

Regularization Networks

9.520 Class 18, 2004

Tomaso Poggio

Plan

- Radial Basis Functions and their extensions
- Additive Models
- Regularization Networks
- Dual Kernels
- Conclusions

About this class

We describe a family of regularization techniques based on radial kernels K and called RBFs. We introduce RBF extensions (somewhat less rigorous) such as Hy-per Basis Functions and characterize their relation with other techniques including MLPs and splines.

Radial Basis Functions

Radial Basis Functions, as MLPs, have the universal approximation property.

Theorem: Let K be a Radial Basis Function and I_i the n -dimensional cube $[0, 1]^n$. Then finite sums of the form

$$f(\mathbf{x}) = \sum_N^{i=1} c_i K(\mathbf{x} - \mathbf{x}_i)$$

are dense in $C[I_i]$. In other words, given a function $h \in C[I_i]$ and $\epsilon > 0$, there is a sum, $f(\mathbf{x})$, of the above form, for which:

$$|f(\mathbf{x}) - h(\mathbf{x})| < \epsilon \text{ for all } \mathbf{x} \in I_n.$$

Notice that RBF correspond to RKHS defined on an infinite domain. Notice also that RKHS do not in general have the same approximation property: RKHS generated by a K with an infinite countable number of strictly positive eigenvalues are dense in L_2 but not necessarily in $C(X)$, though they can be embedded in $C(X)$. See Zhou results.

Density of a RKHS on a bounded domain (the non-RBF case)

We first ask under which condition is a RKHS dense in $L_2(X, \nu)$.

1. when L_K is strictly positive the RKHS is infinite dimensional and dense in $L_2(X, \nu)$.
2. in the degenerate case the RKHS is finite dimensional and not dense in $L_2(X, \nu)$.
3. in the conditionally strictly positive case the RKHS is not dense in $L_2(X, \nu)$ but when completed with a finite number of polynomials of appropriate degree can be made to be dense in $L_2(X, \nu)$.

Density of a RKHS on a bounded domain (cont)

Density of RKHS – defined on a compact domain X – in $C(X)$ (in the sup norm) is a trickier issue that has been answered very recently by Zhou (in preparation). It is however guaranteed for radial kernels K for K continuous and integrable, if density in $L^2(X, \nu)$ holds (with X the infinite domain). These are facts for radial kernels and unrelated to RKHS properties

- $\text{span } K(x - y) : y \in \mathbb{R}^n$ is dense in $L^2(\mathbb{R}^n)$ iff the Fourier transform of K goes not vanish on set of positive Lebesgue measure (N. Wiener).

- $\text{span } K(x - y) : y \in \mathbb{R}^n$ is dense in $C(\mathbb{R}^n)$ (topology of uniform convergence) if $K \in C(\mathbb{R}^n)$, $K \in L^1_+(\mathbb{R}^n)$.

