smooth" level crossings. Thus both approaches described above suggest to impose smoothness of  $f$  or  $g$ , that is to approximate  $f$  or  $g$  with a regularization network.

**9.2 Complexity of the Approximation Problem.** So far we have discussed several approximation techniques only from the point of view of the representation and architecture, and we did not discuss how well they perform in approximating functions of different functions spaces. Since these techniques are derived under different a priori smoothness assumptions, we clearly expect them to perform optimally when those a priori assumptions are satisfied. This makes it difficult to compare their performances, since we expect each technique to work best on a different class of functions. However, if we measure performances by how quickly the approximation error goes to zero when the number of parameters of the approximation scheme goes to infinity, very general results from the theory of linear and nonlinear widths (Timan 1963; Pinkus 1986; Lorentz 1962, 1986; DeVore *et al.* 1989; DeVore 1991; DeVore and Yu 1991) suggest that all techniques share the same limitations. For example, when approximating an  $s$  times continuously differentiable function in  $d$  variables with some function parameterized by  $n$  parameters, one can prove that even the "best" nonlinear parameterization cannot achieve an accuracy that is better than the *Jackson type* bound, that is  $O(n^{-s/d})$ . Here the adjective "best" is used in the sense defined by DeVore *et al.* (1989)

in their work on nonlinear  $n$ -widths, which restricts the sets of nonlinear parameterization to those for which the optimal parameters depend continuously on the function that has to be approximated. Notice that, although this is a desirable property, not all the approximation techniques may have it, and therefore these results may not always be applicable. However, the basic intuition is that a class of functions has an intrinsic complexity that increases exponentially in the ratio  $d/s$ , where  $s$  is a smoothness index, that is a measure of the amount of constraints imposed on the functions of the class. Therefore, if the smoothness index is kept constant, we expect that the number of parameters needed in order to achieve a certain accuracy increases exponentially with the number of dimensions, irrespectively of the approximation technique, showing the phenomenon known as "the curse of dimensionality" (Bellman 1961). Clearly, if we consider classes of functions with a smoothness index that increases when the number of variables increases, then a rate of convergence independent of the dimensionality can be obtained, because the increase in complexity due to the larger number of variables is compensated by the decrease due to the stronger smoothness constraint. To make this concept clear, we summarized in Table 3 a number of different approximation techniques, and the constraints that can be imposed on them in order to make the approximation error to be  $O(1/\sqrt{n})$ , that is "independent of the dimension," and therefore immune to the curse of dimensionality. Notice that since these techniques are derived under different a priori assumptions, the explicit form of the constraints are different. For example in entries 5 and 6 of Table 3 (Girosi and Anzellotti 1992, 1993; Girosi 1993) the result holds in  $H^{2m,1}(R^d)$ , that is the Sobolev space of functions whose derivatives up to order  $2m$  are integrable (Ziemer 1989). Notice that the number of derivatives that are integrable has to increase with the dimension  $d$  in order to keep the rate of convergence constant. A similar phenomenon appears in entries 2 and 3 (Barron 1991, 1993; Breiman 1993), but in a less obvious way. In fact, it can be shown (Girosi and Anzellotti 1992, 1993) that, for example, the spaces of functions considered by Barron (entry 2) and Breiman (entry 3) are the set of functions that can be written respectively as  $f(\mathbf{x}) = \|\mathbf{x}\|^{1-d} * \lambda$  and  $f(\mathbf{x}) = \|\mathbf{x}\|^{2-d} * \lambda$ , where  $\lambda$  is any function whose Fourier transform is integrable, and  $*$  stands for the convolution operator. Notice that, in this way, it becomes more apparent that these space of functions become more and more constrained as the dimensions increase, due to the more and more rapid fall-off of the terms  $\|\mathbf{x}\|^{1-d}$  and  $\|\mathbf{x}\|^{2-d}$ . The same phenomenon is also very clear in the results of Mhaskar (1993a), who proved that the rate of convergence of approximation of functions with  $s$  continuous derivatives by multilayered feedforward neural networks is  $O(n^{-s/d})$ : if the number of continuous derivatives  $s$  increases linearly with the dimension  $d$ , the curse of dimensionality disappears, leading to a rate of convergence independent of the dimension.

It is important to emphasize that in practice the parameters of the

Table 3: Approximation Schemes and Corresponding Functions Spaces with the Same Rate of Convergence  $O(n^{1/2})$ .<sup>a</sup>

Function space	Norm	Approximation scheme
$\int_{R^d} ds  \tilde{f}(s)  < +\infty$ (Jones 1992)	$L_2(\Omega)$	$f(x) = \sum_{i=1}^n c_i \sin(x \cdot w_i + \theta_i)$
$\int_{R^d} ds \ s\   \tilde{f}(s)  < +\infty$ (Barron 1991)	$L_2(\Omega)$	$f(x) = \sum_{i=1}^n c_i \sigma(x \cdot w_i + \theta_i)$
$\int_{R^d} ds \ s\ ^2  \tilde{f}(s)  < +\infty$ (Breiman 1993)	$L_2(\Omega)$	$f(x) = \sum_{i=1}^n c_i  x \cdot w_i + \theta_i _+ + x \cdot a + b$
$e^{-\ x\ ^2} * \lambda, \lambda \in L_1(R^d)$ (Girosi and Anzellotti 1992)	$L_\infty(R^2)$	$f(x) = \sum_{\alpha=1}^n c_\alpha e^{-\ x-t_\alpha\ ^2}$
$H^{2m,1}(R^d), 2m > d$ (Girosi and Anzellotti 1992)	$L_\infty(R^2)$	$f(x) = \sum_{\alpha=1}^n c_\alpha G_m(\ x-t_\alpha\ ^2)$
$H^{2m,1}(R^d), 2m > d$ (Girosi 1993)	$L_2(R^2)$	$f(x) = \sum_{\alpha=1}^n c_\alpha e^{-\ x-t_\alpha\ ^2/\sigma_\alpha^2}$

