

# Knowledge About the Margin of an Example

Marcin Orchel

## **Abstract**

In this report, we propose a method called  $\varphi$  support vector classification ( $\varphi$ -SVC) for incorporating knowledge about margin of an example for classification and regression problems. We propose two applications for  $\varphi$ -SVC: decreasing the generalization error of reduced models while preserving the similar number of support vectors, and incorporating the nonlinear constraint of a special type to the problem. The method was tested for support vector machines (SVM) classifier and  $\varepsilon$ -insensitive support vector regression ( $\varepsilon$ -SVR). Experiments on real world data sets show decreased generalization error of reduced models for linear and polynomial kernels.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Support Vector Classification Basics . . . . .	3
1.2	Support Vector Regression Basics . . . . .	6
1.3	SVM with $C_i$ Weights . . . . .	6
1.4	Introduction to $\delta$ -SVR . . . . .	8
1.4.1	Weighting the Translation Parameter . . . . .	10
<b>2</b>	<b>Knowledge About a Margin</b>	<b>11</b>
2.1	Introduction to $\varphi$ -SVC . . . . .	11
2.2	Knowledge About the Margin as Dynamic Hyperspheres . . . . .	14
2.3	Solving $\varphi$ -SVC Optimization Problem . . . . .	14
2.4	New Types of Support Vectors . . . . .	14
2.5	Reformulation of $\varepsilon$ -SVR as $\varphi$ -SVC . . . . .	15
2.6	Using Knowledge About a Margin with $\varepsilon$ -SVR . . . . .	17
2.7	Using Knowledge About a Margin with $\delta$ -SVR . . . . .	18
2.8	Changing the Output Curve . . . . .	18
2.9	Incorporating Linear Dependency of Function Values . . . . .	20
2.9.1	Incorporating Linear Dependency on Function Values to $\varphi$ -SVC . . . . .	21
2.9.2	Incorporating the Linear Dependency on Function Values to $\varepsilon$ -SVR . . . . .	22
2.9.3	Incorporating Linear Dependency on Function Values to $\delta$ -SVR . . . . .	23
2.10	Incorporating Inequalities with Function Values . . . . .	24
2.11	Reduce a Model with $\varphi$ -SVC . . . . .	27
2.12	Generation of Prior Knowledge . . . . .	27
2.13	Experiments . . . . .	29
2.13.1	First Experiment . . . . .	29
2.13.2	Second experiment . . . . .	31
2.14	Summary . . . . .	36
	<b>Appendix A</b>	<b>41</b>
A.1	Derivation of the Dual Form of OP 11 . . . . .	41
A.2	Derivation of $\varepsilon$ -SVR Reformulation as $\varphi$ -SVC . . . . .	43
A.3	Derivation of the Dual Form of OP 16 . . . . .	44
A.4	Incorporation of the Linear Dependency to $\varphi$ -SVC . . . . .	46

A.5 Incorporation of the Linear Dependency to $\delta$ -SVR . . . . .	47
<b>References</b>	<b>50</b>
<b>Notation and Symbols</b>	<b>52</b>
<b>Abbreviations</b>	<b>53</b>

# Chapter 1

## Introduction

### 1.1 Support Vector Classification Basics

For a classification problem, we consider a set of  $n$  training vectors  $\vec{x}_i$  for  $i \in \{1, \dots, n\}$ , where  $\vec{x}_i = (x_i^1, \dots, x_i^m)$ . The  $i$ -th training vector is mapped to  $y_c^i \in \{-1, 1\}$ . The  $m$  is a dimension of the problem. The support vector classification (SVC) optimization problem for hard margin case with  $\|\cdot\|_1$  norm is

**OP 1.**

$$\min_{\vec{w}_c, b_c} f(\vec{w}_c, b_c) = \|\vec{w}_c\|^2 \quad (1.1)$$

subject to

$$y_c^i h(\vec{x}_i) \geq 1 \quad (1.2)$$

for  $i \in \{1, \dots, n\}$ , where

$$h(\vec{x}_i) = \vec{w}_c \cdot \vec{x}_i + b_c \quad (1.3)$$

All points must be correctly classified Fig. 1.1a. The SVC soft margin case optimization problem with  $\|\cdot\|_1$  norm is

**OP 2.**

$$\min_{\vec{w}_c, b_c, \vec{\xi}_c} f(\vec{w}_c, b_c, \vec{\xi}_c) = \frac{1}{2} \|\vec{w}_c\|^2 + C_c \sum_{i=1}^n \xi_c^i \quad (1.4)$$

subject to

$$y_c^i h(\vec{x}_i) \geq 1 - \xi_c^i \quad (1.5)$$

$$\xi_c^i \geq 0 \quad (1.6)$$

for  $i \in \{1, \dots, n\}$ , where

$$h(\vec{x}_i) = \vec{w}_c \cdot \vec{x}_i + b_c \quad (1.7)$$

$$C_c > 0 \quad (1.8)$$

The  $h^*(\vec{x}) = \vec{w}_c^* \cdot \vec{x} + b_c^* = 0$  is a decision curve of the classification problem. Some of training points can be incorrectly classified Fig. 1.1b.

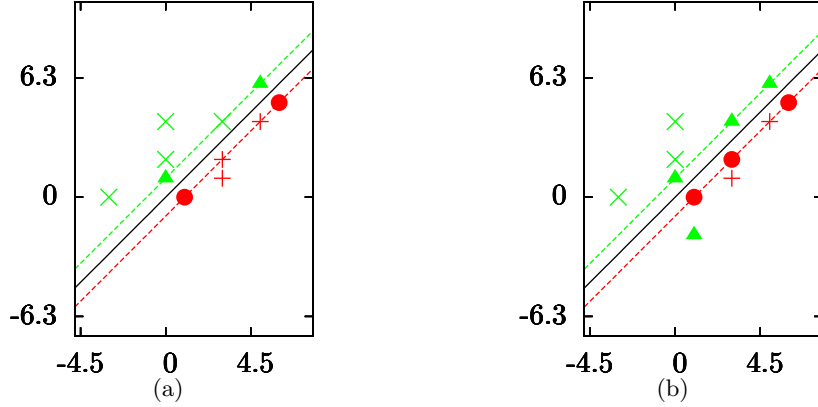


Figure 1.1: Two types of margin classifiers: hard and soft. Example points, support vectors (triangles and circles), solutions (solid lines), margin lines (dashed lines). (a) Hard. (b) Soft. A misclassified point is in (1, -2)

### SVC Dual Optimization Problem

The OP 2 optimization problem after transformation into an equivalent dual optimization problem becomes

**OP 3.**

$$\max_{\vec{\alpha}} f(\vec{\alpha}) = 1 \cdot \vec{\alpha} - \frac{1}{2} \vec{\alpha}^T \mathbf{Q} \vec{\alpha} \quad (1.9)$$

subject to

$$\vec{\alpha} \cdot \vec{y} = 0 \quad (1.10)$$

$$0 \leq \alpha_i \leq C_c \quad (1.11)$$

where

$$Q_{ij} = y_i y_j K(\vec{x}_i, \vec{x}_j) \quad (1.12)$$

for all  $i, j \in \{1, \dots, n\}$ .

The  $\vec{w}_c^* \cdot \vec{x}$  can be computed as

$$\vec{w}_c^* \cdot \vec{x} = \sum_{i=1}^n y_c^i \alpha_i^* K(\vec{x}_i, \vec{x}) . \quad (1.13)$$

Therefore, the decision curve is

$$h^*(\vec{x}) = \sum_{i=1}^n y_c^i \alpha_i^* K(\vec{x}_i, \vec{x}) + b_c^* = 0 , \quad (1.14)$$

where  $\alpha_i$  are Lagrange multipliers of the dual problem,  $K(\cdot, \cdot)$  is a kernel function, which appears only in the dual problem. *Margin boundaries* are defined as the two hyperplanes

$h(\vec{x}) = -1$  and  $h(\vec{x}) = 1$ . *Optimal margin boundaries* are defined as the two hyperplanes  $h^*(\vec{x}) = -1$  and  $h^*(\vec{x}) = 1$ . *Geometric margin of the hyperplane  $h$*  is defined as  $1/\|\vec{w}_c\|$ . The  $i$ -th training example is a *support vector*, when  $\alpha_i^* \neq 0$ . A set of support vectors contains all training examples lying below optimal margin boundaries ( $y_c^i h^*(\vec{x}_i) < 1$ ), and part of the examples lying exactly on the optimal margin boundaries ( $y_c^i h^*(\vec{x}_i) = 1$ ), Fig. 1.1b.

### SVC Without the Offset

Another variant of SVC is the SVC without the offset  $b_c$ , analyzed recently in [13]. The optimization problem is the same except missing  $b_c$  term, for the soft case it is

**OP 4.**

$$\min_{\vec{w}_c, \vec{\xi}_c} f(\vec{w}_c, \vec{\xi}_c) = \frac{1}{2} \|\vec{w}_c\|^2 + C_c \cdot \vec{\xi}_c \quad (1.15)$$

subject to

$$y_c^i h(\vec{x}_i) \geq 1 - \xi_c^i \quad (1.16)$$

$$\vec{\xi} \geq 0 \quad (1.17)$$

for  $i \in \{1, \dots, n\}$ , where

$$h(\vec{x}_i) = \vec{w}_c \cdot \vec{x}_i, \quad (1.18)$$

$$C_c > 0. \quad (1.19)$$

The dual problem is

**OP 5.**

$$\max_{\vec{\alpha}} d(\vec{\alpha}) = \vec{\alpha} - \frac{1}{2} \vec{\alpha}^T Q \vec{\alpha} \quad (1.20)$$

subject to

$$0 \leq \vec{\alpha} \leq C_c \quad (1.21)$$

where  $Q_{ij} = y_i y_j K(\vec{x}_i, \vec{x}_j)$ , for all  $i, j \in \{1, \dots, n\}$ .