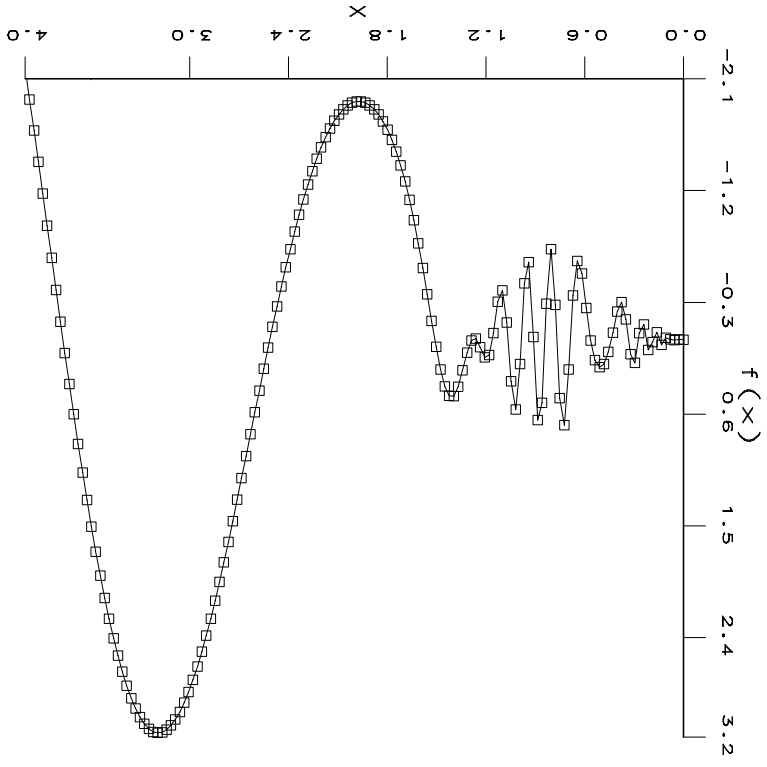
Some good properties of RBF

- Well motivated in the framework of regularization theory;
- The solution is unique and equivalent to solving a linear system;
- Degree of smoothness is tunable (with λ);
- Universal approximation property;
- Large body of applied math literature on the subject;
- Interpretation in terms of *neural networks*(?);
- Biologically plausible;
- Simple interpretation in terms of *smooth look-up table*;
- Similar to other non-parametric techniques, such as nearest neighbor and kernel regression (see end of this class).

Some not-so-good properties of RBF

- Computationally expensive ($O(\ell^3)$);
- Linear system to be solved for finding the coefficients often badly ill-conditioned;
- The same degree of smoothness is imposed on different regions of the domain (we will see how to deal with this problem in the class on wavelets);

This function has different smoothness properties in different regions of its domain.



A first extension: less centers than data points

We look for an *approximation* to the regularization solution:

$$\sum_{\ell} c_{\ell} K(\mathbf{x} - \mathbf{x}_{\ell}) = f(\mathbf{x})$$

↑

$$\sum_{m} c_m K(\mathbf{x} - \mathbf{t}_m) = f_*(\mathbf{x})$$

where $m > \ell$ and the vectors \mathbf{t}_m are called **centers**.

Homework: show that the interpolation problem is still well-posed when $m > \ell$.

(Broomhead and Lowe, 1988; Moody and Darken, 1989; Poggio and Girosi, 1989)

Least Squares Regularization Networks

$$f_*(\mathbf{x}) = \sum_m^{v=1} c^v K(\mathbf{x} - \mathbf{t}^v)$$

Suppose the centers \mathbf{t}^v have been fixed.

How do we find the coefficients c^v ?



Least Squares

Finding the coefficients

Define

$$E(c_1, \dots, c_m) = \sum_{i=1}^{\ell} (y_i - f_*(\mathbf{x}_i))^2$$

The least squares criterion is

$$\min_{c_1, \dots, c_m} E(c_1, \dots, c_m)$$

The problem is convex and quadratic in the c_α , and the solution satisfies:

$$\frac{\partial E}{\partial c_\alpha} = 0$$

Finding the centers

Given the centers t_α we know how to find the c_α .

How do we choose the t_α ?

1. a subset of the examples (random);

2. by a clustering algorithm (k-means, for example);

3. by least squares (*moving centers*);

4. a subset of the examples: Support Vector Machines;

Centers as a subset of the examples

Fair technique. The subset is a random subset, which should reflect the distribution of the data.

Not many theoretical results available (but we proved that solution exists since matrix is full rank).

Main problem: how many centers?

Main answer: we don't know. Cross validation techniques seem a reasonable choice.

Finding the centers by clustering

Very common. However it makes sense only if the input data points are clustered.

No theoretical results.

Not clear that it is a good idea, especially for pattern classification cases.

Moving centers

Define

$$E(c_1, \dots, c_m, t_1, \dots, t_m) = \sum_{i=1}^{\ell} (y_i - f_*(\mathbf{x}_i))^2$$

The least squares criterion is

$$\min_{c, t} E(c_1, \dots, c_m, t_1, \dots, t_m).$$

The problem is not convex and quadratic anymore: expect multiple local minima.

Moving centers

(-) Very flexible, in principle very powerful (more than SVMs);

(-) Some theoretical understanding;

(-) Very expensive computationally due to the local minima problem;

(-) Centers sometimes move in "weird" ways;

Connection with MLP

Radial Basis Functions with moving centers is a particular case of a function approximation technique of the form:

$$f(\mathbf{x}) = \sum_{i=1}^N c_i H(\mathbf{x}, \mathbf{p}_i)$$

where the parameters \mathbf{p}_i can be estimated by least squares techniques.