<sup>a</sup>The function  $\sigma$  is the standard sigmoidal function, the function  $|x|_+$  in the third entry is the ramp function, and the function  $G_m$  in the fifth entry is a Bessel potential, that is the Fourier transform of  $(1 + \|s\|^2)^{-m/2}$  (Stein 1970).  $H^{2m,1}(R^d)$  is the Sobolev space of functions whose derivatives up to order  $2m$  are integrable (Ziemer 1989).

approximation scheme have to be estimated using a finite amount of data (Vapnik and Chervonenkis 1971, 1981, 1991; Vapnik 1982; Pollard 1984; Geman *et al.* 1992; Haussler 1989; Baum and Haussler 1989; Baum 1988; Moody 1991a,b). In fact, what one does in practice is to minimize the *empirical risk* (see equation 2.14), while what one would really like to do is to minimize the *expected risk* (see equation 2.13). This introduces an additional source of error, sometimes called "estimation error," that usually depends on the dimension  $d$  in a much milder way than the approximation error, and can be estimated using the theory of uniform convergence of relative frequencies to probabilities (Vapnik and Chervonenkis 1971, 1981, 1991; Vapnik 1982; Pollard 1984).

Specific results on the generalization error, that combine both approximation and estimation error, have been obtained by Barron (1991, 1994) for sigmoidal neural networks, and by Niyogi and Girosi (1994) for gaussian radial basis functions. Although these bounds are different, they all have the same qualitative behavior: for a fixed number of data points the generalization error first decreases when the number of parameters increases, then reaches a minimum and starts increasing again, revealing the well known phenomenon of *overfitting*. For a general description of how the approximation and estimation error combine together to bound the generalization error see Niyogi and Girosi (1994).

**9.3 Additive Structure and the Sensory World.** In this last section we address the surprising relative success of additive schemes of the ridge approximation type in real world applications. As we have seen, ridge approximation schemes depend on priors that combine additivity of one-dimensional functions with the usual assumption of smoothness. Do such priors capture some fundamental property of the physical world? Consider, for example, the problem of object recognition, or the problem of motor control. We can recognize almost any object from any of many small subsets of its features, visual and nonvisual. We can perform many motor actions in several different ways. In most situations, our sensory and motor worlds are *redundant*. In terms of GRN this means that instead of high-dimensional centers, any of several lower-dimensional centers, that is components, are often sufficient to perform a given task. This means that the “and” of a high-dimensional conjunction can be replaced by the “or” of its components (low-dimensional conjunctions)—a face may be recognized by its eyebrows alone, or a mug by its color. To recognize an object, we may use not only templates comprising all its features, but also subtemplates, comprising subsets of features and in some situations the latter, by themselves, may be fully sufficient. Additive, small centers—in the limit with dimensionality one—with the appropriate  $W$  are of course associated with stabilizers of the additive type.

Splitting the recognizable world into its additive parts may well be preferable to reconstructing it in its full multidimensionality, because a system composed of several independent, additive parts is inherently more robust than a whole simultaneously dependent on each of its parts. The small loss in uniqueness of recognition is easily offset by the gain against noise and occlusion. There is also a possible meta-argument that we mention here only for the sake of curiosity. It may be argued that humans would not be able to understand the world if it were not additive because of the too-large number of necessary examples (because of high dimensionality of any sensory input such as an image). Thus one may be tempted to conjecture that our sensory world is biased towards an “additive structure.”

#### Appendix A: Derivation of the General Form of Solution of the Regularization Problem

---

We have seen in Section 2 that the regularized solution of the approximation problem is the function that minimizes a cost functional of the following form:

$$H[f] = \sum_{i=1}^N [y_i - f(x_i)]^2 + \lambda \phi[f] \quad (\text{A.1})$$

where the smoothness functional  $\phi[f]$  is given by

$$\phi[f] = \int_{R^d} ds \frac{|\tilde{f}(s)|^2}{\tilde{G}(s)}$$

The first term measures the distance between the data and the desired solution  $f$ , and the second term measures the cost associated with the deviation from smoothness. For a wide class of functionals  $\phi$  the solutions of the minimization problem A.1 all have the same form. A detailed and rigorous derivation of the solution of the variational principle associated with equation A.1 is outside the scope of this paper. We present here a simple derivation and refer the reader to the current literature for the mathematical details (Wahba 1990; Madych and Nelson 1990; Dyn 1987).

We first notice that, depending on the choice of  $G$ , the functional  $\phi[f]$  can have a nonempty null space, and therefore there is a certain class of functions that are "invisible" to it. To cope with this problem we first define an equivalence relation among all the functions that differ for an element of the null space of  $\phi[f]$ . Then we express the first term of  $H[f]$  in terms of the Fourier transform of  $f$ :<sup>2</sup>

$$f(\mathbf{x}) = \int_{R^d} ds \tilde{f}(s) e^{i\mathbf{x} \cdot \mathbf{s}}$$

obtaining the functional

$$H[\tilde{f}] = \sum_{i=1}^N \left[ y_i - \int_{R^d} ds \tilde{f}(s) e^{i\mathbf{x}_i \cdot \mathbf{s}} \right]^2 + \lambda \int_{R^d} ds \frac{|\tilde{f}(s)|^2}{\tilde{G}(s)}$$

Then we notice that since  $f$  is real, its Fourier transform satisfies the constraint:

$$\tilde{f}^*(\mathbf{s}) = \tilde{f}(-\mathbf{s})$$

so that the functional can be rewritten as

$$H[\tilde{f}] = \sum_{i=1}^N \left[ y_i - \int_{R^d} ds \tilde{f}(s) e^{i\mathbf{x}_i \cdot \mathbf{s}} \right]^2 + \lambda \int_{R^d} ds \frac{\tilde{f}(-\mathbf{s}) \tilde{f}(\mathbf{s})}{\tilde{G}(\mathbf{s})}$$

In order to find the minimum of this functional we take its functional derivatives with respect to  $\tilde{f}$  and set it to zero:

$$\frac{\delta H[\tilde{f}]}{\delta \tilde{f}(\mathbf{t})} = 0 \quad \forall \mathbf{t} \in R^d \quad (\text{A.2})$$

<sup>2</sup>For simplicity of notation we take all the constants that appear in the definition of the Fourier transform to be equal to 1.