We can notice missing linear constraint. The decision curve is

$$h^*(\vec{x}) = \sum_{i=1}^n y_c^i \alpha_i^* K(\vec{x}_i, \vec{x}) = 0. \quad (1.22)$$

## 1.2 Support Vector Regression Basics

In a regression problem, we consider a set of training vectors  $\vec{x}_i$  for  $i \in \{1, \dots, n\}$ , where  $\vec{x}_i = (x_i^1, \dots, x_i^m)$ . The  $i$ -th training vector is mapped to  $y_r^i \in \mathbb{R}$ . The  $m$  is a dimension of the problem. The  $\varepsilon$ -SVR soft case optimization problem is

**OP 6.**

$$\min_{\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*} f(\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*) = \frac{1}{2} \|\vec{w}_r\|^2 + C_r \sum_{i=1}^n (\xi_r^i + \xi_r^{*i}) \quad (1.23)$$

subject to

$$y_r^i - g(\vec{x}_i) \leq \varepsilon + \xi_r^i \quad (1.24)$$

$$g(\vec{x}_i) - y_r^i \leq \varepsilon + \xi_r^{*i} \quad (1.25)$$

$$\vec{\xi}_r \geq 0 \quad (1.26)$$

$$\vec{\xi}_r^* \geq 0 \quad (1.27)$$

for  $i \in \{1, \dots, n\}$ , where

$$g(\vec{x}_i) = \vec{w}_r \cdot \vec{x}_i + b_r, \quad (1.28)$$

$$\varepsilon \in \mathbb{R}. \quad (1.29)$$

The  $g^*(\vec{x}) = \vec{w}_r^* \cdot \vec{x} + b_r^*$  is a regression function. Optimization problem 6 is transformed into an equivalent dual problem. The regression function becomes

$$g^*(\vec{x}) = \sum_{i=1}^n (\alpha_i^* - \beta_i^*) K(\vec{x}_i, \vec{x}) + b_r^*, \quad (1.30)$$

where  $\alpha_i, \beta_i$  are Lagrange multipliers,  $K(\cdot, \cdot)$  is a kernel function. The  $\varepsilon$  boundaries are defined as  $g(\vec{x}) - \varepsilon$  and  $g(\vec{x}) + \varepsilon$ . The  $i$ -th training example is a *support vector*, when  $\alpha_i^* - \beta_i^* \neq 0$ . For  $\varepsilon \geq 0$ , a set of support vectors contains all training examples lying outside  $\varepsilon$  boundaries, and part of the examples lying exactly on  $\varepsilon$  boundaries, Fig. 1.2. The number of support vectors can be controlled by  $\varepsilon$  parameter.

## 1.3 SVM with $C_i$ Weights

A 1-norm soft margin SVC optimization problem for training examples  $\vec{x}_i$  with weights  $C_i$  is

**OP 7.**

$$\min_{\vec{w}, b, \vec{\xi}} f(\vec{w}, b, \vec{\xi}) = \frac{1}{2} \|\vec{w}\|^2 + \vec{C}_c \cdot \vec{\xi} \quad (1.31)$$

subject to

$$y_i h(\vec{x}_i) \geq 1 - \xi_i \quad (1.32)$$

$$\vec{\xi} \geq 0 \quad (1.33)$$



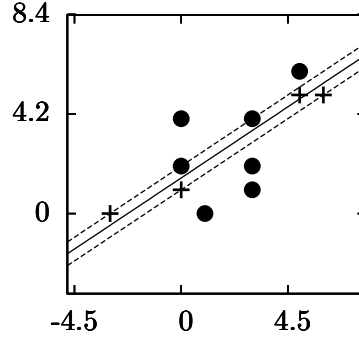


Figure 1.2: The idea of  $\varepsilon$ -SVR. Points - examples, circles - support vectors, a solid line - a solution, dashed lines -  $\varepsilon$  boundaries

for  $i \in \{1, \dots, n\}$ , where

$$\vec{C}_c \gg 0 \quad (1.34)$$

$$h(\vec{x}_i) = \vec{w} \cdot \vec{x}_i + b \quad (1.35)$$

The  $C_i$  weights were also used with  $\varepsilon$ -SVR for predicting time series data, [14]. The other types of weights that were used with  $\varepsilon$ -SVR are  $\varepsilon_i$  weights per example replacing the parameter  $\varepsilon$ . They were used for density estimation, [15]. Sometimes we use different  $\varepsilon$  weights for inequalities (1.24), (1.25),  $\varepsilon_u$  and  $\varepsilon_d$  respectively. And finally we can combine above changes and use  $\varepsilon_u^i$  and  $\varepsilon_d^i$  weights

**OP 8.**

$$\min_{\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*} f(\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*) = \frac{1}{2} \|\vec{w}_r\|^2 + \vec{C}_r \sum_{i=1}^n (\xi_r^i + \xi_r^{*i}) \quad (1.36)$$

subject to

$$y_r^i - g(\vec{x}_i) \leq \varepsilon_u^i + \xi_r^i \quad (1.37)$$

$$g(\vec{x}_i) - y_r^i \leq \varepsilon_d^i + \xi_r^{*i} \quad (1.38)$$

$$\vec{\xi}_r \geq 0 \quad (1.39)$$

$$\vec{\xi}_r^* \geq 0 \quad (1.40)$$

for  $i \in \{1, \dots, n\}$ , where

$$g(\vec{x}_i) = \vec{w}_r \cdot \vec{x}_i + b_r \quad (1.41)$$

We can notice that changing  $y_r^i$  value by  $\Delta y_r^i$  is equivalent to changing  $\varepsilon_u^i$  by  $-\Delta y_r^i$  and changing  $\varepsilon_d^i$  by  $\Delta y_r^i$ .

## 1.4 Introduction to $\delta$ -SVR

We consider a set of training vectors  $\vec{x}_i$  for  $i \in \{1, \dots, n\}$ , where  $\vec{x}_i = (x_i^1, \dots, x_i^m)$ . The  $i$ -th training vector is mapped to  $y_r^i \in \mathbb{R}$ . The  $\delta$  support vector regression ( $\delta$ -SVR) method finds a regression function by the following procedure.

1. Every training example  $\vec{x}_i$  is duplicated, an output value  $y_r^i$  is increased by a value of a parameter  $\delta \geq 0$  for original training examples, and decreased by  $\delta$  for duplicated training examples.
2. Every training example  $\vec{x}_i$  is converted to a classification example by incorporating the output to the input vector as an additional feature and setting class 1 for original training examples, class  $-1$  for duplicated training examples.
3. The SVC method is launched for a classification problem.
4. The solution of SVC method is converted back to function form.

The result of the first step is a set of training mappings for  $i \in \{1, \dots, 2n\}$

$$\begin{cases} \vec{x}_i \rightarrow y_r^i + \delta & \text{for } i \in \{1, \dots, n\} \\ \vec{x}_i \rightarrow y_r^i - \delta & \text{for } i \in \{n+1, \dots, 2n\} \end{cases} \quad (1.42)$$

where  $x_{n+i}^1 = \vec{x}_i$ ,  $y_r^{n+i} = y_r^i$  for  $i \in \{1, \dots, n\}$ ,  $\delta \in \mathbb{R}$ . The  $\delta$  is called *the translation parameter*. The result of the second step is a set of training mappings for  $i \in \{1, \dots, 2n\}$

$$\vec{c}_i = (x_i^1, \dots, x_i^m, y_r^i + y_c^i \delta) \rightarrow y_c^i \quad (1.43)$$

where  $y_c^i = 1$  for  $i \in \{1, \dots, n\}$ , and  $y_c^i = -1$  for  $i \in \{n+1, \dots, 2n\}$ . The dimension of the  $\vec{c}_i$  vectors is equal to  $m+1$ . The set of  $\vec{x}_i$  mappings before duplication is called *a regression data setting*, the set of  $\vec{c}_i$  ones is called *a classification data setting*. In the third step, OP 2 is solved with  $\vec{c}_i$  examples, so we can write OP 2 as

**OP 9.**

$$\min_{\vec{w}_c, b_c, \xi_c} f(\vec{w}_c, b_c, \xi_c) = \frac{1}{2} \|\vec{w}_c\|^2 + C_c \sum_{i=1}^{2n} \xi_c^i \quad (1.44)$$

subject to

$$y_c^i (w_{c,\text{red}} \cdot \vec{x}_i + w_c^{m+1} (y_r^i + y_c^i \delta) + b_c) \geq 1 - \xi_c^i \quad (1.45)$$

$$\xi_c^i \geq 0 \quad (1.46)$$

for  $i \in \{1, \dots, 2n\}$ .

The  $w_{c,\text{red}}$  is defined as  $w_{c,\text{red}} = (w_1, \dots, w_m)$ . The  $h^*(\vec{x}) = \vec{w}_c^* \cdot \vec{x} + b_c^* = 0$  is a decision curve of the classification problem. Note that  $h^*(\vec{x})$  is in the implicit form of the last coordinate of  $\vec{x}$ . In the fourth step, an explicit form of the last coordinate needs

to be find. The explicit form is needed for example for testing new examples. The  $\vec{w}_c$  variable of the primal problem for a simple linear kernel is found in the following way

$$\vec{w}_c = \sum_{i=1}^{2n} y_c^i \alpha_i \vec{c}_i . \quad (1.47)$$

For a simple linear kernel the explicit form of (1.14) is

$$x_{m+1} = \frac{-\sum_{j=1}^m w_c^j x_j - b_c}{w_c^{m+1}} . \quad (1.48)$$

The regression solution is  $g^*(\vec{x}) = \vec{w}_r \cdot \vec{x} + b_r$ , where  $w_r^i = -w_c^i/w_c^{m+1}$ ,  $b_r = -b_c/w_c^{m+1}$  for  $i = 1, \dots, m$ . For nonlinear kernels, a conversion to the explicit form has some limitations. First, a decision curve could have more than one value of the last coordinate for specific values of remaining coordinates of  $\vec{x}$  and therefore it cannot be converted unambiguously to the function (e.g. a polynomial kernel with a dimension equals to 2). Second, even when the conversion to the function is possible, there is no explicit analytical formula (e.g. a polynomial kernel with a dimension greater than 4), hence a special method for finding the explicit formula of the coordinate should be used, e.g. a bisection method. The disadvantage of this solution is a longer time of testing new examples. To overcome these problems, we proposed to incorporate prior knowledge to the classification problem, that the solution will be always in the form of the function in the chosen direction. Thus, we proposed in [10] a new kernel type in which the last coordinate is placed only inside a linear term. The new kernel is constructed from an original kernel by removing the last coordinate, and adding the linear term with the last coordinate

$$K(\vec{x}, \vec{y}) = K_o(x_{\text{red}}, y_{\text{red}}) + x_{m+1}y_{m+1} , \quad (1.49)$$

where  $\vec{x}$  and  $\vec{y}$  are  $m+1$  dimensional vectors,  $x_{\text{red}} = (x_1, \dots, x_m)$ ,  $y_{\text{red}} = (y_1, \dots, y_m)$ ,  $K_o(\cdot, \cdot)$  is the original kernel from which the new one was constructed. For the most popular kernels polynomial, radial basis function (RBF) and sigmoid, the conversions are respectively

$$(\vec{x} \cdot \vec{y})^d \rightarrow \left( \sum_{i=1}^m x_i y_i \right)^d + x_{m+1} y_{m+1} , \quad (1.50)$$

$$\exp - \frac{\|\vec{x} - \vec{y}\|^2}{2\sigma^2} \rightarrow \exp - \frac{\sum_{i=1}^m (x_i - y_i)^2}{2\sigma^2} + x_{m+1} y_{m+1} , \quad (1.51)$$

$$\tanh \vec{x} \vec{y} \rightarrow \tanh \sum_{i=1}^m x_i y_i + x_{m+1} y_{m+1} , \quad (1.52)$$

where  $\vec{x}$  and  $\vec{y}$  are  $m+1$  dimensional vectors. For the new kernel type, the explicit form of (1.14) for  $\delta$ -SVR is

$$x_{m+1} = \frac{-\sum_{i=1}^{2n} y_c^i \alpha_i K_o(\vec{x}_i, x_{\text{red}}) - b_c}{\sum_{i=1}^{2n} y_c^i \alpha_i c_i^{m+1}} . \quad (1.53)$$