Radial Basis Functions corresponds to the choice $N = m$ and $\mathbf{p}_i = \mathbf{t}_i$, and

$$H(\mathbf{x}, \mathbf{p}_i) = K(\|\mathbf{x} - \mathbf{t}_i\|)$$

Extensions of Radial Basis Functions (much beyond what SVMs can do)

- Different variables can have different scales: $f(x, y) = y^2 \sin(100x)$;

- Different variables could have different units of measure $f = f(\ddot{x}, \dot{x}, \underline{x})$;

- Not all the variables are independent or relevant: $f(x, y, z, t) = g(y, x, z, y)$;

- Only some linear combinations of the variables are relevant: $f(x, y, z) = \sin(x + y + z)$;

Extensions of regularization theory

A priori knowledge:

- the relevant variables are linear combination of the original ones:

$$\mathbf{z} = W\mathbf{x}$$

for some (possibly rectangular) matrix W ;

- $f(\mathbf{x}) = g(W\mathbf{x}) = g(\mathbf{z})$ and the function g is smooth;

The regularization functional is now

$$\sum_{i=1}^n (y_i - g(\mathbf{z}))^2 + \lambda \Phi[g]$$

where $\mathbf{z}_i = W\mathbf{x}_i$.

Extensions of regularization theory (continue)

The solution is

$$\cdot \sum_{\gamma}^{l=1} c_{\gamma} K(z) = g(z)$$

Therefore the solution for f is:

$$(\mathbf{x}_M - \mathbf{x}_M) K \sum_{\gamma}^{l=1} = (\mathbf{x}_M) g = (\mathbf{x}) f$$

Extensions of regularization theory (continue)

If the matrix W were known, the coefficients could be computed as in the radial case:

$$(K + \lambda I)c = y$$

where

$$(y)_i = y_i, \quad (c)_i = c_i, \quad (K)_{ij} = K(Wx_i - Wx_j)$$

and the same argument of the Regularization Networks technique apply, leading to *Generalized Regularization Net-*

works:

$$\sum_{m=1}^{\infty} c_m K(Wx - Wt_m) = f_*(x)$$

Extensions of regularization theory (continue)

Since W is usually not known, it could be found by *least squares*. Define

$$E(c_1, \dots, c_m, W) = \sum_{i=1}^{\ell} (y_i - f_*(\mathbf{x}_i))^2$$

Then we can solve:

$$\min_{c_\alpha, W} E(c_1, \dots, c_m, W)$$

The problem is not convex and quadratic anymore: expect multiple local minima.

From RBF to HyperBF

When the basis function K is radial the Generalized Regularization Networks becomes

$$f(\mathbf{x}) = \sum_m^{\alpha=1} c_m K(\|\mathbf{x} - \mathbf{t}_m\|)$$

that is a *non radial basis function* technique (we define $\|\mathbf{x}\|_2 = \mathbf{x}^\top W^\top W \mathbf{x}$).

Least Squares

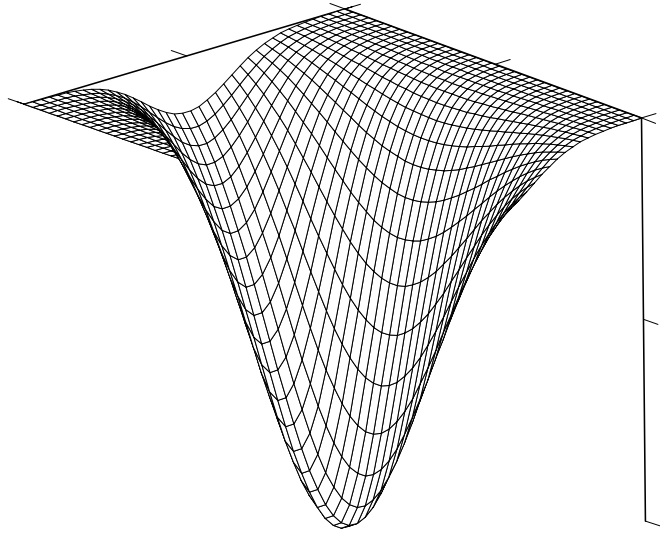
$$1. \min_{c_\alpha} E(c_1, \dots, c_m)$$