We now proceed to compute the functional derivatives of the first and second term of  $H[f]$ . For the first term we have

$$\begin{aligned} \frac{\delta}{\delta \tilde{f}(\mathbf{t})} \sum_{i=1}^N \left[ y_i - \int_{\mathbb{R}^d} d\mathbf{s} \tilde{f}(\mathbf{s}) e^{i\mathbf{x}_i \cdot \mathbf{s}} \right]^2 &= 2 \sum_{i=1}^N [y_i - f(\mathbf{x}_i)] \int_{\mathbb{R}^d} d\mathbf{s} \frac{\delta \tilde{f}(\mathbf{s})}{\delta \tilde{f}(\mathbf{t})} e^{i\mathbf{x}_i \cdot \mathbf{s}} \\ &= 2 \sum_{i=1}^N [y_i - f(\mathbf{x}_i)] \int_{\mathbb{R}^d} d\mathbf{s} \delta(\mathbf{s} - \mathbf{t}) e^{i\mathbf{x}_i \cdot \mathbf{s}} \\ &= 2 \sum_{i=1}^N [y_i - f(\mathbf{x}_i)] e^{i\mathbf{x}_i \cdot \mathbf{t}} \end{aligned}$$

For the smoothness functional we have

$$\begin{aligned} \frac{\delta}{\delta \tilde{f}(\mathbf{t})} \int_{\mathbb{R}^d} d\mathbf{s} \frac{\tilde{f}(-\mathbf{s}) \tilde{f}(\mathbf{s})}{\tilde{G}(\mathbf{s})} &= 2 \int_{\mathbb{R}^d} d\mathbf{s} \frac{\tilde{f}(-\mathbf{s})}{\tilde{G}(\mathbf{s})} \frac{\delta \tilde{f}(\mathbf{s})}{\delta \tilde{f}(\mathbf{t})} \\ &= 2 \int_{\mathbb{R}^d} d\mathbf{s} \frac{\tilde{f}(-\mathbf{s})}{\tilde{G}(\mathbf{s})} \delta(\mathbf{s} - \mathbf{t}) \\ &= 2 \frac{\tilde{f}(-\mathbf{t})}{\tilde{G}(\mathbf{t})} \end{aligned}$$

Using these results we can now write equation A.2 as

$$\sum_{i=1}^N [y_i - f(\mathbf{x}_i)] e^{i\mathbf{x}_i \cdot \mathbf{t}} + \lambda \frac{\tilde{f}(-\mathbf{t})}{\tilde{G}(\mathbf{t})} = 0$$

Changing  $\mathbf{t}$  in  $-\mathbf{t}$  and multiplying by  $\tilde{G}(-\mathbf{t})$  on both sides of this equation we get

$$\tilde{f}(\mathbf{t}) = \tilde{G}(-\mathbf{t}) \sum_{i=1}^N \frac{[y_i - f(\mathbf{x}_i)]}{\lambda} e^{-i\mathbf{x}_i \cdot \mathbf{t}}$$

We now define the coefficients

$$c_i = \frac{[y_i - f(\mathbf{x}_i)]}{\lambda} \quad i = 1, \dots, N$$

assume that  $\tilde{G}$  is symmetric (so that its Fourier transform is real), and take the Fourier transform of the last equation, obtaining

$$f(\mathbf{x}) = \sum_{i=1}^N c_i \delta(\mathbf{x}_i - \mathbf{x}) * G(\mathbf{x}) = \sum_{i=1}^N c_i G(\mathbf{x} - \mathbf{x}_i)$$

We now recall that we had defined as equivalent all the functions differing by a term that lies in the null space of  $\phi[f]$ , and therefore the most general solution of the minimization problem is

$$f(\mathbf{x}) = \sum_{i=1}^N c_i G(\mathbf{x} - \mathbf{x}_i) + p(\mathbf{x})$$

where  $p(\mathbf{x})$  is a term that lies in the null space of  $\phi[f]$ , that is a set of polynomials for most common choices of stabilizer  $\phi[f]$ .

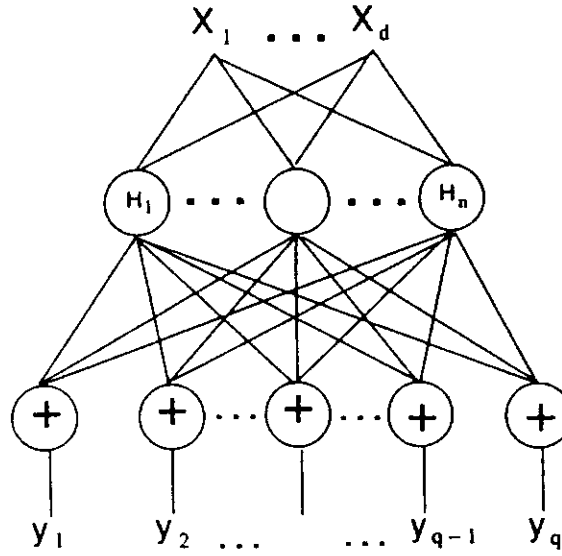


Figure 7: The most general network with one hidden layer and vector output. Notice that this approximation of a  $q$ -dimensional vector field has, in general, fewer parameters than the alternative representation consisting of  $q$  networks with one-dimensional outputs. If the only free parameters are the weights from the hidden layer to the output (as for simple RBF with  $n = N$ , where  $N$  is the number of examples) the two representations are equivalent.

## Appendix B: Approximation of Vector Fields with Regularization Networks

Consider the problem of approximating a  $q$ -dimensional vector field  $y(x)$  from a set of sparse data, the examples, which are pairs  $(x_i, y_i)$  for  $i = 1, \dots, N$ . Choose a *generalized regularization network* as the approximation scheme, that is, a network with one “hidden” layer and linear output units. Consider the case of  $N$  examples,  $n \leq N$  centers, input dimensionality  $d$  and output dimensionality  $q$  (see Fig. 7). Then the approximation is

$$y(x) = \sum_{\alpha=1}^n c_{\alpha} G(x - t_{\alpha}) \quad (B.1)$$

where  $G$  is the chosen basis function and the coefficients  $c_{\alpha}$  are now  $q$ -dimensional vectors:<sup>3</sup>  $c_{\alpha} = (c_{\alpha}^1, \dots, c_{\alpha}^q)$ .