### 1.4.1 Weighting the Translation Parameter

We can consider incorporating prior knowledge by setting different values of the translation parameter for each example, so we can have  $\delta_i$  parameters, for  $i \in \{1, \dots, n\}$  and the same parameters for  $i \in \{n+1, \dots, 2n\}$ . We can also consider setting different values of  $\delta$  for up and down translations, so we can have two parameters:  $\delta_u$  and  $\delta_d$ . And finally we can also consider the parameters  $\delta_u^i$  and  $\delta_d^i$  (additionally with the  $\vec{C}_c$  weight)

**OP 10.**

$$\min_{\vec{w}_c, b_c, \vec{\xi}_c} f(\vec{w}_c, b_c, \vec{\xi}_c) = \frac{1}{2} \|\vec{w}_c\|^2 + \vec{C}_c \sum_{i=1}^{2n} \xi_c^i \quad (1.54)$$

subject to

$$y_c^i (w_{c,\text{red}} \cdot \vec{x}_i + w_c^{m+1} (y_r^i + y_c^i \delta_i) + b_c) \geq 1 - \xi_c^i \quad (1.55)$$

$$\vec{\xi}_c \geq 0 \quad (1.56)$$

for  $i \in \{1, \dots, 2n\}$ , where  $\delta_i = \delta_u^i$  for  $i \in \{1, \dots, n\}$  and  $\delta_i = \delta_d^i$  for  $i \in \{n+1, \dots, 2n\}$ .

## Chapter 2

# Knowledge About a Margin

We define the margin of an example (sometimes called just a margin) in the following way:

**Definition 2.0.1.** Given some curve  $h(\vec{x}) = 0$ , the margin of the  $\vec{x}_p$  example is defined as a value  $|h(\vec{x}_p)|$ .

For hard margin SVC, OP 1, the margin of the closest examples is equal to 1. *The knowledge about the margin of an example* (sometimes called the knowledge about a margin) is defined as prior information about the margins of particular examples.

### 2.1 Introduction to $\varphi$ -SVC

The  $\varphi$ -SVC method is a recently proposed method of incorporating knowledge about the margin of an example to SVC, [8, 9, 11]. The  $\varphi$ -SVC optimization problem is defined with an additional parameter per example added in the right side of the inequality (1.5). Another modification of inequality constraints was proposed in [17]. The authors modify the inequalities by multiplying the left side of the inequalities by some monotonically decreasing function of additional example weights. The  $\varphi$ -SVC is a more general concept of weights per example with any values possible and with different interpretation.

Now, we will closely look at  $\varphi$ -SVC optimization problem. We define  $\varphi$ -SVC optimization problem based on SVC with cost weights per example OP 7 as

**OP 11.**

$$\min_{\vec{w}_c, b_c, \vec{\xi}_c} f(\vec{w}_c, b_c, \vec{\xi}_c) = \frac{1}{2} \|\vec{w}_c\|^2 + \vec{C}_c \cdot \vec{\xi}_c \quad (2.1)$$

subject to

$$y_c^i h(\vec{x}_i) \geq 1 + \varphi_i - \xi_c^i \quad (2.2)$$

$$\vec{\xi}_c \geq 0 \quad (2.3)$$

for  $i \in \{1, \dots, n\}$ , where

$$\vec{C}_c \gg 0 \quad (2.4)$$

$$\varphi_i \in \mathbb{R} \quad (2.5)$$

$$h(\vec{x}_i) = \vec{w}_c \cdot \vec{x}_i + b_c. \quad (2.6)$$

The new weights  $\varphi_i$  are present only in (2.2). When  $\vec{\varphi} = 0$ , the OP 11 is equivalent to OP 7. When all  $\varphi_i$  are equal to some constant  $\varphi > -1$ , we will get the same decision boundary as for  $\varphi = 0$  when we change  $\vec{C}_c$  to  $\vec{C}_c / (1 + \varphi)$ .

*Proof.* We will prove that we get the same decision boundary when we replace 1 with  $1/d$ , for some  $d > 0$ ,  $\varphi_i = 0$ . Let's replace  $\xi_i$  with  $\xi_i/d$  and we get the inequalities  $y_c^i h(\vec{x}_i) \geq 1/d - \xi_c^i/d$ . After multiplying by  $d$  we get  $y_c^i d h(\vec{x}_i) \geq 1 - \xi_c^i$ . The objective function can be multiplied by  $d^2$  and we get  $\frac{1}{2} \|\vec{d} \vec{w}_c\|^2 + d \vec{C}_c \cdot \vec{\xi}_c$ , the inequalities  $\xi_i/d \geq 0$  can be replaced by  $\xi_i \geq 0$ . So we get the same optimization problem as the original one with the new  $\vec{C}_c = d \vec{C}_c$  and with the new decision curve  $dh(\vec{x}) = 0$ .  $\square$

We can notice that when neglecting  $\xi_c^i$ , we get different bounds for the margin: when  $y_c^i = 1$  and  $g(\vec{x}_i) \geq 0$ , then we get a lower bound  $1 + \varphi_i$ , when  $y_c^i = -1$  and  $g(\vec{x}_i) \geq 0$ , then we get an upper bound  $-(1 + \varphi_i)$ , when  $y_c^i = 1$  and  $g(\vec{x}_i) < 0$ , then we get an upper bound  $-(1 + \varphi_i)$ , when  $y_c^i = -1$  and  $g(\vec{x}_i) < 0$ , then we get a lower bound  $1 + \varphi_i$ . We can distinguish three cases:  $1 + \varphi_i > 0$ ,  $1 + \varphi_i < 0$  and  $1 + \varphi_i = 0$ . For the first one, we get a lower bound on the margin equal to  $1 + \varphi_i$ . For the second, we get an upper bound on the margin equal to  $-(1 + \varphi_i)$ , for the third, we get an upper bound on the margin equal to 0. Therefore, by using  $\varphi_i$  weights, we can incorporate knowledge about the margin of an example.

The next property of  $\varphi_i$  weights is that they have impact on the distance between the curve  $h(\vec{x}) = 0$  and the  $i$ -th example. We can conduct the same analysis as above for distances by dividing both sides of (2.2) by  $\|\vec{w}_c\|$ , so we get  $y_c^i h(\vec{x}_i) / \|\vec{w}_c\| \geq 1 / \|\vec{w}_c\| + \varphi_i / \|\vec{w}_c\|$ . The conclusion is similar: for the case when  $1 + \varphi_i > 0$ , we get a lower bound on the distance equal to  $(1 + \varphi_i) / \|\vec{w}_c\|$ . For the case when  $1 + \varphi_i < 0$ , we get an upper bound on the distance equal to  $-(1 + \varphi_i) / \|\vec{w}_c\|$ . For the case when  $1 + \varphi_i = 0$ , we get an upper bound on the distance equal to 0.

Note that when we take into account  $\xi_c^i$ , we can see that knowledge about the margin of an example is incorporated as imperfect prior knowledge. The violation of knowledge about the margin of an example is controlled by the  $C_c^i$  parameters. Comparing loosely  $\varphi_i$  weights with slack variables:  $\varphi_i$  weights are constant, they are absent in the objective function, whereas a sum of slack variables is minimized as part of the objective function.

We can also derive the equivalent optimization problem to OP 11, where  $\varphi_i$  weights are present in the constraints with slack variables

**OP 12.**

$$\min_{\vec{w}_c, b_c, \vec{\xi}_c} f(\vec{w}_c, b_c, \vec{\xi}_c) = \frac{1}{2} \|\vec{w}_c\|^2 + \vec{C}_c \cdot \vec{\xi}_c \quad (2.7)$$

subject to

$$y_c^i h(\vec{x}_i) \geq 1 - \xi_c^i \quad (2.8)$$

$$\vec{\xi} \geq \varphi_i \quad (2.9)$$

for  $i \in \{1, \dots, n\}$ .

In order to construct an efficient algorithm for the OP 11 its dual form was derived (derivation in A.1). The final form of the dual problem is

**OP 13.**

$$\max_{\vec{\alpha}} \quad d(\vec{\alpha}) = \vec{\alpha} \cdot (1 + \vec{\varphi}) - \frac{1}{2} \vec{\alpha}^T Q \vec{\alpha} \quad (2.10)$$

subject to

$$\vec{\alpha} \cdot \vec{y}_c = 0 \quad (2.11)$$

$$0 \leq \vec{\alpha} \leq \vec{C} \quad , \quad (2.12)$$

where

$$Q_{ij} = y_c^i y_c^j (\vec{x}_i \cdot \vec{x}_j) \quad (2.13)$$

for all  $i, j \in \{1, \dots, n\}$ .

It differs from the original SVC dual form OP 3 by only  $\vec{\alpha} \cdot \vec{\varphi}$  term (also here we use  $C_c^i$  weights instead of  $C_c$ ). In the above formulation, similarly as for the original SVC, it is possible to introduce nonlinear decision functions by using a kernel function instead of a scalar product. The final decision boundary has a form

$$h^*(\vec{x}) = \sum_{i=1}^n y_c^i \alpha_i^* K(\vec{x}_i, \vec{x}) + b^* = 0 \quad , \quad (2.14)$$

where  $K(\cdot, \cdot)$  is a kernel function. The Karush-Kuhn-Tucker (KKT) complementary condition is

$$\alpha_i (y_c^i h(\vec{x}_i) - 1 - \varphi_i + \xi_c^i) = 0 \quad (2.15)$$

$$(C_i - \alpha_i) \xi_c^i = 0 \quad . \quad (2.16)$$

The conclusions from (2.15) and (2.16) are: when  $\alpha_i = 0$ , then  $\xi_c^i = 0$ , when  $0 < \alpha_i < C_c$ , then  $\xi_c^i = 0$  and  $y_c^i h(\vec{x}_i) = 1 + \varphi_i$ , when  $\alpha_i = C_c$ , then  $y_c^i h(\vec{x}_i) = 1 + \varphi_i - \xi_c^i$ . Moreover, when  $\xi_c^i > 0$ , then  $\alpha_i = C_c$  and  $y_c^i h(\vec{x}_i) = 1 + \varphi_i - \xi_c^i$ , when  $y_c^i h(\vec{x}_i) > 1 + \varphi_i - \xi_c^i$ , then  $\alpha_i = 0$ ,  $\xi_c^i = 0$  and  $y_c^i h(\vec{x}_i) > 1 + \varphi_i$ . We can find  $\xi_c^i$  parameters from the solution of the dual form as following. When

