# APPLICATION OF NON-PARAMETRIC REGRESSION IN ENGINEERING OPTIMIZATION

**Imre KOCSIS[1], Tamas MANKOVITS[2]**
[1]University of Debrecen Faculty of Engineering, kocsisi@eng.unideb.hu
[2] University of Debrecen Faculty of Engineering, tamas.mankovits@eng.unideb.hu

*Abstract*—Non-parametric regression models (e.g. kernel smoothing techniques, support vector regression model) have been widely used in statistics and recently in econometrics and engineering as well. Regression functions can be effective tools in the solution of engineering optimization problems. In the investigation of complex technical systems the class of the function describing the connection between the input and output data is generally unknown, moreover the "classical" least square fitting method is not flexible enough to provide an effective regression function in the case of higher dimensional optimization problems. In this paper regression functions obtained in different models are compared graphically in one dimensional case. Calculations were made using Maple, R, and MS Excel software.

*Keywords*—kernel function, regression

## I. INTRODUCTION

THE smoothing techniques and the non-parametric approach have a long tradition in empirical analysis. The intensive interest in smoothing over the previous decades had two reasons: "statisticians realized that pure parametric thinking in curve estimations often does not meet the need for flexibility in data analysis and the development of hardware created the demand for theory of now computable nonparametric estimates." [1]

Kernel smoothing is an effective regression technique for engineering optimization as well. It has the advantage that the function defining the non-linear relation between the input and output data does not need to be given in explicit form. As an example the shape optimization problem of rubber bumpers can be mentioned that was investigated by the second author in his doctoral thesis [7].

To present the flexibility of the kernel smoothing and the support vector regression we give the regression function for certain input data (one dimensional case).

## II. KERNEL FUNCTIONS

Let us consider a data set (frequently called training date set)

$$\{(\bar{t}_1, d_1), ..., (\bar{t}_P, d_P)\} \subset R^{N+1}, \quad (1)$$

where $\bar{t}_i$ are (possibly one dimensional) input vectors, $d_i$ are associated targets ($R$ is the set of reals). The goal is to fit a function $f : R^N \to R$ which approximates the relation between the data sets. Any regression algorithm has a loss function $L(f(\bar{t}), d)$ which describes how the estimated function deviated from the true one. In this note we refer to the following two types of loss function: the quadratic loss function (least square fitting)

$$\mathbf{L}(\mathbf{f}(\bar{\mathbf{t}}_i), \mathbf{d}_i) = (\mathbf{f}(\bar{\mathbf{t}}_i) - \mathbf{d}_i)^2 \quad (2)$$

and the $\varepsilon$-insensitive (Vapnik's) loss function

$$\mathbf{L}(\mathbf{f}(\bar{\mathbf{t}}_i), \mathbf{d}_i) = \begin{cases} \mathbf{0} & \text{if } |\mathbf{f}(\bar{\mathbf{t}}_i) - \mathbf{d}_i| \le \varepsilon \\ |\mathbf{f}(\bar{\mathbf{t}}_i) - \mathbf{d}_i| - \varepsilon & \text{otherwise} \end{cases}, \quad (3)$$

where $\varepsilon$ is a fixed positive parameter.

It is well-know that in the linear model (with the quadratic loss function)

$$\sum_{i=1}^{P} (\overline{\mathbf{w}} \cdot \bar{\mathbf{t}}_i + \mathbf{b} - \mathbf{d}_i)^2 \xrightarrow[\overline{\mathbf{w}}]{} \mathbf{min} \quad (4)$$

the regression function can be expressed as a linear combination of so-called kernel functions

$$\mathbf{k}_i(\bar{\mathbf{x}}) = \bar{\mathbf{x}} \cdot \bar{\mathbf{t}}_i \quad (5)$$

in the form

$$\mathbf{f}(\bar{\mathbf{x}}) = \overline{\mathbf{w}}^{\mathbf{opt}} \cdot \bar{\mathbf{x}} + \mathbf{b} = \sum_{i=1}^{P} \alpha_i \mathbf{k}_i(\bar{\mathbf{x}}) + \mathbf{b}, \quad (6)$$

where · denotes the inner product [5].