<sup>3</sup>The components of an output vector will always be denoted by superscript, Greek indices.

Here we assume, for simplicity, that  $G$  is positive definite in order to avoid the need of additional polynomial terms in the previous equation. Equation B.1 can be rewritten in matrix notation as

$$\mathbf{y}(\mathbf{x}) = \mathbf{C}\mathbf{g}(\mathbf{x}) \quad (\text{B.2})$$

where the matrix  $\mathbf{C}$  is defined by  $(\mathbf{C})_{\mu,\alpha} = c_{\alpha}^{\mu}$  and  $\mathbf{g}$  is the vector with elements  $[\mathbf{g}(\mathbf{x})]_{\alpha} = G(\mathbf{x} - \mathbf{t}_{\alpha})$ . Assuming, for simplicity, that there is no noise in the data [that is equivalent to choosing  $\lambda = 0$  in the regularization functional (2.1)], the equations for the coefficients  $c_{\alpha}$  can be found imposing the interpolation conditions:

$$\mathbf{y}_i = \mathbf{C}\mathbf{g}(\mathbf{x}_i)$$

Introducing the following notation

$$(\mathbf{Y})_{i,\mu} = y_i^{\mu}(\mathbf{x}_i), \quad (\mathbf{C})_{\mu,\alpha} = c_{\alpha}^{\mu}, \quad (\mathbf{G})_{\alpha,i} = G(\mathbf{x}_i - \mathbf{t}_{\alpha})$$

the matrix of coefficients  $\mathbf{C}$  is given by

$$\mathbf{C} = \mathbf{Y}\mathbf{G}^{+}$$

where  $\mathbf{G}^{+}$  is the pseudoinverse of  $\mathbf{G}$  (Penrose 1955; Albert 1972). Substituting this expression in equation B.2, the following expression is obtained:

$$\mathbf{y}(\mathbf{x}) = \mathbf{Y}\mathbf{G}^{+}\mathbf{g}(\mathbf{x})$$

After some algebraic manipulations, this expression can be rewritten as

$$\mathbf{y}(\mathbf{x}) = \sum_{i=1}^N b_i(\mathbf{x})\mathbf{y}_i$$

where the functions  $b_i(\mathbf{x})$ , that are the elements of the vector  $\mathbf{b}(\mathbf{x})$ , depend on the chosen  $\mathbf{G}$ , according to

$$\mathbf{b}(\mathbf{x}) = \mathbf{G}^{+}\mathbf{g}(\mathbf{x})$$

Therefore, it follows (though it is not so well known) that the vector field  $\mathbf{y}(\mathbf{x})$  is approximated by the network as the linear combination of the example fields  $\mathbf{y}_i$ .

Thus for any choice of the regularization network and any choice of the (positive definite) basis function the estimated output vector is always a linear combination of the output example vectors with coefficients  $\mathbf{b}$  that depend on the input value. The result is valid for all networks with one hidden layer and linear outputs, provided that the mean square error criterion is used for training.

## Acknowledgments

---

We are grateful to P. Niyogi, H. Mhaskar, J. Friedman, J. Moody, V. Tresp, and one of the (anonymous) referees for useful discussions and suggestions. This paper describes research done within the Center for Biological and Computational Learning in the Department of Brain and Cognitive Sciences and at the Artificial Intelligence Laboratory at MIT. This research is sponsored by grants from the Office of Naval Research under contracts N00014-91-J-0385 and N00014-92-J-1879 and by a grant from the National Science Foundation under contract ASC-9217041 (which includes funds from ARPA provided under the HPCC program). Support for the A.I. Laboratory's artificial intelligence research is provided by ARPA-ONR contract N00014-91-J-4038. Tomaso Poggio is supported by the Uncas and Helen Whitaker Chair at the Whitaker College, Massachusetts Institute of Technology.

## References

---

- Aidu, F. A., and Vapnik, V. N. 1989. Estimation of probability density on the basis of the method of stochastic regularization. *Avtom. Telemek. (4)*, 84-97.
- Albert, A. 1972. *Regression and the Moore-Penrose Pseudoinverse*. Academic Press, New York.
- Allen, D. 1974. The relationship between variable selection and data augmentation and a method for prediction. *Technometrics* 16, 125-127.
- Aronszajn, N. 1950. Theory of reproducing kernels. *Trans. Am. Math. Soc.* 686, 337-404.
- Ballard, D. H. 1986. Cortical connections and parallel processing: structure and function. *Behav. Brain Sci.* 9, 67-120.
- Barron, A. R., and Barron, R. L. 1988. Statistical learning networks: A unifying view. In *Symposium on the Interface: Statistics and Computing Science*, Reston, Virginia.
- Barron, A. R. 1991. *Approximation and estimation bounds for artificial neural networks*. Tech. Rep. 59, Department of Statistics, University of Illinois at Urbana-Champaign, Champaign, IL.
- Barron, A. R. 1993. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Transact. Inform. Theory* 39(3), 930-945.
- Barron, A. R. 1994. Approximation and estimation bounds for artificial neural networks. *Machine Learn.* 14, 115-133.
- Baum, E. B. 1988. On the capabilities of multilayer perceptrons. *J. Complex.* 4, 193-215.
- Baum, E. B., and Haussler, D. 1989. What size net gives valid generalization? *Neural Comp.* 1, 151-160.
- Bellman, R. E. 1961. *Adaptive Control Processes*. Princeton University Press, Princeton, NJ.
- Bertero, M. 1986. Regularization methods for linear inverse problems. In *Inverse Problems*, C. G. Talenti, ed. Springer-Verlag, Berlin.