$$2. \min_{c_\alpha, t_\alpha} E(c_1, \dots, c_m, t_1, \dots, t_m)$$

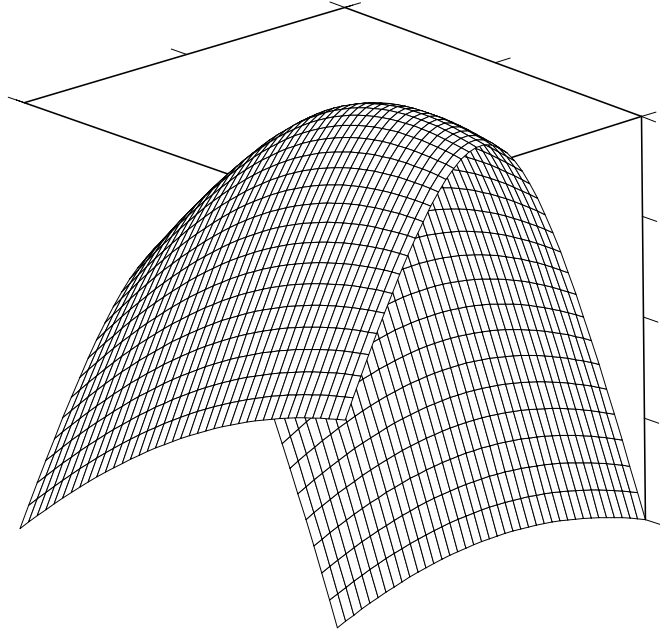
$$3. \min_{c_\alpha, W} E(c_1, \dots, c_m, W)$$

$$4. \min_{c_\alpha, t_\alpha, W} E(c_1, \dots, c_m, t_1, \dots, t_m, W)$$

A nonradial Gaussian function



A nonradial multiquadric function



Additive models

In statistics an additive model has the form

$$(\mathbf{x})f = \sum_{p=1}^n (\mathbf{x})f_p$$

where

$$(\mathbf{x})f_p = \sum_{j=1}^p (\mathbf{x})f_{pj}$$

In other words

$$(\mathbf{x})f = \sum_{j=1}^p (\mathbf{x})f_{pj} = \sum_{j=1}^p \sum_{l=1}^n (\mathbf{x})f_{pjl}$$

Additive stabilizers

To obtain an approximation of the form

$$f(\mathbf{x}) = \sum_{p=1}^n f_p(\mathbf{x})$$

We choose a stabilizer corresponding to an additive basis function

$$K(\mathbf{x}) = \sum_{p=1}^n \theta_p K_p(\mathbf{x})$$

This scheme leads to an approximation scheme of the additive form with

$$f(\mathbf{x}) = \sum_{l=1}^n \theta_l K_l(\mathbf{x})$$

Notice that the additive components are not independent since there is only one set of c_l – which makes sense since I have only l data points to determine the c_l .

Extensions of Additive Models

We start from the non-independent additive component formulation obtained from additive stabilizers

$$f(\mathbf{x}) = \sum_{\ell=1}^i c_{\ell} \theta^{\ell} K^{\ell} (x_{\ell} - \theta^{\ell} x)$$

We assume now that the parameters θ^{ℓ} are free, eg different for each i . We now have to fit

$$f(\mathbf{x}) = \sum_{\ell=1}^i \sum_{p=1}^{\ell} c_{\ell}^p (x_{\ell} - \theta^{\ell} x)$$

with $\ell \times d$ independent c_{ℓ}^p . In order to avoid overfitting we

reduce the number of centers ($m \gg \ell$):

$$f(\mathbf{x}) = \sum_{m=1}^{\ell \times d} \sum_{p=1}^{\ell} c_{\ell}^p (x_{\ell} - \theta^{\ell} x)$$

Extensions of Additive Models

If we now allow for an arbitrary linear transformation of the inputs:

$$\mathbf{x} \mapsto W\mathbf{x}$$

where W is a $d' \times d$ matrix, we obtain:

$$f(\mathbf{x}) = \sum_{p=1}^{d'} \sum_{u=1}^n c_{pu} K_{pu}(\mathbf{x}^\top W \mathbf{w}^u - t_{pu})$$

where \mathbf{w}^u is the u -th row of the matrix W .