$$y_c^i h^*(\vec{x}_i) \geq 1 + \varphi_i \quad , \quad (2.17)$$

then  $\xi_c^i = 0$ , else

$$\xi_c^i = 1 + \varphi_i - y_c^i h(\vec{x}_i) \quad . \quad (2.18)$$

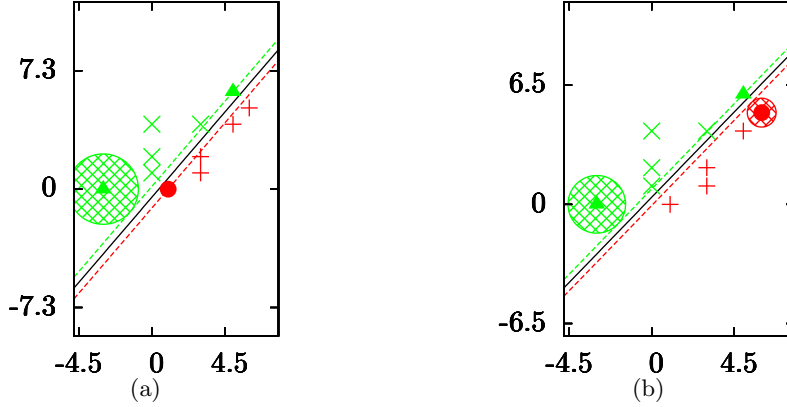


Figure 2.1: Interpretation of knowledge about a margin as dynamic hyperspheres. Points - examples, solid lines - solutions, triangles and circles - support vectors, circles filled with grid pattern - dynamic hyperspheres. The example in  $(-3, 0)$  has  $\varphi_i = 5.0$  and it is a marginal support vector. The hyperspheres must not cross the appropriate margin boundaries. (a)  $C = 100$ . (b)  $C = 10$  and  $\varphi_i = 2.5$  in  $(6, 5)$ . A radius of a dynamic hypersphere in  $(-3, 0)$  differs in both cases, in (a) it is 2.2, and in (b) it is 1.6

## 2.2 Knowledge About the Margin as Dynamic Hyperspheres

Knowledge about the margin of the  $p$ -th example can be visualized as a hypersphere with the center in the  $p$  point and a radius equal to  $\varphi_p / \|\vec{w}_c\|$ . We call it a *dynamic hypersphere*, because a radius depends on  $\|\vec{w}_c\|$ . For the two solution candidates  $h_1(\vec{x}) = 0$  and  $h_2(\vec{x}) = 0$ , where  $h_2(\vec{x}) = ah_1(\vec{x})$  and  $a \neq 0$  (both hyperplanes have the same geometric location), the hyperspheres are respectively  $S_1(\vec{x}_p, r)$ , and  $S_2(\vec{x}_p, r/a)$ , Fig. 2.1. Knowledge about the margin of an example in the form of a lower bound on the margin can be interpreted as a constraint that the hypersphere must not cross the appropriate margin boundary, Fig. 2.1. Knowledge about the margin of an example in the form of an upper bound on the margin can be interpreted as a constraint that the hypersphere must cross the appropriate margin boundary, Fig. 2.2.

## 2.3 Solving $\varphi$ -SVC Optimization Problem

For solving OP 13, a decomposition method similar to sequential minimal optimization (SMO), [12] was derived. In every step, two parameter subproblems are solved. Heuristic and stopping criterion are based on KKT conditions.

## 2.4 New Types of Support Vectors

The  $i$ -th example is a support vector, when  $\alpha_i^* \neq 0$ . From KKT (2.15), (2.16), we can conclude which examples could be support vectors. In the original SVC, only the example that lies on the optimal margin boundaries ( $y_i h^*(\vec{x}_i) = 1$ ) or below optimal



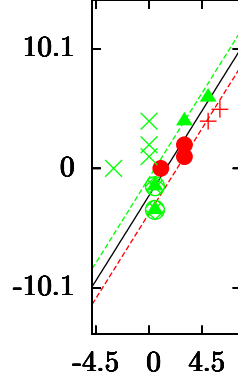


Figure 2.2: Interpretation of knowledge about a margin as dynamic hyperspheres. Points - examples, a solid line - a solution, triangles and circles - support vectors, circles filled with grid pattern - dynamic hyperspheres. The example in  $(0.5, -1.5)$  has  $\varphi_i = -1.0$ . Its hypersphere must cross the appropriate margin boundary

margin boundaries ( $y_i h^*(\vec{x}_i) < 1$ ) could be a support vector. In  $\varphi$ -SVC, also the example satisfying  $\varphi_i > 0$  and lying above margin boundaries ( $y_i h^*(\vec{x}_i) > 1$ ) could be a support vector. Such example is called *marginal support vector*, Fig. 2.1. Because an output model is formulated based on support vectors, introducing the new type of support vectors leads to richer models with additional examples lying above optimal margin boundaries.

Examples with upper bounded margins could be the examples that lie below the margin boundary and are not support vectors. Moreover, such examples could lead to the new type of support vectors – that have slack variables equal to zero and lie below the margin boundaries.

The better flexibility of  $\varphi$ -SVC in choosing support vectors suggests that we can achieve solutions with fewer support vectors for the similar generalization performance.

## 2.5 Reformulation of $\varepsilon$ -SVR as $\varphi$ -SVC

We have found that the  $\varepsilon$ -SVR is a special case of  $\varphi$ -SVC, [11]. The similar reformulation was previously implemented in LibSVM, [1] for solving ordinal classification problems – without prior knowledge, and  $\varepsilon$ -SVR. The consequence of this reformulation is that we can apply all the applications for  $\varphi$ -SVC also for  $\varepsilon$ -SVR.

Here we present a reformulation of the OP 6 (derivation in A.2). Every regression training example is duplicated, Fig. 2.3. Every original training example gets 1 class, and the duplicated training example gets -1 class and therefore we get

**OP 14.**

$$\min_{\vec{w}_c, b_c, \vec{\xi}_c} f(\vec{w}_c, b_c, \vec{\xi}_c) = \frac{1}{2} \|\vec{w}_c\|^2 + C_c \sum_{i=1}^{2n} \xi_c^i \quad (2.19)$$

subject to

$$y_c^i h(\vec{x}_i) \geq 1 - \xi_c^i + \varphi_i \quad (2.20)$$

$$\vec{\xi}_c \geq 0 \quad (2.21)$$

for  $i \in \{1, \dots, 2n\}$ , where

$$h(\vec{x}_i) = \vec{w}_c \cdot \vec{x}_i + b_c \quad (2.22)$$

and

$$\varphi_i = y_c^i y_r^i - \varepsilon - 1. \quad (2.23)$$

The OP 14 is a special case of OP 11. Instead of using a decision curve of OP 14 we use a regression function

$$\sum_{i=1}^{2n} y_c^i \alpha_{c,i}^* K(\vec{x}_i, \vec{x}) + b_c^* = 0 \rightarrow g^*(\vec{x}) = \sum_{i=1}^n y_c^i \alpha_{r,i}^* K(\vec{x}_i, \vec{x}) + \sum_{i=n+1}^{2n} y_c^i \beta_{r,i-n}^* K(\vec{x}_i, \vec{x}) + b_r^*, \quad (2.24)$$

where  $\alpha_{r,i}^* = \alpha_{c,i}^*$ ,  $\beta_{r,i-n}^* = \alpha_{c,i}^*$ ,  $b_r^* = b_c^*$ . Because data are duplicated we can merge the sums in the final regression and we get

$$g^*(\vec{x}) = \sum_{i=1}^n (\alpha_{r,i}^* - \beta_{r,i}^*) K(\vec{x}_i, \vec{x}) + b_r^*. \quad (2.25)$$

In a typical scenario  $\varphi_i < 0$ , because  $\varepsilon$  is close to 0 and  $y_r^i$  is less than 1. We can notice the following property of the OP 14. Because every training example is duplicated, for every possible solution,  $n$  training examples will be always incorrectly classified except those lying on a classification decision boundary, Fig. 2.3. The KKT complementary condition for  $\varepsilon$ -SVR is the same as for  $\varphi$ -SVC after reformulation. From (2.15) and (2.16), we get

$$\alpha_i (h(\vec{x}_i) - y_r^i + \varepsilon + \xi_i) = 0 \quad (2.26)$$

$$\beta_i (-h(\vec{x}_i) + y_r^i + \varepsilon + \xi_i^*) = 0 \quad (2.27)$$

$$(C_c - \alpha_i) \xi_i = 0. \quad (2.28)$$

$$(C_c - \beta_i) \xi_i^* = 0. \quad (2.29)$$

We can also analyze some properties of the variable  $\gamma_r^i = \alpha_{r,i} - \beta_{r,i}$ . When  $\gamma_r^i = -C_c$ , then  $\alpha_i = 0$  and  $\beta_i = C_c$ , so  $g(\vec{x}_i) - y_r^i \geq \varepsilon$ . When  $\gamma_r^i = C_c$ , then  $\alpha_i = C_c$  and  $\beta_i = 0$ , so  $g(\vec{x}_i) - y_r^i \leq \varepsilon$ . When  $\gamma_r^i = 0$ , then  $\alpha_i = \beta_i$ . When  $\alpha_i = \beta_i < C_c$ , then  $\xi_i = 0$  and  $\xi_i^* = 0$ . Because  $-\varepsilon - \xi_r^* \leq g(\vec{x}_i) - y_r^i \leq \varepsilon + \xi_r^*$ , so  $-\varepsilon \leq g(\vec{x}_i) - y_r^i \leq \varepsilon$ . The case when  $\alpha_i = \beta_i = C_c$  is possible only when  $\xi_i + \xi_i^* = -2\varepsilon$ , then  $g(\vec{x}_i) - y_r^i = -\varepsilon - \xi_i$ . When  $-C_c < \gamma_r^i < 0$ , then  $\alpha_i < \beta_i$ , so  $\alpha_i < C_c$  and  $\beta_i > 0$ . From KKT we have

$$-g(\vec{x}_i) + y_r^i + \varepsilon + \xi_i^* = 0$$

and

$$\alpha_i (g(\vec{x}_i) - y_r^i + \varepsilon) = 0$$

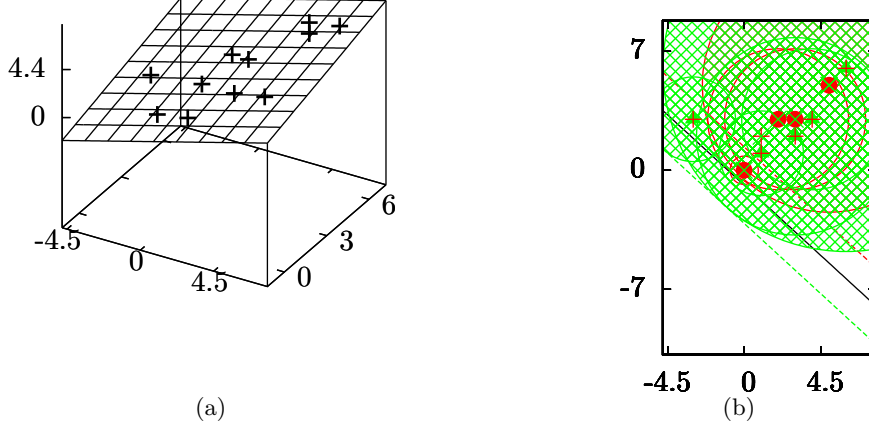


Figure 2.3: The idea of reformulating  $\varepsilon$ -SVR as  $\varphi$ -SVC. (a) Points - regression examples, a plane - a solution. (b) The transformed problem. Points - classification examples, triangles and circles - support vectors, circles filled with grid pattern - dynamic hyperspheres, a solid line - a solution, dashed lines - margin boundaries