- Bertero, M., Poggio, T., and Torre, V. 1988. Ill-posed problems in early vision. *Proc. IEEE* 76, 869–889.
- Bottou, L., and Vapnik, V. 1992. Local learning algorithms. *Neural Comp.* 4(6), 888–900.
- Breiman, L. 1993. Hinging hyperplanes for regression, classification, and function approximation. *IEEE Trans. Inform. Theory* 39(3), 999–1013.
- Broomhead, D. S., and Lowe, D. 1988. Multivariable functional interpolation and adaptive networks. *Complex Syst.* 2, 321–355.
- Buhmann, M. D. 1990. Multivariate cardinal interpolation with radial basis functions. *Construct. Approx.* 6, 225–255.
- Buhmann, M. D. 1991. On quasi-interpolation with radial basis functions. Numerical Analysis Reports DAMPT 1991/NA3, Department of Applied Mathematics and Theoretical Physics, Cambridge, England.
- Buja, A., Hastie, T., and Tibshirani, R. 1989. Linear smoothers and additive models. *Ann. Statist.* 17, 453–555.
- Caprile, B., and Girosi, F. 1990. A nondeterministic minimization algorithm. A.I. Memo 1254, Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA.
- Cox, D. D. 1984. Multivariate smoothing spline functions. *SIAM J. Numer. Anal.* 21, 789–813.
- Craven, P., and Wahba, G. 1979. Smoothing noisy data with spline functions: Estimating the correct degree of smoothing by the method of generalized cross validation. *Numer. Math.* 31, 377–403.
- Cybenko, G. 1989. Approximation by superposition of a sigmoidal function. *Math. Control Systems Signals* 2(4), 303–314.
- de Boor, C. 1978. *A Practical Guide to Splines*. Springer-Verlag, New York.
- de Boor, C. 1990. Quasi-interpolants and approximation power of multivariate splines. In *Computation of Curves and Surfaces*, M. Gasca and C. A. Micchelli, eds., pp. 313–345. Kluwer Academic Publishers, Dordrecht, Netherlands.
- DeVore, R. A. 1991. Degree of nonlinear approximation. In *Approximation Theory, VI*, C. K. Chui, L. L. Schumaker, and D. J. Ward, eds., pp. 175–201. Academic Press, New York.
- DeVore, R. A., and Yu, X. M. 1991. Nonlinear  $n$ -widths in Besov spaces. In *Approximation Theory, VI*, C. K. Chui, L. L. Schumaker, and D. J. Ward, eds., pp. 203–206. Academic Press, New York.
- DeVore, R., Howard, R., and Micchelli, C. 1989. Optimal nonlinear approximation. *Manuskrip. Math.*
- Devroye, L. P., and Wagner, T. J. 1980. Distribution-free consistency results in nonparametric discrimination and regression function estimation. *Ann. Statist.* 8, 231–239.
- Diaconis, P., and Freedman, D. 1984. Asymptotics of graphical projection pursuit. *Ann. Statist.* 12(3), 793–815.
- Donoho, D. L., and Johnstone, I. M. 1989. Projection-based approximation and a duality with kernel methods. *Ann. Statist.* 17(1), 58–106.
- Duchon, J. 1977. Spline minimizing rotation-invariant semi-norms in Sobolev spaces. In *Constructive Theory of Functions of Several Variables, Lecture Notes in Mathematics*, 571, W. Schempp and K. Zeller, eds. Springer-Verlag, Berlin.

- Dyn, N. 1987. Interpolation of scattered data by radial functions. In *Topics in Multivariate Approximation*, C. K. Chui, L. L. Schumaker, and F. I. Utreras, eds. Academic Press, New York.
- Dyn, N. 1991. Interpolation and approximation by radial and related functions. In *Approximation Theory, VI*, C. K. Chui, L. L. Schumaker, and D. J. Ward, eds., pp. 211–234. Academic Press, New York.
- Dyn, N., Levin, D., and Rippa, S. 1986. Numerical procedures for surface fitting of scattered data by radial functions. *SIAM J. Sci. Stat. Comput.* 7(2), 639–659.
- Dyn, N., Jackson, I. R. H., Levin, D., and Ron, A. 1989. On multivariate approximation by integer translates of a basis function. Computer Sciences Tech. Rep. 886, University of Wisconsin–Madison.
- Eubank, R. L. 1988. *Spline Smoothing and Nonparametric Regression*, Vol. 90 of *Statistics, Textbooks and Monographs*. Marcel Dekker, Basel.
- Franke, R. 1982. Scattered data interpolation: Tests of some method. *Math. Comp.* 38(5), 181–200.
- Franke, R. 1987. Recent advances in the approximation of surfaces from scattered data. In *Topics in Multivariate Approximation*, C. K. Chui, L. L. Schumaker, and F. I. Utreras, eds. Academic Press, New York.
- Friedman, J. H., and Stuetzle, W. 1981. Projection pursuit regression. *J. Am. Statist. Assoc.* 76(376), 817–823.
- Funahashi, I. 1989. On the approximate realization of continuous mappings by neural networks. *Neural Networks* 2, 183–192.
- Gasser, Th., and Müller, H. G. 1985. Estimating regression functions and their derivatives by the kernel method. *Scand. J. Statist.* 11, 171–185.
- Gelfand, I. M., and Shilov, G. E. 1964. *Generalized Functions. Vol. 1: Properties and Operations*. Academic Press, New York.
- Geman, S., Bienenstock, E., and Doursat, R. 1992. Neural networks and the bias/variance dilemma. *Neural Comp.* 4, 1–58.
- Girosi, F. 1991. Models of noise and robust estimates. A.I. Memo 1287, Artificial Intelligence Laboratory, Massachusetts Institute of Technology.
- Girosi, F. 1992. On some extensions of radial basis functions and their applications in artificial intelligence. *Comput. Math. Applic.* 24(12), 61–80.
- Girosi, F. 1993. Regularization theory, radial basis functions and networks. In *From Statistics to Neural Networks. Theory and Pattern Recognition Applications*, V. Cherkassky, J. H. Friedman, and H. Wechsler, eds. Subseries F, Computer and Systems Sciences. Springer-Verlag, Berlin.
- Girosi, F., and Anzellotti, G. 1992. Rates of convergence of approximation by translates. A.I. Memo 1288, Artificial Intelligence Laboratory, Massachusetts Institute of Technology.
- Girosi, F., and Anzellotti, G. 1993. Rates of convergence for radial basis functions and neural networks. In *Artificial Neural Networks for Speech and Vision*, R. J. Mammone, ed., pp. 97–113. Chapman & Hall, London.
- Girosi, F., and Poggio, T. 1990. Networks and the best approximation property. *Biol. Cybernet.* 63, 169–176.
- Girosi, F., Poggio, T., and Caprile, B. 1991. Extensions of a theory of networks for approximation and learning: Outliers and negative examples. In *Ad-*