Extensions of Additive Models

The expression

$$\sum_{u=1}^{\infty} \sum_{p=1}^n K_{np}^{\alpha} (\mathbf{x} - \mathbf{w}_{\perp}^n)^{\alpha} = f(\mathbf{x})$$

can be written as

$$\sum_{p=1}^n (\mathbf{w}_{\perp}^n)^{\alpha} \eta^{\alpha} = f(\mathbf{x})$$

where

$$\sum_{u=1}^{\infty} K_{np}^{\alpha} (\eta - \hat{\eta})^{\alpha} = (\hat{\eta})^{\alpha} \eta^{\alpha}$$

This form of approximation is called **ridge approximation**

Gaussian MLP network

From the extension of additive models we can therefore justify an approximation technique of the form

$$f(\mathbf{x}) = \sum_{p=1}^n \sum_{m=1}^m c_{pm} G_{pm}(\mathbf{x} - \mathbf{t}_{pm})$$

Particular case: $m = 1$ (one center per dimension). Then we derive the following technique:

$$f(\mathbf{x}) = \sum_{p=1}^n c_p G_p(\mathbf{x} - \mathbf{t}_p)$$

which is a Multilayer Perceptron with a Radial Basis Functions G instead of the sigmoid function. One can argue rather formally that for normalized inputs the weight vectors of MLPs are equivalent to the centers of RBFs.

Notice that the sigmoid function cannot be derived – directly and formally – from regularization but...