In non-linear models the situation is similar: applying a suitable transformation $\bar{\phi} : R^N \to R^M$ the problem will be linear in $R^M$ and the regression function is a linear combination of kernel functions

$$\mathbf{k}_i(\bar{\mathbf{x}}) = \bar{\phi}(\bar{\mathbf{x}}) \cdot \bar{\phi}(\bar{\mathbf{t}}_i) \quad (7)$$

[3], [4] in the form

$$\mathbf{f}(\bar{\mathbf{x}}) = \overline{\mathbf{w}}^{\mathbf{opt}} \cdot \bar{\phi}(\bar{\mathbf{x}}) + \mathbf{b} = \sum_{i=1}^{P} \alpha_i \mathbf{k}_i(\bar{\mathbf{x}}) + \mathbf{b}. \quad (8)$$

Since function $\bar{\phi}$ and, consequently, functions $k_i$ are generally unknown in practice, one of the recommended

kernel functions is chosen in the model to build $f$. The most used kernel in technical literature is the radial base Gaussian kernel function

$$\mathbf{k_i(\overline{x})} = \mathbf{e}^{-\gamma\|\overline{x}-\overline{t}_i\|^2}, \qquad (9)$$

where $\gamma$ is a parameter determining the shape of the kernel function. The role of $\gamma$ is presented graphically in section III.

### III. APPLICATION OF GAUSSIAN KERNEL FUNCTION

As an input data set consider the set of the following points $\{(t_i,d_i)/i=1,..,P\}=\{(1,16),(2,19),(3,9),(4,16),(5,25),(6,10),(7,5),(8,4),(9,24),(10,20)\}$ (Fig.5.) and solve the quadratic optimization problem

$$\mathbf{f(x)} = \sum_{i=1}^{P}\alpha_i\mathbf{k_i(x)} + \mathbf{b} = \sum_{i=1}^{P}\alpha_i \cdot \mathbf{e}^{-\gamma|x-t_i|^2} + \mathbf{b}$$

$$\sum_{i=1}^{P}\left(\mathbf{f(t_i)}-\mathbf{d_i}\right)^2 \xrightarrow[\overline{\alpha},b]{} \mathbf{min} \qquad (10)$$

for $\overline{\alpha} = (\alpha_1,...\alpha_P) \in R^P$ and $b \in R$, where $\gamma$ is a fixed positive number.

The following figures show graphically the regression functions obtained for certain values of $\gamma$. (Calculations and plotting were carried out using Maple).
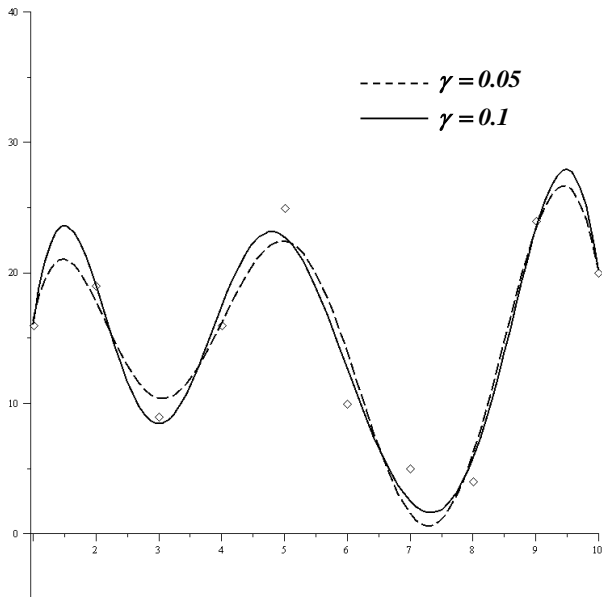


Fig. 1. The regression functions (solutions of the optimization problem (10)) when $\gamma=0.05$ and $\gamma=0.1$
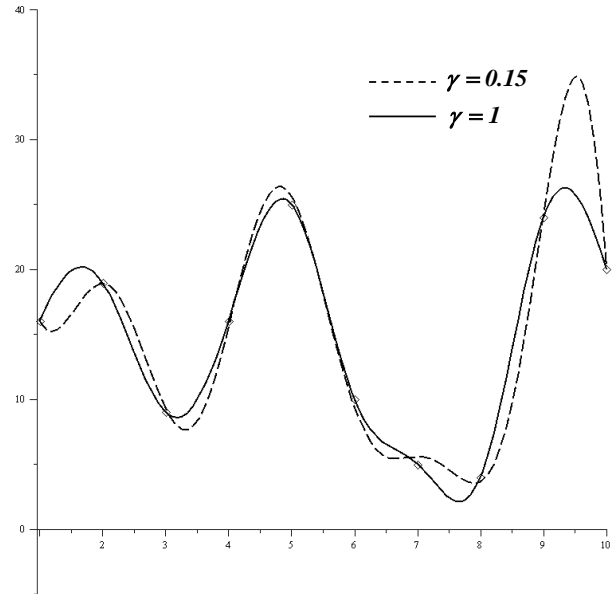