- vances in Neural Information Processings Systems 3*, R. Lippmann, J. Moody, and D. Touretzky, eds. Morgan Kaufmann, San Mateo, CA.
- Girosi, F., Jones, M., and Poggio, T. 1993. Priors, stabilizers and basis functions: From regularization to radial, tensor and additive splines. A.I. Memo No. 1430, Artificial Intelligence Laboratory, Massachusetts Institute of Technology.
- Golub, G., Heath, M., and Wahba, G. 1979. Generalized cross validation as a method for choosing a good ridge parameter. *Technometrics* 21, 215–224.
- Grimson, W. E. L. 1982. A computational theory of visual surface interpolation. *Proc. R. S. London B* 298, 395–427.
- Harder, R. L., and Desmarais, R. M. 1972. Interpolation using surface splines. *J. Aircraft* 9, 189–191.
- Härdle, W. 1990. *Applied Nonparametric Regression*, Vol. 19 of *Econometric Society Monographs*. Cambridge University Press, Cambridge.
- Hardy, R. L. 1971. Multiquadric equations of topography and other irregular surfaces. *J. Geophys. Res.* 76, 1905–1915.
- Hardy, R. L. 1990. Theory and applications of the multiquadric-biharmonic method. *Computers Math. Applic.* 19(8/9), 163–208.
- Hastie, T., and Tibshirani, R. 1986. Generalized additive models. *Statist. Sci.* 1, 297–318.
- Hastie, T., and Tibshirani, R. 1987. Generalized additive models: Some applications. *J. Am. Statist. Assoc.* 82, 371–386.
- Hastie, T., and Tibshirani, R. 1990. *Generalized Additive Models*, Vol. 43 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, London.
- Hausser, D. 1989. Decision theoretic generalizations of the PAC model for neural net and other learning applications. Tech. Rep. UCSC-CRL-91-02, University of California, Santa Cruz.
- Hertz, J. A., Krogh, A., and Palmer, R. 1991. *Introduction to the Theory of Neural Computation*. Addison-Wesley, Redwood City, CA.
- Hornik, K., Stinchcombe, M., and White, W. 1989. Multilayer feedforward networks are universal approximators. *Neural Networks* 2, 359–366.
- Huber, P. J. 1985. Projection pursuit. *Ann. Statist.* 13(2), 435–475.
- Hurlbert, A., and Poggio, T. 1988. Synthesizing a color algorithm from examples. *Science* 239, 482–485.
- Irie, B., and Miyake, S. 1988. Capabilities of three-layered perceptrons. *IEEE Int. Conf. Neural Networks* 1, 641–648.
- Jackson, I. R. H. 1988. *Radial basis functions methods for multivariate approximation*. Ph.D. thesis, University of Cambridge, U.K.
- Jones, L. K. 1992. A simple lemma on greedy approximation in Hilbert space and convergence rates for projection pursuit regression and neural network training. *Ann. Statist.* 20(1), 608–613.
- Kansa, E. J. 1990a. Multiquadrics—A scattered data approximation scheme with applications to computational fluid dynamics—I. *Comput. Math. Applic.* 19(8/9), 127–145.
- Kansa, E. J. 1990b. Multiquadrics—A scattered data approximation scheme with applications to computational fluid dynamics—II. *Comput. Math. Applic.* 19(8/9), 147–161.

- Kimeldorf, G. S., and Wahba, G. 1971. A correspondence between Bayesian estimation on stochastic processes and smoothing by splines. *Ann. Math. Statist.* 2, 495–502.
- Kohonen, T. 1990. The self-organizing map. *Proc. IEEE* 78(9), 1464–1480.
- Kung, S. Y. 1993. *Digital Neural Networks*. Prentice Hall, Englewood Cliffs, NJ.
- Lancaster, P., and Salkauskas, K. 1986. *Curve and Surface Fitting*. Academic Press, London.
- Lapedes, A., and Farber, R. 1988. How neural nets work. In *Neural Information Processing Systems*, D. Z. Anderson, ed., pp. 442–456. American Institute of Physics, New York.
- Lippmann, R. P. 1989. Review of neural networks for speech recognition. *Neural Comp.* 1, 1–38.
- Lippmann, R. P., and Lee, Y. 1991. A critical overview of neural network pattern classifiers. Presented at Neural Networks for Computing Conference, Snowbird, UT.
- Lorentz, G. G. 1962. Metric entropy, widths, and superposition of functions. *Am. Math. Monthly* 69, 469–485.
- Lorentz, G. G. 1986. *Approximation of Functions*. Chelsea, New York.
- Madych, W. R., and Nelson, S. A. 1990a. Multivariate interpolation and conditionally positive definite functions. II. *Math. Comput.* 54(189), 211–230.
- Madych, W. R., and Nelson, S. A. 1990b. Polyharmonic cardinal splines: A minimization property. *J. Approx. Theory* 63, 303–320.
- Marroquin, J. L., Mitter, S., and Poggio, T. 1987. Probabilistic solution of ill-posed problems in computational vision. *J. Am. Stat. Assoc.* 82, 76–89.
- Maruyama, M., Girosi, F., and Poggio, T. 1992. A connection between HBF and MLP. A.I. Memo No. 1291, Artificial Intelligence Laboratory, Massachusetts Institute of Technology.
- Meinguet, J. 1979. Multivariate interpolation at arbitrary points made simple. *J. Appl. Math. Phys.* 30, 292–304.
- Mel, B. W. 1988. MURPHY: A robot that learns by doing. In *Neural Information Processing Systems*, D. Z. Anderson, ed. American Institute of Physics, New York.
- Mel, B. W. 1990. The sigma-pi column: A model of associative learning in cerebral neocortex. Tech. Rep. 6, California Institute of Technology.
- Mel, B. W. 1992. NMDA-based pattern-discrimination in a modeled cortical neuron. *Neural Comp.* 4, 502–517.
- Mhaskar, H. N. 1993a. Approximation properties of a multilayered feedforward artificial neural network. *Adv. Comp. Math.* 1, 61–80.
- Mhaskar, H. N. 1993b. Neural networks for localized approximation of real functions. In *Neural Networks for Signal Processing III, Proceedings of the 1993 IEEE-SP Workshop*, C. A. Kamm et al., eds., pp. 190–196. IEEE Signal Processing Society, New York.
- Mhaskar, H. N., and Micchelli, C. A. 1992. Approximation by superposition of a sigmoidal function. *Adv. Appl. Math.* 13, 350–373.
- Mhaskar, H. N., and Micchelli, C. A. 1993. How to choose an activation function. In *Advances in Neural Information Processing Systems 5*, S. J. Hanson, J. D. Cowan, and C. L. Giles, eds. Morgan Kaufmann, San Mateo, CA.