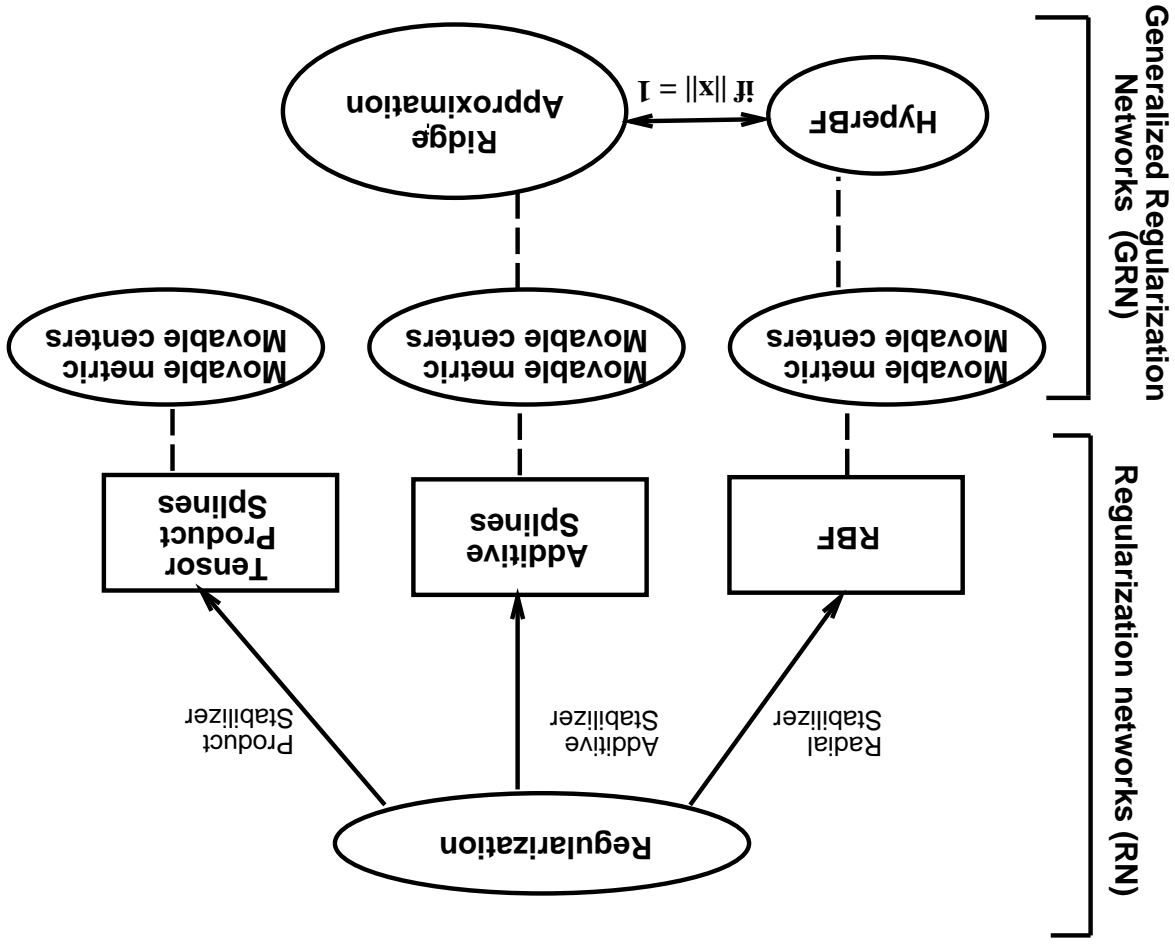
Sigmoids and Regularization

Suppose to have learned the representation

$$f(\mathbf{x}) = \sum_{p=1}^n \mathcal{C}_p K_p(\mathbf{x}) \mathbf{w}_p^\top$$

where $K_p(x) = |x|$. Notice that a finite linear combination of translates of a sigmoidal, piece-wise linear basis function can be written as a linear combination of translates of $|x|$. There is a very close relationship between 1-D radial and sigmoidal functions.

Regularization Networks



Regularization networks and Kernel regression

- Kernel regression: no complex global model of the world is assumed. Many simple local models instead (a case of *kernel methods*)

$$f(\mathbf{x}) = \frac{\sum_{j=1}^m w_j(\mathbf{x}) y_j}{\sum_{j=1}^m w_j(\mathbf{x})}$$

- Regularization networks: fairly complex global model of the world (a case of *dictionary methods*)

$$f(\mathbf{x}) = \sum_{j=1}^m c_j K(\mathbf{x} - \mathbf{x}_j)$$

Are these two techniques related? Can you say something about the apparent dichotomy of "local" vs. "global"?

Least square Regularization networks

A model of the form

$$f(\mathbf{x}) = \sum_m c_\alpha K(\mathbf{x} - \mathbf{t}_\alpha)$$

is assumed and the parameters c_α and \mathbf{t}_α are found by

$$\min_{c_\alpha, \mathbf{t}_\alpha} E[\{c_\alpha\}, \{\mathbf{t}_\alpha\}]$$

where

$$E[\{c_\alpha\}, \{\mathbf{t}_\alpha\}] = \sum_{i=1}^l (y_i - f(\mathbf{x}_i))^2$$

Least square Regularization networks

The coefficients c_α and the centers t_α have to satisfy the conditions:

$$\frac{\partial E}{\partial c_\alpha} = 0, \quad \frac{\partial E}{\partial t_\alpha} = 0 \quad \alpha = 1, \dots, m$$

The equation for the coefficients gives:

$$c_\alpha \sum_j H^{\alpha_j} y_j$$

where

$$H = (K^T K)^{-1} K^T \mathbf{t} = K^{-1} K^T \mathbf{t}$$

Dual representation

Substituting the expression for the coefficients in the regularization network we obtain

$$(\mathbf{x} - \mathbf{t}^\alpha) K_J^{\alpha} H \sum_m^{l=1} h_m \sum_j^{l=1} = f(\mathbf{x})$$