Fig. 2. The regression functions (solutions of the optimization problem (10)) when $\gamma=0.15$ and $\gamma=1$
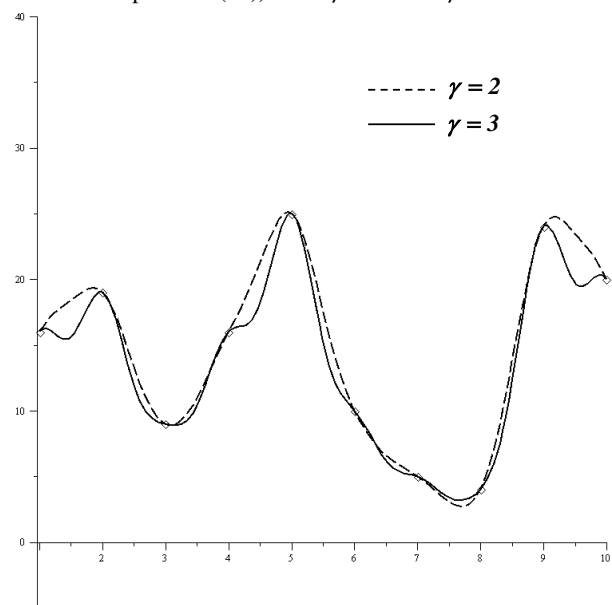


Fig. 3. The regression functions (solutions of the optimization problem (10)) when $\gamma=2$ and $\gamma=3$

It can be seen that the smaller values of $\gamma$ result flatter regression functions (see e.g. Fig.1.) but the deviation of these "flat" functions at the learning points can be relatively high. Large values of $\gamma$ lead to regression functions varying too fast and the shape of the regression functions do not seem to be reasonable (see e.g. Fig.4.). Analyzing the pictures it can be seen that the values $\gamma = 1..2$ mean an acceptable compromise between the flatness of the function and the accuracy at the learning points (Fig.2., Fig.3.).
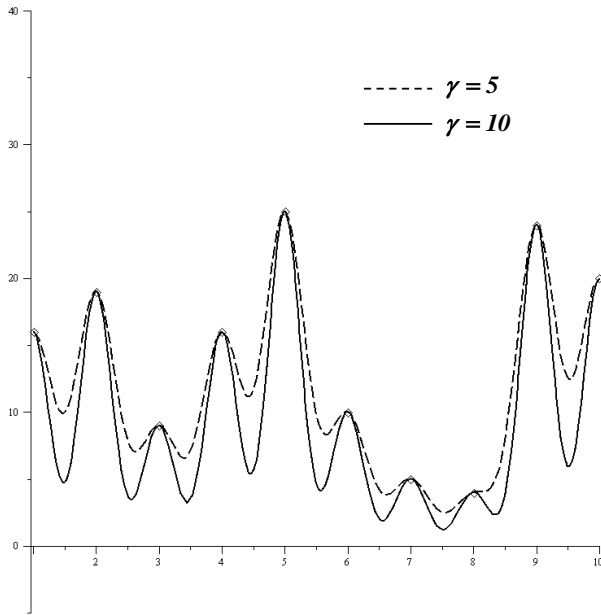
Fig. 4. The regression functions (solutions of the optimization problem (10)) when $\gamma=5$ and $\gamma=6$

## IV. SUPPORT VECTOR REGRESSION MODEL

Application of the $\varepsilon$-insensitive loss function (3) leads to a quadratic convex optimization problem and the regression function can be expressed by kernel functions.