- Micchelli, C. A. 1986. Interpolation of scattered data: Distance matrices and conditionally positive definite functions. *Construct. Approx.* 2, 11–22.
- Moody, J. 1991a. Note on generalization, regularization, and architecture selection in nonlinear learning systems. In *Proceedings of the First IEEE-SP Workshop on Neural Networks for Signal Processing*, pp. 1–10. IEEE Computer Society Press, Los Alamitos, CA.
- Moody, J. 1991b. The effective number of parameters: an analysis of generalization and regularization in nonlinear learning systems. In *Advances in Neural Information Processing Systems 4*, J. Moody, S. Hanson, and R. Lippmann, eds., pp. 847–854. Morgan Kaufmann, Palo Alto, CA.
- Moody, J., and Darken, C. 1988. Learning with localized receptive fields. In *Proceedings of the 1988 Connectionist Models Summer School*, G. Hinton, T. Sejnowski, and D. Touretzky, eds., pp. 133–143. Palo Alto, CA.
- Moody, J., and Darken, C. 1989. Fast learning in networks of locally-tuned processing units. *Neural Comp.* 1(2), 281–294.
- Moody, J., and Yarvin, N. 1991. Networks with learned unit response functions. In *Advances in Neural Information Processing Systems 4*, J. Moody, S. Hanson, and R. Lippmann, eds., pp. 1048–1055. Morgan Kaufmann, Palo Alto, CA.
- Morozov, V. A. 1984. *Methods for Solving Incorrectly Posed Problems*. Springer-Verlag, Berlin.
- Nadaraya, E. A. 1964. On estimating regression. *Theor. Prob. Appl.* 9, 141–142.
- Niyogi, P., and Girosi, F. 1994. On the relationship between generalization error, hypothesis complexity, and sample complexity for radial basis functions. A.I. Memo 1467, Artificial Intelligence Laboratory, Massachusetts Institute of Technology.
- Omohundro, S. 1987. Efficient algorithms with neural network behaviour. *Complex Syst.* 1, 273.
- Parisi, G. 1988. *Statistical Field Theory*. Addison-Wesley, Reading, MA.
- Parzen, E. On estimation of a probability density function and mode. *Ann. Math. Statist.* 33, 1065–1076.
- Penrose, R. 1955. A generalized inverse for matrices. *Proc. Cambridge Philos. Soc.* 51, 406–413.
- Pinkus, A. 1986. *N-widths in Approximation Theory*. Springer-Verlag, New York.
- Poggio, T. 1975. On optimal nonlinear associative recall. *Biol. Cybernet.* 19, 201–209.
- Poggio, T. 1990. A theory of how the brain might work. *Cold Spring Harbor Symp. Quantit. Biol.* 899–910.
- Poggio, T., and Girosi, F. 1989. A theory of networks for approximation and learning. A.I. Memo No. 1140, Artificial Intelligence Laboratory, Massachusetts Institute of Technology.
- Poggio, T., and Girosi, F. 1990a. Networks for approximation and learning. *Proc. IEEE* 78(9).
- Poggio, T., and Girosi, F. 1990b. Extension of a theory of networks for approximation and learning: dimensionality reduction and clustering. In *Proceedings Image Understanding Workshop*, pp. 597–603, Pittsburgh, Pennsylvania, September 11–13. Morgan Kaufmann, Palo Alto, CA.