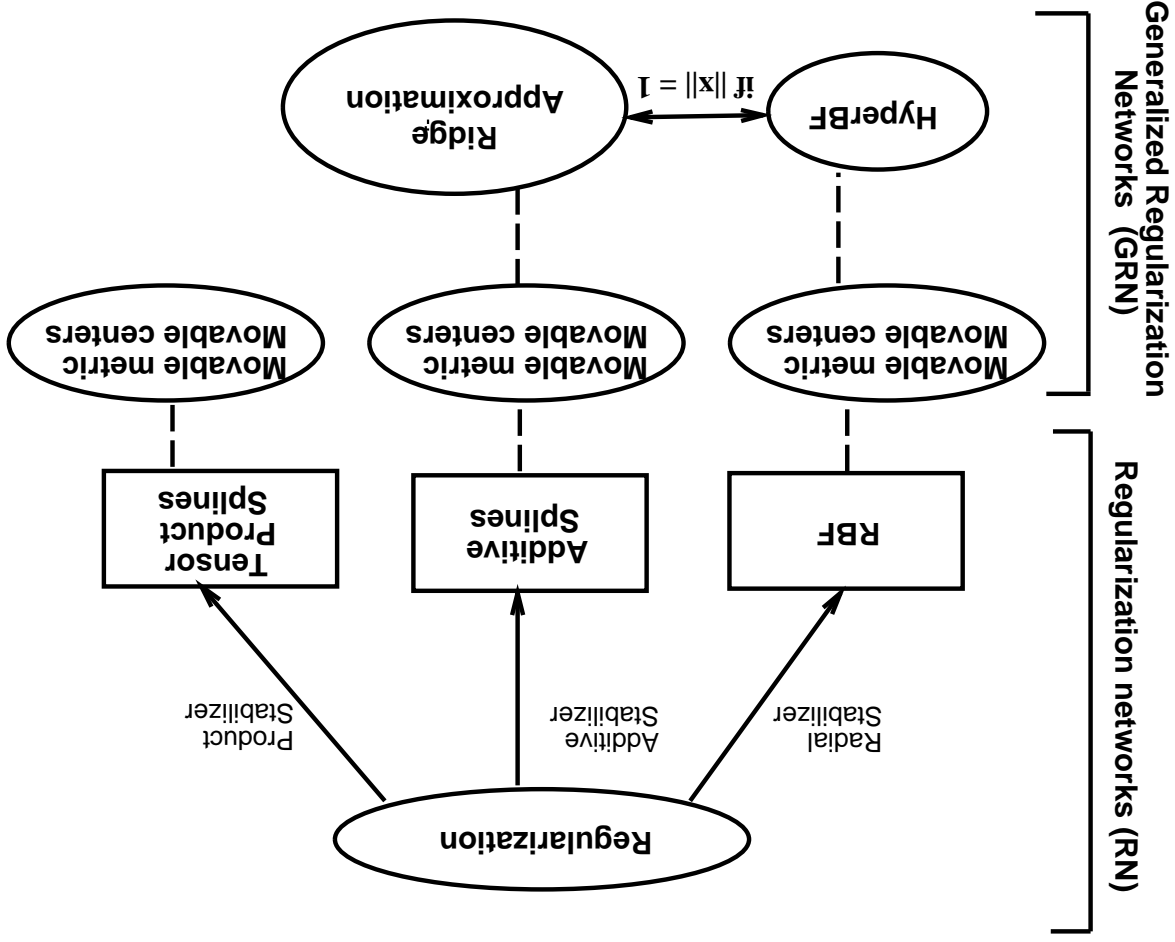
$$\sum_j^{l=1} h_j q_j(\mathbf{x}) = f(\mathbf{x})$$