In the linear model the regression function is sought in the form

$$\mathbf{f}(\overline{\mathbf{x}}) = \overline{\mathbf{w}} \cdot \overline{\mathbf{x}} + \mathbf{b}, \quad \overline{\mathbf{w}} \in \mathbf{R}^N, \mathbf{b} \in \mathbf{R}. \quad (11)$$

The flatness of $f$ in this case means that one seeks a small $\overline{w}$, that is, the norm $\|\overline{w}\|$ is to be minimized. The related convex optimization problem is the following:

$$\text{minimize} \quad \frac{1}{2}\|\overline{\mathbf{w}}\|^2$$

$$\text{subject to} \quad \begin{cases} \mathbf{d}_i - (\overline{\mathbf{w}} \cdot \overline{\mathbf{t}}_i + \mathbf{b}) \le \varepsilon \\ (\overline{\mathbf{w}} \cdot \overline{\mathbf{t}}_i + \mathbf{b}) - \mathbf{d}_i \le \varepsilon \end{cases}, \quad i = 1,...,\mathbf{P}.$$

$$(12)$$

Inequalities in (12) mean that the regression function is required to approximate the pairs $(\bar{t}_i, d_i)$ with $\varepsilon$ precision. Sometimes the convex optimization problem is infeasible with the given constraints and some "errors" must be allowed. For this purpose so-called slack variables was introduced by Vapnik in [2]. The modified convex optimization problem with slack variables $\xi_i$ and $\xi_i'$ is the following

$$\text{minimize} \quad \frac{1}{2}\|\overline{\mathbf{w}}\|^2 + \mathbf{C}\sum_{i=1}^{\mathbf{P}}(\xi_i + \xi_i')$$

$$\text{subject to} \quad \begin{cases} \mathbf{d}_i - (\overline{\mathbf{w}} \cdot \overline{\mathbf{t}}_i + \mathbf{b}) \le \varepsilon + \xi_i \\ (\overline{\mathbf{w}} \cdot \overline{\mathbf{t}}_i + \mathbf{b}) - \mathbf{d}_i \le \varepsilon + \xi_i', \quad i = 1,...,\mathbf{P}. \\ \xi_i \ge 0, \xi_i' \ge 0 \end{cases} \quad (13)$$

The positive constant $C$ determines the trade-off between the flatness of $f$ and the amount up to which

deviations larger than $\varepsilon$ are tolerated.

According to the standard dualization method [6] we introduce the Lagrange function

$$\mathbf{L}(\overline{\mathbf{w}}, \mathbf{b}, \overline{\lambda}_i, \overline{\lambda}_i', \overline{\alpha}_i, \overline{\alpha}_i') = \frac{1}{2}\|\overline{\mathbf{w}}\|^2 + \mathbf{C}\sum_{i=1}^{\mathbf{P}}(\xi_i + \xi_i')$$

$$- \sum_{i=1}^{\mathbf{P}}(\lambda_i\xi_i + \lambda_i'\xi_i') - \sum_{i=1}^{\mathbf{P}}\alpha_i(\varepsilon + \xi_i + (\overline{\mathbf{w}} \cdot \overline{\mathbf{x}} + \mathbf{b}) - \mathbf{d}_i) \quad (14)$$

$$- \sum_{i=1}^{\mathbf{P}}\alpha_i'(\varepsilon + \xi_i' + \mathbf{d}_i - (\overline{\mathbf{w}} \cdot \overline{\mathbf{x}} + \mathbf{b}))$$

where $\alpha_i, \alpha_i', \lambda_i, \lambda_i' \ge 0$. The dual optimization problem is the following:

$$\text{maximize} \quad -\frac{1}{2}\sum_{i,j=1}^{\mathbf{P}}(\alpha_i - \alpha_i')(\alpha_j - \alpha_j')\bar{\mathbf{t}}_i \cdot \bar{\mathbf{t}}_j$$

$$- \varepsilon\sum_{i=1}^{\mathbf{P}}(\alpha_i + \alpha_i') + \sum_{i=1}^{\mathbf{P}}(\alpha_i - \alpha_i')\bar{\mathbf{t}}_i \quad (15)$$

$$\text{subject to} \quad \sum_{i=1}^{\mathbf{P}}(\alpha_i - \alpha_i') = 0 \quad \text{and} \quad \alpha_i, \alpha_i' \in [0, \mathbf{C}].$$