- Poggio, T., and Girosi, F. 1990c. Regularization algorithms for learning that are equivalent to multilayer networks. *Science* **247**, 978–982.
- Poggio, T., and Hurlbert, A. 1994. Observation on cortical mechanisms for object recognition and learning. In *Large-Scale Neuronal Theories of the Brain*, C. Koch and J. Davis, eds. In press.
- Poggio, T., Torre, V., and Koch, C. 1985. Computational vision and regularization theory. *Nature* **317**, 314–319.
- Poggio, T., Voorhees, H., and Yuille, A. 1988. A regularized solution to edge detection. *J. Complex.* **4**, 106–123.
- Pollard, D. 1984. *Convergence of Stochastic Processes*. Springer-Verlag, Berlin.
- Powell, M. J. D. 1987. Radial basis functions for multivariable interpolation: A review. In *Algorithms for Approximation*, J. C. Mason and M. G. Cox, eds. Clarendon Press, Oxford.
- Powell, M. J. D. 1992. The theory of radial basis functions approximation in 1990. In *Advances in Numerical Analysis Volume II: Wavelets, Subdivision Algorithms and Radial Basis Functions*, W. A. Light, ed., pp. 105–210. Oxford University Press, Oxford.
- Priestley, M. B., and Chao, M. T. 1972. Non-parametric function fitting. *J. R. Statist. Soc. B* **34**, 385–392.
- Rabut, C. 1991. How to build quasi-interpolants. applications to polyharmonic B-splines. In *Curves and Surfaces*, P.-J. Laurent, A. Le Mehaute, and L. L. Schumaker, eds., pp. 391–402. Academic Press, New York.
- Rabut, C. 1992. An introduction to Schoenberg's approximation. *Comput. Math. Applic.* **24**(12), 149–175.
- Ripley, B. D. 1994. Neural networks and related methods for classification. *Proc. R. Soc. London*, in press.
- Rissanen, J. 1978. Modeling by shortest data description. *Automatica* **14**, 465–471.
- Rosenblatt, M. 1971. Curve estimates. *Ann. Math. Statist.* **64**, 1815–1842.
- Rumelhart, D. E., Hinton, G. E., and Williams, R. J. 1986. Learning representations by back-propagating errors. *Nature (London)* **323**(9), 533–536.
- Schoenberg, I. J. 1946a. Contributions to the problem of approximation of equidistant data by analytic functions, part a: On the problem of smoothing of graduation, a first class of analytic approximation formulae. *Quart. Appl. Math.* **4**, 45–99.
- Schoenberg, I. J. 1969. Cardinal interpolation and spline functions. *J. Approx. Theory* **2**, 167–206.
- Schumaker, L. L. 1981. *Spline Functions: Basic Theory*. John Wiley, New York.
- Sejnowski, T. J., and Rosenberg, C. R. 1987. Parallel networks that learn to pronounce English text. *Complex Syst.* **1**, 145–168.
- Silverman, B. W. 1984. Spline smoothing: The equivalent variable kernel method. *Ann. Statist.* **12**, 898–916.
- Sivakumar, N., and Ward, J. D. 1991. On the best least square fit by radial functions to multidimensional scattered data. Tech. Rep. 251, Center for Approximation Theory, Texas A&M University.
- Solomonoff, R. J. 1978. Complexity-based induction systems: Comparison and convergence theorems. *IEEE Trans. Inform. Theory* **24**.

- Specht, D. F. 1991. A general regression neural network. *IEEE Trans. Neural Networks* 2(6), 568–576.
- Stein, E. M. 1970. *Singular Integrals and Differentiability Properties of Functions*. Princeton University Press, Princeton, NJ.
- Stewart, J. 1976. Positive definite functions and generalizations, an historical survey. *Rocky Mountain J. Math.* 6, 409–434.
- Stone, C. J. 1985. Additive regression and other nonparametric models. *Ann. Statist.* 13, 689–705.
- Tikhonov, A. N. 1963. Solution of incorrectly formulated problems and the regularization method. *Soviet Math. Dokl.* 4, 1035–1038.
- Tikhonov, A. N., and Arsenin, V. Y. 1977. *Solutions of Ill-Posed Problems*. W. H. Winston, Washington, DC.
- Timan, A. F. 1963. *Theory of Approximation of Functions of a Real Variable*. Macmillan, New York.
- Tresp, V., Hollatz, J., and Ahmad, S. 1993. Network structuring and training using rule-based knowledge. In *Advances in Neural Information Processing Systems 5*, S. J. Hanson, J. D. Cowan, and C. L. Giles, eds. Morgan Kaufmann, San Mateo, CA.
- Utreras, F. 1979. Cross-validation techniques for smoothing spline functions in one or two dimensions. In *Smoothing Techniques for Curve Estimation*, T. Gasser and M. Rosenblatt, eds., pp. 196–231. Springer-Verlag, Heidelberg.
- Vapnik, V. N. 1982. *Estimation of Dependences Based on Empirical Data*. Springer-Verlag, Berlin.
- Vapnik, V. N., and Chervonenkis, A. Y. 1971. On the uniform convergence of relative frequencies of events to their probabilities. *Th. Prob. Applic.* 17(2), 264–280.
- Vapnik, V. N., and Chervonenkis, A. Y. 1981. The necessary and sufficient conditions for the uniform convergence of averages to their expected values. *Teor. Veroyat. Primen.* 26(3), 543–564.
- Vapnik, V. N., and Chervonenkis, A. Y. 1991. The necessary and sufficient conditions for consistency in the empirical risk minimization method. *Pattern Recog. Image Anal.* 1(3), 283–305.
- Vapnik, V. N., and Stefanyuk, A. R. 1978. Nonparametric methods for restoring probability densities. *Avtomat. Telemek.* 8, 38–52.
- Wahba, G. 1975. Smoothing noisy data by spline functions. *Numer. Math* 24, 383–393.
- Wahba, G. 1979. Smoothing and ill-posed problems. In *Solutions Methods for Integral Equations and Applications*, M. Golberg, ed., pp. 183–194. Plenum Press, New York.
- Wahba, G. 1980. Spline bases, regularization, and generalized cross-validation for solving approximation problems with large quantities of noisy data. In *Proceedings of the International Conference on Approximation Theory in Honour of George Lorenz*, J. Ward and E. Cheney, eds., Austin, TX, January 8–10, 1980. Academic Press, New York.
- Wahba, G. 1985. A comparison of GCV and GML for choosing the smoothing parameter in the generalized splines smoothing problem. *Ann. Statist.* 13, 1378–1402.

- Wahba, G. 1990. *Splines Models for Observational Data*. Series in Applied Mathematics, Vol. 59, SIAM, Philadelphia.
- Wahba, G., and Wold, S. 1975. A completely automatic French curve. *Commun. Statist.* 4, 1-17.
- Watson, G. S. 1964. Smooth regression analysis. *Sankhya A* 26, 359-372.
- White, H. 1989. Learning in artificial neural networks: A statistical perspective. *Neural Comp.* 1, 425-464.
- White, H. 1990. Connectionist nonparametric regression: Multilayer perceptrons can learn arbitrary mappings. *Neural Networks* 3, 535-549.
- Yuille, A., and Grzywacz, N. 1988. The motion coherence theory. In *Proceedings of the International Conference on Computer Vision*, pp. 344-354, IEEE Computer Society Press, Washington, DC.
- Ziemer, W. P. 1989. *Weakly Differentiable Functions: Sobolev Spaces and Functions of Bounded Variation*. Springer-Verlag, New York.

---

Received February 2, 1994; accepted June 22, 1994.