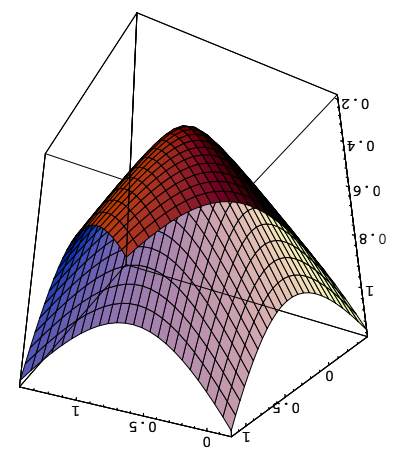
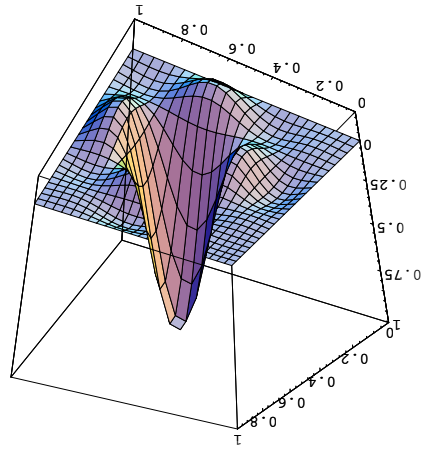
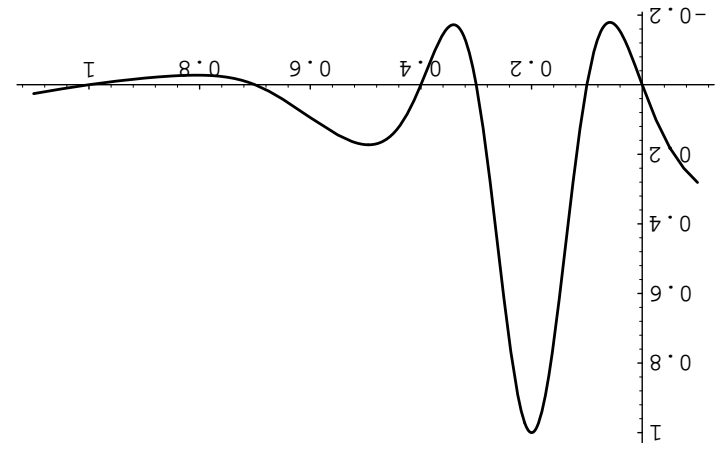
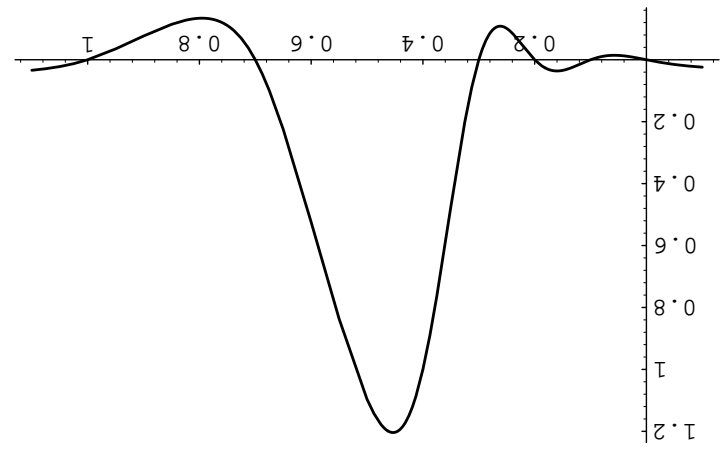
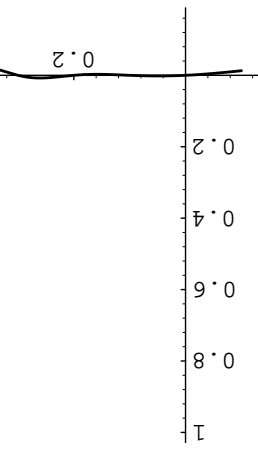
where we have defined

$$q_m(\mathbf{x}) = H \sum_m^{l=1} K_J^{\alpha}(\mathbf{x} - \mathbf{t}^\alpha)$$

The basis functions $q_m(\mathbf{x})$ are called "dual kernels".

Equivalent kernels for multiquadric basis functions





Dual formulation of Regularization networks and Kernel regression

$$\boxed{f(\mathbf{x}) = \sum_{j=1}^m y_j b_j(\mathbf{x})} \quad \text{Regularization networks}$$



$$\boxed{f(\mathbf{x}) = \frac{\sum_{j=1}^m y_j}{\sum_{j=1}^m w_j} w_j b_j(\mathbf{x})} \quad \text{Kernel regression}$$

In both cases the value of f at point \mathbf{x} is a weighted average of the values at the data points.

Project: is this true for SVMs? Can it be generalized?

Conclusions

- We have extended – with some hand waving – classical, quadratic Regularization Networks including RBF into a number of schemes that are inspired by regularization though do not strictly follow from it.
- The extensions described seem to work well in practice. Main problem – for schemes involving moving centers and or learning the metric – is efficient optimization.