Solving the dual problem we get

$$\mathbf{f}(\overline{\mathbf{x}}) = \sum_{i=1}^{\mathbf{P}}(\alpha_i - \alpha_i')\overline{\mathbf{x}} \cdot \bar{\mathbf{t}}_i + \mathbf{b}, \quad (16)$$

that is, the solution is a linear combination of kernel functions. This form of $f$ says that the explicit form of $\overline{w}$ does not need to be computed. Furthermore it can be proved that for the learning points inside the $\varepsilon$-tube $\alpha_i - \alpha_i' = 0$, that is, $f$ is determined by the learning points having nonvanishing coefficients. These pairs are called support vectors.

As it mentioned before, in non-linear model the application of a suitable transformation function $\overline{\phi}: R^N \to R^M$ leads to a linear problem. It can be proved that in this case regression function is a linear combination of kernel functions $k_i(\bar{x}) = \overline{\phi}(\bar{x}) \cdot \overline{\phi}(\bar{t}_i)$.

The convex optimization problem is the following

$$\text{minimize} \quad \frac{1}{2}\|\overline{\mathbf{w}}\|^2 + \mathbf{C}\sum_{i=1}^{\mathbf{P}}(\xi_i + \xi_i')$$

$$\text{subject to} \quad \begin{cases} \mathbf{d}_i - (\overline{\mathbf{w}} \cdot \phi(\bar{\mathbf{t}}_i) + \mathbf{b}) \le \varepsilon + \xi_i \\ (\overline{\mathbf{w}} \cdot \phi(\bar{\mathbf{t}}_i) + \mathbf{b}) - \mathbf{d}_i \le \varepsilon + \xi_i', i = 1,...,\mathbf{P}. \\ \xi_i \ge 0, \xi_i' \ge 0 \end{cases} \quad (17)$$

The dual problem is

$$\text{maximize} \quad -\frac{1}{2}\sum_{i,j=1}^{\mathbf{P}}(\alpha_i - \alpha_i')(\alpha_j - \alpha_j')\overline{\phi}(\bar{\mathbf{t}}_i) \cdot \overline{\phi}(\bar{\mathbf{t}}_j)$$

$$- \varepsilon\sum_{i=1}^{\mathbf{P}}(\alpha_i + \alpha_i') + \sum_{i=1}^{\mathbf{P}}(\alpha_i - \alpha_i')\bar{\mathbf{t}}_i \quad (18)$$

$$\text{subject to} \quad \sum_{i=1}^{\mathbf{P}}(\alpha_i - \alpha_i') = 0 \quad \text{and} \quad \alpha_i, \alpha_i' \in [0, \mathbf{C}].$$

Solving the dual problem we get

$$\mathbf{f}(\overline{\mathbf{x}}) = \sum_{i=1}^{\mathbf{P}}(\alpha_i - \alpha_i')\overline{\phi}(\overline{\mathbf{x}}) \cdot \overline{\phi}(\bar{\mathbf{t}}_i) + \mathbf{b}, \quad (19)$$

that is, the solution is a linear combination of kernel

functions.

Using the SVR package of the R software we present the role of parameters $\varepsilon$, $\gamma$ and $C$ in the SVR model. We consider the same input data set as in section II: *{($t_i,d_i$)/i=1,..,P}={(1,16),(2,19),(3,9),(4,16),(5,25),(6,10), (7,5),(8,4),(9,24),(10,20)}* (Fig.5.) and choose the Gaussian kernel function (9). (Plotting was carried out using MS Excel.)
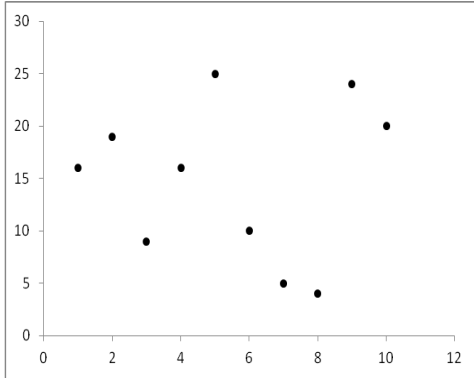