Merging both we get

$$\alpha_i (2\varepsilon + \xi_i^*) = 0$$

So  $\xi_i^* = -2\varepsilon$  or  $\alpha_i = 0$ . For the second case,  $\beta < C_c$ , because  $\gamma_r^i > -C_c$ , so

$$g(\vec{x}_i) - y_r^i = \varepsilon .$$

Finally for the case when  $0 < \gamma_r^i < C_c$ , then  $\alpha_i > \beta_i$ , so  $\alpha_i > 0$  and  $\beta_i < C_c$ . From KKT we have

$$g(\vec{x}_i) - y_r^i + \varepsilon + \xi_i = 0$$

and

$$\beta_i (-g(\vec{x}_i) + y_r^i + \varepsilon) = 0$$

Merging both we get

$$\beta_i (2\varepsilon + \xi_i) = 0$$

So  $\xi_i = -2\varepsilon$  or  $\beta_i = 0$ . For the second case,  $\alpha < C_c$ , because  $\gamma_r^i < C_c$ , so

$$g(\vec{x}_i) - y_r^i = -\varepsilon .$$

## 2.6 Using Knowledge About a Margin with $\varepsilon$ -SVR

The  $\varepsilon$ -SVR can be reformulated as  $\varphi$ -SVC. We can additionally modify margin weights for incorporating knowledge about a margin. There are two options of incorporation

available, either we can modify  $\varphi_p$  and  $\varphi_{p+n}$  for some vector  $p$  after transforming the problem into  $\varphi$ -SVC or we can modify  $\varepsilon_u^p$  or  $\varepsilon_d^p$  in OP 8 before the transformation. In the first option, the  $\varphi_p$  and  $\varphi_{p+n}$  weights are already set according to (2.23). So adding knowledge about a margin means the manipulation of these weights. In the second option,  $\varepsilon_u^p$  or  $\varepsilon_d^p$  weights are set to some  $\varepsilon$  value for standard  $\varepsilon$ -SVR. Both approaches are equivalent, in the sense that modification of  $\varepsilon_u^p$  and  $\varepsilon_d^p$  is equivalent to the modification

$$\Delta\varphi_p = -\Delta\varepsilon_d^p , \quad (2.30)$$

$$\Delta\varphi_{p+n} = -\Delta\varepsilon_u^p . \quad (2.31)$$

## 2.7 Using Knowledge About a Margin with $\delta$ -SVR

The  $\delta$ -SVR is transformed into the classification problem, so we can use  $\varphi$ -SVC instead of standard SVC for incorporating knowledge about a margin. From some point  $p$ , we set  $\varphi_p$  for the original example, and  $\varphi_{p+n}$  for the duplicated one.

## 2.8 Changing the Output Curve

After modification of  $\varphi_i$  weights, either the output curve stays the same, or it is changed. For example, let's assume that we modify only a one example  $\vec{p}$  and  $\varphi_{\vec{p}}$  is equal to zero before the modification. When  $y_{\vec{p}}h^*(\vec{p}) > 1$ , then after setting  $0 < \varphi_{\vec{p}} \leq y_{\vec{p}}h^*(\vec{p}) - 1$  the solution remains the same. When we set  $\varphi_{\vec{p}} > y_{\vec{p}}h^*(\vec{p}) - 1$ , the solution will be different, but not necessarily a decision boundary. Particularly, setting  $\varphi_{\vec{p}} > 0$  can increase a slack variable, when a value of  $C_{\vec{p}}$  is small, or it can cause decrease of the geometric margin (increase of  $\|\vec{w}_c\|$ ), and in both cases the solution can stay the same. We can see the example of changed output curve for  $\varphi$ -SVC, Fig. 2.4, for  $\delta$ -SVR, Fig. 2.5, for  $\varepsilon$ -SVR, Fig. 2.6.

Let's now analyze in details changing the output curve. First consider changing the  $\vec{C}_c$  parameter. We can see that for OP 7, if we increase value of  $C_c^p$  for some point  $p$  for which  $\xi_c^p = 0$ , then the solution remains the same. Consider now OP 11 and the solution with all  $\varphi_i = 0$ . When we set  $\varphi_p > 0$  for some example  $p$ , such as

$$\varphi_p \leq y_c^p h(x_c^p) - 1 + \xi_c^p , \quad (2.32)$$

the solution remains the same, in the other case variables  $\vec{w}, \vec{\xi}, b$  will be adjusted, but the decision curve can stay the same. It will stay the same in the following cases:

1. The  $\xi_c^p$  is only adjusted, to

$$\xi_{cnew}^p = \varphi_p - y_c^p h(\vec{x}_p) + 1 . \quad (2.33)$$

The objective function will raise by

$$\Delta f = C_p \Delta \xi_c^p . \quad (2.34)$$

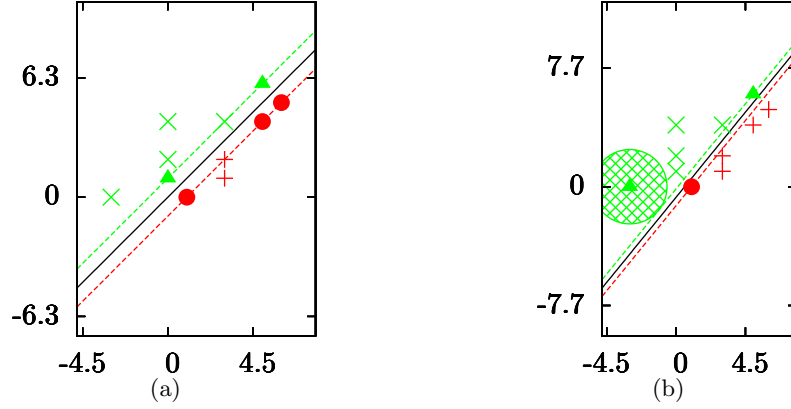


Figure 2.4: Direct curve manipulation for SVC. Points - examples, triangles and circles - support vectors, circles filled with grid pattern - dynamic hyperspheres, solid lines - solutions, dashed lines - margin boundaries. (b) The example with  $\varphi_i$  weight caused lowering the solution from (a)

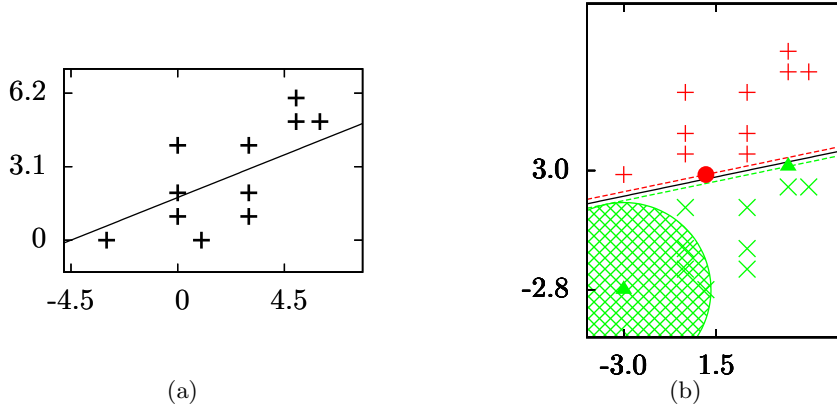


Figure 2.5: Direct curve manipulation for  $\delta$ -SVR. (a) Points - regression examples, a solid line - a solution. (b) The transformed problem. Points - classification examples, triangles and circles - support vectors, circles filled with grid pattern - dynamic hyperspheres, a solid line - a solution, dashed lines - margin boundaries. The example with  $\varphi_i$  weight caused changing the solution from (a)

2. The  $h(\vec{x}_c)$  is modified

$$h_{new}(\vec{x}_c) = ch(\vec{x}_c) , \quad (2.35)$$

where  $c > 1$ . The objective function will raise by

$$\Delta f = \frac{1}{2} \|\vec{w}\|^2 (c^2 - 1) . \quad (2.36)$$

The modification parameter  $c$  is set in a way that the following equation is satisfied

$$cy_c^p h(\vec{x}_p) = 1 - \xi_c^p + \varphi_p . \quad (2.37)$$

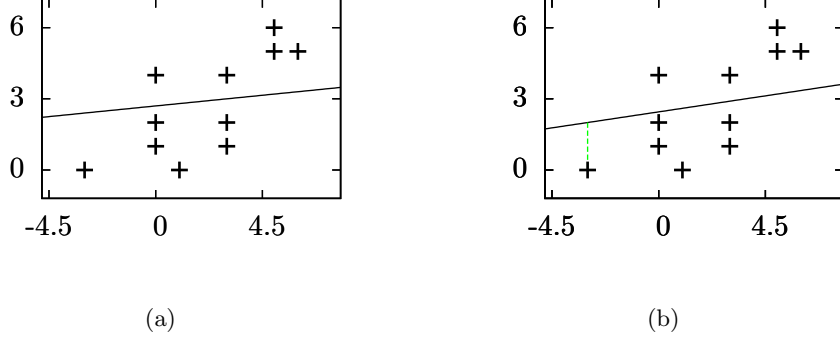


Figure 2.6: Direct curve manipulation for  $\varepsilon$ -SVR. Points - regression examples, solid lines - solutions. (b) A dotted line must cross the solution. The example with  $\varphi_i$  weight caused changing the solution from (a)

Note that in this case other non-zero  $\xi_c^i$  values have to be modified

$$\xi_c^{iold} = 1 - y_c^i h(\vec{x}_i) \quad (2.38)$$

$$\xi_c^{inew} = 1 - cy_c^i h(\vec{x}_i) \quad (2.39)$$

$$\Delta \xi_c^i = y_c^i h(\vec{x}_i) (1 - c) \quad (2.40)$$

For misclassified examples the change of an objective function will be positive, otherwise negative

$$\Delta f = \frac{1}{2} \|\vec{w}\|^2 (c^2 - 1) + \sum_{i=1, i \neq p}^n C_c^i \Delta \xi_c^i \quad (2.41)$$

- Both above changes are present simultaneously. The objective function will be changed by

$$\Delta f = \frac{1}{2} \|\vec{w}\|^2 (c^2 - 1) + C_c^p \Delta \xi_c^p + \sum_{i=1, i \neq p}^n C_c^i \Delta \xi_c^i \quad (2.42)$$

Changing the decision curve is possible, especially when  $\Delta f$  is high, because other solutions, which are able to change the decision curve, are more likely to have a better value of  $f$ . We can see that we can increase a chance for changing the decision boundary by increasing the parameter  $C_c^p$ .

## 2.9 Incorporating Linear Dependency of Function Values

Another example of application for knowledge about the margin of an example is the incorporation of the additional constraint in the form of a constant value of linear de-

pendency on function values, that the solution must satisfy, for regression it is

$$\sum_{i=1}^s s_i g(\vec{d}_i) = e, \quad (2.43)$$

where  $s_i$  are some parameters, for which  $\sum_{i=1}^s s_i \neq 0$ ,  $d_i$  are some points,  $s$  is the number of  $d_i$  points,  $e$  is a parameter,  $g$  is defined in (1.28), and for classification it is

$$\sum_{i=1}^s s_i h(\vec{d}_i) = e. \quad (2.44)$$

where  $s_i$  are some parameters, for which  $\sum_{i=1}^s s_i \neq 0$ ,  $d_i$  are some points,  $s$  is the number of  $d_i$  points,  $e$  is a parameter.

We will show how to incorporate this constraint to  $\varphi$ -SVC,  $\delta$ -SVR, and  $\varepsilon$ -SVR. Knowledge about the margin of an example is used in incorporation of this constraint to  $\varphi$ -SVC and  $\varepsilon$ -SVR.

### 2.9.1 Incorporating Linear Dependency on Function Values to $\varphi$ -SVC

Incorporating (2.44) to OP 11 leads to the  $\varphi$ -SVC optimization problem without the offset, which is the SVC optimization problem without the offset OP 4 with additional margin weights (and  $\vec{C}_c$  weights for completeness)

**OP 15.**

$$\min_{\vec{w}_c, \vec{\xi}_c} f(\vec{w}_c, \vec{\xi}_c) = \frac{1}{2} \|\vec{w}_c\|^2 + \vec{C}_c \cdot \vec{\xi}_c \quad (2.45)$$

subject to

$$y_c^i h(\vec{x}_i) \geq 1 - \xi_c^i + \varphi_i \quad (2.46)$$

$$\vec{\xi} \geq 0 \quad (2.47)$$

for  $i \in \{1, \dots, n\}$ , where

$$\vec{C}_c \gg 0 \quad (2.48)$$

$$\varphi_c^i \in \mathbb{R} \quad (2.49)$$

$$h(\vec{x}_i) = \vec{w}_c \cdot \vec{x}_i. \quad (2.50)$$

The dual problem compared to OP 5 contains additionally margin weights in an objective function (derivation in Appendix A.3)

**OP 16.**

$$\max_{\vec{\alpha}} d(\vec{\alpha}) = \vec{\alpha} \cdot (1 + \vec{\varphi}) - \frac{1}{2} \vec{\alpha}^T Q \vec{\alpha} \quad (2.51)$$

subject to

$$0 \leq \vec{\alpha} \leq \vec{C}_c \quad (2.52)$$

where  $Q_{ij} = y_i y_j K(\vec{x}_i, \vec{x}_j)$ , for all  $i, j \in \{1, \dots, n\}$ .

Incorporating (2.44) to OP 11 also leads to the new kernel in the form of transformation of the input vectors (derivation in Appendix A.4)

$$\vec{x} \rightarrow \vec{x} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i, \quad (2.53)$$

$$K(\vec{x}, \vec{y}) = \left( \vec{x} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i \right) \left( \vec{y} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i \right). \quad (2.54)$$

This is a symmetrical kernel. We set the margin weights as

$$\varphi_i = \varphi_{old} - y_i \frac{e}{\sum_{i=1}^s s_i}. \quad (2.55)$$

We propose also a nonlinear kernel by further kernelization of (2.54), we get

$$K(\vec{x}, \vec{y}) = K_o(\vec{x}, \vec{y}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(\vec{x}, \vec{d}_i) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(\vec{y}, \vec{d}_i) \quad (2.56)$$

$$+ \frac{1}{(\sum_{i=1}^s s_i)^2} \sum_{i=1}^s \sum_{j=1}^s s_i s_j K_o(\vec{d}_i, \vec{d}_j). \quad (2.57)$$

This is also a symmetrical kernel. In the final solution, we use the transformation of the input vectors (2.53) only for the training vectors present in the term  $\vec{w}_c$ , therefore for the solution, we get the following kernel

$$K(\vec{x}_j, \vec{x}) = K_o(\vec{x}_j, \vec{x}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(\vec{d}_i, \vec{x}), \quad (2.58)$$

where  $j \in \{1, \dots, n\}$ . Note that this is a kernel between a training vector and  $\vec{x}$  used only for defining the solution. For the final solution, we compute the offset as

$$b = \frac{1}{\sum_{i=1}^s s_i} \left( e - \sum_{i=1}^s s_i \sum_{j=1}^n \alpha_j y_c^j K_o(\vec{x}_j, \vec{d}_i) \right) \quad (2.59)$$

$$+ \frac{1}{(\sum_{i=1}^s s_i)^2} \sum_{i=1}^s s_i \sum_{j=1}^n \sum_{k=1}^s \alpha_j y_c^j s_k K_o(\vec{d}_k, \vec{d}_i). \quad (2.60)$$

### 2.9.2 Incorporating the Linear Dependency on Function Values to $\varepsilon$ -SVR

We incorporated (2.43) to OP 6. Because  $\varepsilon$ -SVR is a special case of  $\varphi$ -SVC, and we have derived already the incorporation for  $\varphi$ -SVC, for  $\varepsilon$ -SVR we set the following weights

$$\varphi_i = y_c^i y_r^i - \varepsilon - 1 - y_c^i \frac{e}{\sum_{i=1}^s s_i} \quad (2.61)$$

for  $i \in \{1, \dots, 2n\}$ . We also use input space transformation. We use  $\varphi$ -SVC without the offset.



### 2.9.3 Incorporating Linear Dependency on Function Values to $\delta$ -SVR

For completeness we show here the incorporation of linear dependency on function values to  $\delta$ -SVR. We do not use margin weights in this case. We use kernels developed for  $\delta$ -SVR, (1.49). Incorporating (2.43) to  $\delta$ -SVR leads to the SVC optimization problem without the offset OP 4 with the new kernel (derivation in Appendix A.5)

$$K(\vec{x}, \vec{y}) = \left( x_{\text{red}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}} \right) \left( y_{\text{red}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}} \right) \quad (2.62)$$

$$+ \left( x_{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) \left( y_{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right), \quad (2.63)$$

where  $d_{i,\text{red}} = (d_i^1, \dots, d_i^m)$ . This is a symmetrical kernel. We propose also a nonlinear kernel by further kernelization of (2.62), we get

$$K(\vec{x}, \vec{y}) = K_o(x_{\text{red}}, y_{\text{red}}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(d_{i,\text{red}}, x_{\text{red}}) \quad (2.64)$$

$$- \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(d_{i,\text{red}}, y_{\text{red}}) + \frac{1}{(\sum_{i=1}^s s_i)^2} \sum_{i=1}^s \sum_{j=1}^s s_i s_j K_o(d_{i,\text{red}}, d_{j,\text{red}}) \quad (2.65)$$

$$+ x_{m+1} y_{m+1} - x_{m+1} \frac{e}{\sum_{i=1}^s s_i} - y_{m+1} \frac{e}{\sum_{i=1}^s s_i} + \frac{e^2}{(\sum_{i=1}^s s_i)^2}. \quad (2.66)$$

This is also a symmetrical kernel. In the final solution, we use the kernelization process only for the training vectors present in the term  $\vec{w}_c$ , therefore for the solution, we get the following kernel

$$K(\vec{x}_j, \vec{x}) = K_o(x_{j,\text{red}}, x_{\text{red}}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(d_{i,\text{red}}, x_{\text{red}}) + \left( x_j^{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) x_{m+1}, \quad (2.67)$$

where  $j \in \{1, \dots, n\}$ . Note that this is a kernel between a training vector and  $\vec{x}$  used only for defining the solution. For the final solution, we compute the offset as

$$b_c = - \frac{1}{\sum_{i=1}^s s_i} e \sum_{j=1}^{2n} \alpha_j y_c^j \left( x_j^{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \sum_{j=1}^{2n} \alpha_j y_c^j K_o(x_{j,\text{red}}, d_{i,\text{red}}) \quad (2.68)$$

$$+ \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \sum_{j=1}^{2n} \alpha_j y_c^j \frac{1}{\sum_{i=1}^s s_i} \sum_{k=1}^s s_k K_o(d_{k,\text{red}}, d_{i,\text{red}}). \quad (2.69)$$

## 2.10 Incorporating Inequalities with Function Values

Another example of application for knowledge about the margin of an example is the incorporation of the additional constraints in the form of inequalities with function values for training points for regression case:

$$g(\vec{x}_i) \geq a_i \quad (2.70)$$

for  $i = 1..n$ , where  $a_i$  are some parameters,  $g$  is defined in (1.28). We propose soft incorporation by changing  $\varphi_i$  values. For  $\varepsilon$ -SVR, it leads to the modification of the  $\varepsilon_u^i$  value in OP 8:

$$\varepsilon_{u,\text{new}}^i = \min(\varepsilon_{u,\text{old}}^i, y_r^i - a_i) \quad (2.71)$$

*Proof.* From (1.37) we have

$$g(\vec{x}_i) \geq y_r^i - \varepsilon_u^i - \xi_r^i \quad (2.72)$$

It is a soft incorporation so

$$g(\vec{x}_i) \geq y_r^i - \varepsilon_u^i \quad (2.73)$$

So we should set  $\varepsilon_u^i$  such as

$$y_r^i - \varepsilon_u^i \geq a_i \quad (2.74)$$

$$\varepsilon_u^i \leq y_r^i - a_i \quad (2.75)$$

so

$$\varepsilon_{u,\text{new}}^i = \min(\varepsilon_{u,\text{old}}^i, y_r^i - a_i) \quad (2.76)$$

□

For  $\delta$ -SVR, we set a value of the  $\delta_d^i$  parameter in OP 10:

$$\delta_{d,\text{new}}^i = \min(\delta_{d,\text{old}}^i, y_r^i - a_i) \quad (2.77)$$

*Proof.* After shifting data the barrier for the function is the shifted point. It will be more restrict barrier, because it is a barrier for the margin boundary. It implies the barrier for the function. We have

$$y_r^i - \delta_d^i \geq a_i \quad (2.78)$$

$$\delta_d^i \leq y_r^i - a_i \quad (2.79)$$

So we modify the shifting parameter  $\delta_d^i$  as

$$\delta_{d,\text{new}}^i = \min(\delta_{d,\text{old}}^i, y_r^i - a_i) \quad (2.80)$$

□

Although we do not modify directly  $\varphi_i$  weights, they are modified indirectly for  $\varepsilon$ -SVR, because changing  $\varepsilon_i$  leads to changing  $\varphi_i$ . Changing  $\delta_i$  can be interpreted as changing  $\varepsilon_i$ .