Fig. 5. The set of input data

Fig.6. shows the role of the "penalty" parameter $C$. Application of a higher value of $C$ results a regression function with small deviation even in outlying learning points, while low value of $C$ gives a flat function with large deviation at certain points.
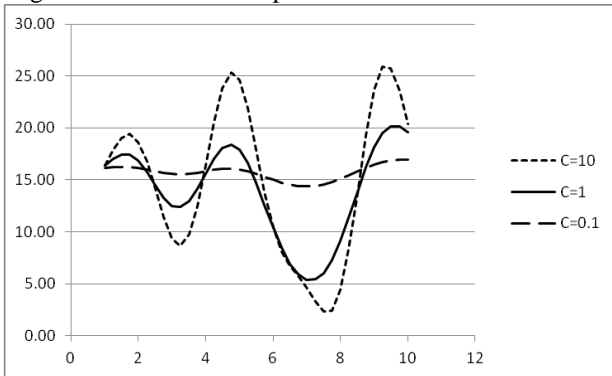


Fig. 6. The role of the "penalty" parameter $C$
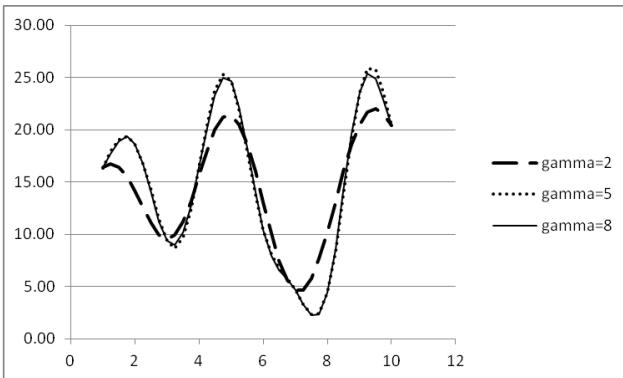($C$=0.1,1,10, $\gamma$=5, $\varepsilon$=0.05)



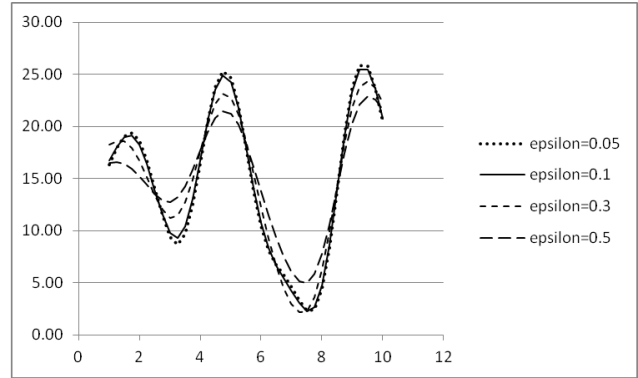Fig. 7. The role of the parameter $\gamma$ ($\gamma$=2,5,8, $C$=10, $\varepsilon$=0.05)



Fig. 8. The role of the parameter $\varepsilon$
($\varepsilon$=0.05,0.1,0.3,0.5, $\gamma$=5, $C$=10)

Fig.7. shows the role of the parameter $\gamma$ (it was described also in section II). Fig.8. shows the role of the parameter $\varepsilon$. The learning points are possibly (depending on the value of the slack variables) in an $\varepsilon$-tube or near the $\varepsilon$-tube.

Examples presented in this note show the flexibility of the regression methods using kernel functions. In engineering applications this flexibility can have an important role when the model has to be adjusted to the special characteristics of the input data set or to the requirements related to the regression function.

REFERENCES

[1] W. Hardle, "Applied non-parametric regression", *Humbolt-Universität*, Berlin, 1994
[2] C. Cortes and V. Vapnik, "Support vector networks", *Machine Learning*, Vol. 20, pp. 273-297, 1995.
[3] B. Schölkopf, and A. J. Smola, "Learning with Kernels", *MIT Press*, 2000
[4] A. Farag and R. M. Mohamed, "Regression Using Support Vector Machines: Basic Foundations", *University of Louisville*, 2004
[5] S. Haykin, "Neural Networks and Learning Machines", *Prentice Hall*, 2009
[6] S. Boyd and L. Vandenberghe, "Convex Optimization", *Cambridge University Press*, 2009
[7] T. Mankovits, "Shape optimization of rubber parts", PhD Thesis, *University of Miskolc*, 2013