Similar incorporation scheme exists for

$$g(\vec{x}_i) \leq a_i . \quad (2.81)$$

For  $\varepsilon$ -SVR, it leads to the modification of the  $\varepsilon_d^i$  value in OP 8:

$$\varepsilon_{d,\text{new}}^i = \min(\varepsilon_{d,\text{old}}^i, a_i - y_r^i) \quad (2.82)$$

*Proof.* From (1.38) we have

$$g(\vec{x}_i) \leq \varepsilon_d^i + y_r^i + \xi_r^i \quad (2.83)$$

It is a soft incorporation so

$$g(\vec{x}_i) \leq \varepsilon_d^i + y_r^i \quad (2.84)$$

So we should set  $\varepsilon_d^i$  such as

$$\varepsilon_d^i + y_r^i \leq a_i \quad (2.85)$$

$$\varepsilon_d^i \leq a_i - y_r^i \quad (2.86)$$

so

$$\varepsilon_{d,\text{new}}^i = \min(\varepsilon_{d,\text{old}}^i, a_i - y_r^i) \quad (2.87)$$

□

For  $\delta$ -SVR, we set a value of the  $\delta_u^i$  parameter in OP 10:

$$\delta_{u,\text{new}}^i = \min(\delta_{u,\text{old}}^i, a_i - y_r^i) \quad (2.88)$$

*Proof.* We have

$$y_r^i + \delta_u^i \leq a_i \quad (2.89)$$

$$\delta_u^i \leq a_i - y_r^i \quad (2.90)$$

So we modify the shifting parameter  $\delta_u^i$  as

$$\delta_{u,\text{new}}^i = \min(\delta_{u,\text{old}}^i, a_i - y_r^i) . \quad (2.91)$$

□

We can also increase values of proper  $C_i$  parameters to improve fulfillment of the constraints.

Another incorporation type is of the same form as (2.70) but defined for a new point  $x_{n+1}^{\rightarrow}$  without defined  $y_r^{n+1}$

$$g(x_{n+1}^{\rightarrow}) \geq a_{n+1} , \quad (2.92)$$

where  $a_{n+1}$  is a parameter,  $g$  is defined in (1.28). We propose soft incorporation. We will use a special version of  $\varepsilon$ -SVR, where we allow presence of only one constraint either (1.37) or (1.38). So for (2.92) we will have only one constraint that is (1.37). We set

$$y_r^{n+1} = a_{n+1} \quad (2.93)$$

After substituting above to (2.71) we get

$$\varepsilon_{u,\text{new}}^{n+1} = \min \left( \varepsilon_{u,\text{old}}^{n+1}, 0 \right) \quad (2.94)$$

One constraint in the formulation means that while transforming to the  $\varphi$ -SVC only one classification point will be created, in this case with  $y_c^{n+1} = 1$ .

For  $\delta$ -SVR, we propose also soft incorporation. We will use a special version of  $\delta$ -SVR, where we allow presence of points which are not duplicated but shifted either up or down. We define such point with  $y_c^{n+1} = -1$  and we set

$$y_r^{n+1} = a_{n+1} \quad (2.95)$$

and we shift it down by

$$\delta_{d,\text{new}}^{n+1} = \min \left( \delta_{d,\text{old}}^{n+1}, 0 \right) . \quad (2.96)$$

For

$$g(x_{n+1}^{\vec{}}) \leq a_{n+1} , \quad (2.97)$$

For (2.97) we will have only one constraint that is (1.38). We set

$$y_r^{n+1} = a_{n+1} \quad (2.98)$$

After substituting above to (2.82), we get

$$\varepsilon_{d,\text{new}}^{n+1} = \min \left( \varepsilon_{d,\text{old}}^{n+1}, 0 \right) \quad (2.99)$$

One constraint in the formulation means that while transforming to the  $\varphi$ -SVC only one classification point will be created, in this case with  $y_c^{n+1} = -1$ .

For  $\delta$ -SVR, we define a point with  $y_c^{n+1} = 1$  and we set

$$y_r^{n+1} = a_{n+1} \quad (2.100)$$

and we shift it up by

$$\delta_{u,\text{new}}^{n+1} = \min \left( \delta_{u,\text{old}}^{n+1}, 0 \right) . \quad (2.101)$$

## 2.11 Reduce a Model with $\varphi$ -SVC

Various methods for reducing the complexity of the output model were widely investigated, [4]. In particular, the reduction by removing support vectors was also analyzed in [3] for regression problems.

Sparse models have an advantage of faster post-processing such as testing new examples. For highly nonseparable data with a linear decision boundary, SVC has multiple support vectors. The general idea of constructing even more sparse solutions than SVC is to find a solution spanned on the given number of support vectors as close to the original SVC solution as possible. The most representative method for this idea is presented in [16]. Another example for this group are reduced support vector machines (RSVM), which randomly selects support vectors to a reduced set, [2, 7]. The alternative approach is to replace SVM with another training method that is designed for returning sparse solutions. The most representative method for this group is a greedy approach that adds basis functions to the solution until no progress in optimizing a cost function is made, [4].

Our proposed approach, [9, 11], is conceptually closer to the first idea. After SVM is trained, we remove randomly selected data vectors, and run again SVM on a reduced training set with incorporated prior knowledge about the original solution, Fig. 2.7, Fig. 2.8, Fig. 2.9. A concept of randomly removing support vectors was presented in [3]. For highly small number of support vectors removing some of them would definitely lead to decreasing generalization performance. The proposed method generates reduced models from the original full model with incorporated knowledge about the margin of an example. The procedure of generating knowledge about the margin of an example is as following. First,  $\varphi_i$  weights are automatically generated from an existing solution as

$$\varphi_i = y_i h^*(\vec{x}_i) - 1 \quad (2.102)$$

for all training examples. For  $\varepsilon$ -SVR, instead of modifying  $\varphi_i$  weights, we can alternatively modify  $\varepsilon_u^i$  and  $\varepsilon_d^i$  weights in OP 8:

$$\begin{aligned} \varepsilon_u^i &= y_r^i - g^*(\vec{x}_i) \quad , \\ \varepsilon_d^i &= g^*(\vec{x}_i) - y_r^i \quad . \end{aligned}$$

After that, a reduced model is generated by removing a bunch of data vectors – randomly selected data vectors, with maximal removal ratio of  $p\%$  off all training vectors, where  $p$  is a configurable parameter. Not that in this set, there are also support vectors that will be removed. Finally, we run  $\varphi$ -SVC with a reduced data set.

## 2.12 Generation of Prior Knowledge

When prior knowledge comes from the external source, we can directly use it for comparison of accuracy of the models. Otherwise, we have to generate prior knowledge from the available data. Mangasarian et al. [6] divide a data set to two sets (about 20% and

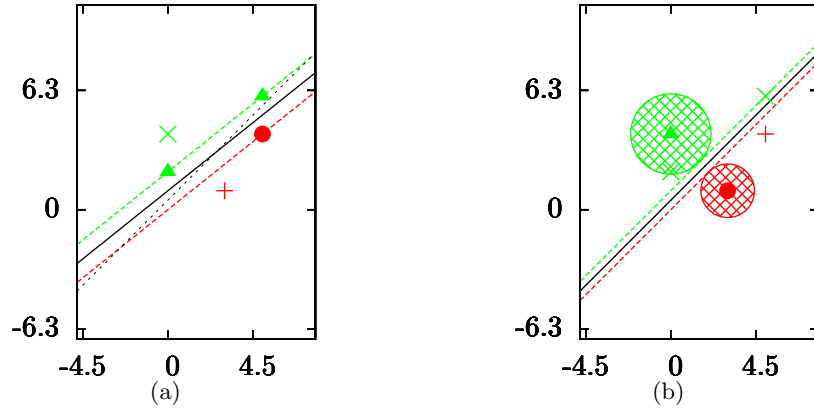


Figure 2.7: The idea of reduced SVC. Points - examples after reducing, triangles and circles - support vectors, solid lines - solutions, dotted lines - original solutions before reduction, dashed lines - margin boundaries. (b) Circles filled with grid pattern - dynamic hyperspheres

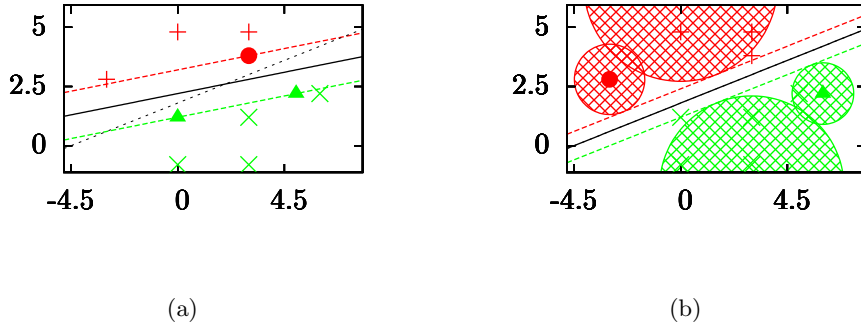


Figure 2.8: The idea of reduced  $\delta$ -SVR. Points - examples after reducing, triangles and circles - support vectors, solid lines - solutions, dotted lines - original solutions before reduction, dashed lines - margin boundaries. (b) Circles filled with grid pattern - dynamic hyperspheres

80%), and they generate manually prior knowledge from the first set, and use it in the other. We proposed a slightly modified procedure. We generate knowledge from the whole training data set, and then we remove some of data vectors. The procedure is as follows:

1. Find a solution of SVM problem without prior knowledge.
2. Extract prior knowledge from the solution.
3. Remove randomly  $p\%$  of input data.
4. Find the solution with and without prior knowledge on a reduced data set.

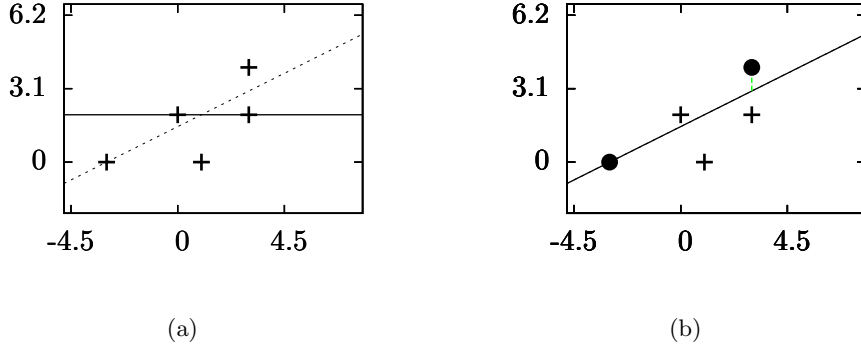


Figure 2.9: The idea of reduced  $\varepsilon$ -SVR. Points - examples after reducing, filled circles - support vectors, solid lines - solutions, dotted lines - original solutions before reduction. (b) The dashed lines must cross the solution

## 2.13 Experiments

In the first experiment, we compare the performance of the SVM without and with knowledge about the margin of an example depends on the removal ratio  $p$ . For comparison purposes, a reduced model is the same for both methods. We use our implementation of all methods. In the second experiment, we compare the performance of both methods related to various number of support vectors.

For all data sets, every feature is scaled linearly to  $[0, 1]$  including an output. For variable parameters like the  $C$ ,  $\sigma$  for the RBF kernel,  $\delta$  for  $\delta$ -SVR, and  $\varepsilon$  for  $\varepsilon$ -SVR, we use a grid search method for finding best values. The number of values searched by the grid method is a trade-off between an accuracy and a speed of simulations. Note that for a particular data set it is possible to use more accurate grid searches than for massive tests with multiple number of simulations.

### 2.13.1 First Experiment

For synthetic data, we compare both methods on data generated from particular functions with added Gaussian noise for output values. For  $p = 70$ , the SVC with knowledge about a margin has smaller generalization error for every test except one, Table 2.1. The testing performance gain varies from 0% to 3%. The number of support vectors is slightly smaller for the method with knowledge about a margin. For the regression case, for  $p = 70$ , the  $\varepsilon$ -SVR with knowledge about a margin has smaller generalization error for every test, Table 2.2. The testing performance gain varies from 0% to 50%. The number of support vectors is bigger for the method with knowledge about a margin.

The real world data sets were taken from the LibSVM site, [1, 5], except stock price data. The stock price data consist of monthly prices of the Dow Jones Industrial Average (DJIA) index from 1898 up to 2010. We generated the training set as follows: for every

Table 2.1: Performance of  $\varphi$ -SVC for reduced models for synthetic data. Column descriptions: *id* – id of the test, *a function* – a function used for generating data  $y_1 = \sum_{i=1}^{\dim-1} x_i$ ,  $y_4, y_5 = \left(\sum_{i=1}^{\dim-1} x_i\right)^{\ker P}$ ,  $y_6 = 0.5 \sum_{i=1}^{\dim-1} \sin 10x_i + 0.5$ , *ker* – a kernel with a kernel parameter (for a polynomial kernel it is a dimension, for the RBF kernel it is  $\sigma$ ), *idRef* – a reference to the test, *te12M* – a percent average difference in correctly classified examples for testing data, if greater than 0 than a method with knowledge about a margin is better, *s1* – the average number of support vectors for a method without knowledge about a margin, *s2* – the average number of support vectors for a method with knowledge about a margin

(a)			(b)			
id	function	ker	idRef	te12M	s1	s2
0	$y_1$	denseLinear 0.0	0	2.4	23	19
1	$y_2 = 3y_1$	denseLinear 0.0	1	3.22	23	19
2	$y_3 = 1/3y_1$	denseLinear 0.0	2	-0.63	24	21
3	$y_4$	densePolynomial 5.0	3	1.19	23	16
4	$y_5$	denseRBF 0.25	4	0.92	25	20
5	$y_6$	denseRBF 0.25	5	0.43	26	25

Table 2.2: Performance of  $\varphi$ -SVC for reduced models for synthetic data, for regression. Column descriptions: *id* – id of the test, *a function* – a function used for generating data  $y_1 = \sum_{i=1}^{\dim-1} x_i$ ,  $y_2 = \left(\sum_{i=1}^{\dim-1} x_i\right)^{\ker P}$ ,  $y_3 = 0.5 \sum_{i=1}^{\dim-1} \sin 10x_i + 0.5$ , *ker* – a kernel with a kernel parameter (for a polynomial kernel it is a dimension, for the RBF kernel it is  $\sigma$ ), *idRef* – a reference to the test, *te12M* – a percent average difference in MSE for testing data, if greater than 0 than a method with knowledge about a margin is better, *s1* – the average number of support vectors for a method without knowledge about a margin, *s2* – the average number of support vectors for a method with knowledge about a margin

(a)			(b)			
id	function	ker	idRef	te12M	s1	s2
0	$y_1$	denseLinear 0.0	0	50.05	11	20
1	$y_2$	densePolynomial 5.0	1	0.98	42	52
2	$y_3$	denseRBF 0.25	2	4.6	41	54

month the output value is a growth/fall comparing to the next month. Every feature  $i$  is a percent price change between the month and the  $i$ -th previous month. In every simulation, training data are randomly chosen, the remaining examples become test data. We can see that for variable  $p$ , the SVM with knowledge about a margin has generally smaller generalization error than without, Fig. 2.10, Fig. 2.11, Fig. 2.12, Fig. 2.13. For



Table 2.3: Performance of  $\varphi$ -SVC for reduced models for real world data. Column descriptions: *id* – id of the test, *dn* – a data set, *ker* – a kernel with a kernel parameter (for a polynomial kernel it is a dimension, for the RBF kernel it is  $\sigma$ ), *idRef* – a reference to the test, *te12M* – a percent average difference in correctly classified examples for testing data, if greater than 0 than a method with knowledge about a margin is better, *s1* – the average number of support vectors for a method without knowledge about a margin, *s2* – the average number of support vectors for a method with knowledge about a margin

(a)			(b)			
id	dn	ker	idRef	te12M	s1	s2
0	a1aAll	denseLinear 0.0	0	15.47	16	16
1	a1aAll	densePolynomial 5.0	1	3.79	11	25
2	a1aAll	denseRBF 0.00813	2	0.0	27	27
3	breast-cancer	denseLinear 0.0	3	14.26	7	8
4	breast-cancer	densePolynomial 3.0	4	24.71	5	6
5	breast-cancer	denseRBF 0.1	5	3.61	24	25
6	diabetes	denseLinear 0.0	6	9.41	20	17
7	diabetes	densePolynomial 3.0	7	7.16	16	15
8	diabetes	denseRBF 0.125	8	0.42	26	26
9	djia	denseLinear 0.0	9	0.67	23	16
10	djia	densePolynomial 5.0	10	1.26	22	16
11	djia	denseRBF 0.08333	11	1.52	29	28

$p = 70$ , the SVC with knowledge about a margin has smaller generalization error for all data sets, with smaller or equal number of support vectors in 8 out of 12 tests, Table 2.3. The testing performance gain varies from 0% to 25%. For a regression case, for  $p = 70$ , the  $\varepsilon$ -SVR with knowledge about a margin has smaller generalization error for every test, Table 2.4. The testing performance gain varies from 0% to 34%. The number of support vectors is bigger for the method with knowledge about a margin, for all regression tests.

We achieve smaller generalization error for the method with knowledge about a margin. We confirmed that knowledge about a margin in deed is able to preserve information about the original solution.

### 2.13.2 Second experiment

In the second experiment, we compare performance of the method without and with knowledge about a margin, but related to the number of support vectors. We performed tests for classification and regression. We can notice that the method with knowledge about a margin tends to achieve smaller generalization error for the similar number of support vectors, although for the RBF kernel the results are similar (Fig. 2.14, Fig. 2.15, Fig. 2.16, Fig. 2.17).

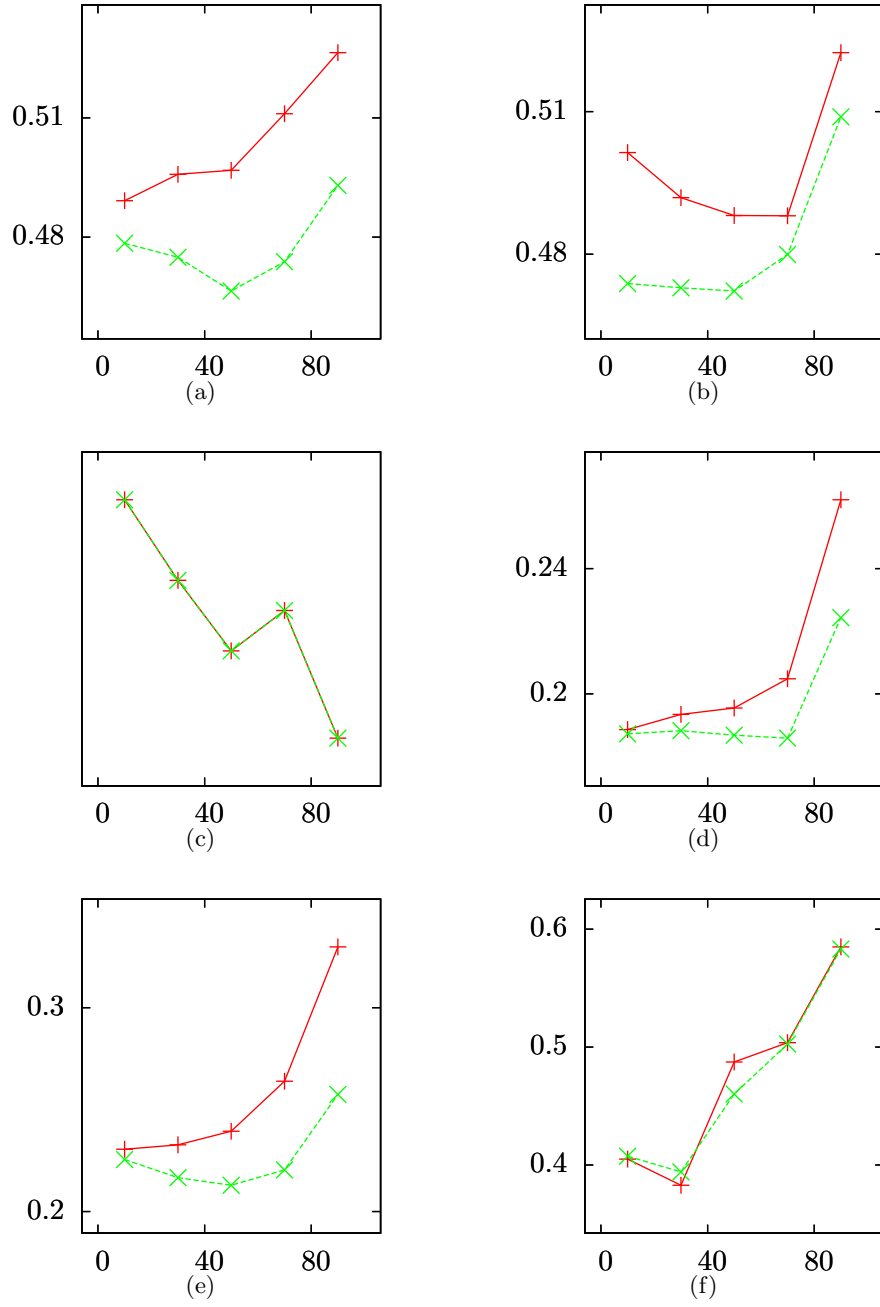


Figure 2.10: Comparison of two methods of removing support vectors for the test cases with ids 0-5 from Table 2.3. The  $x$  axis - the percent of removed data,  $y$  axis - a percent difference in misclassified testing examples, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin

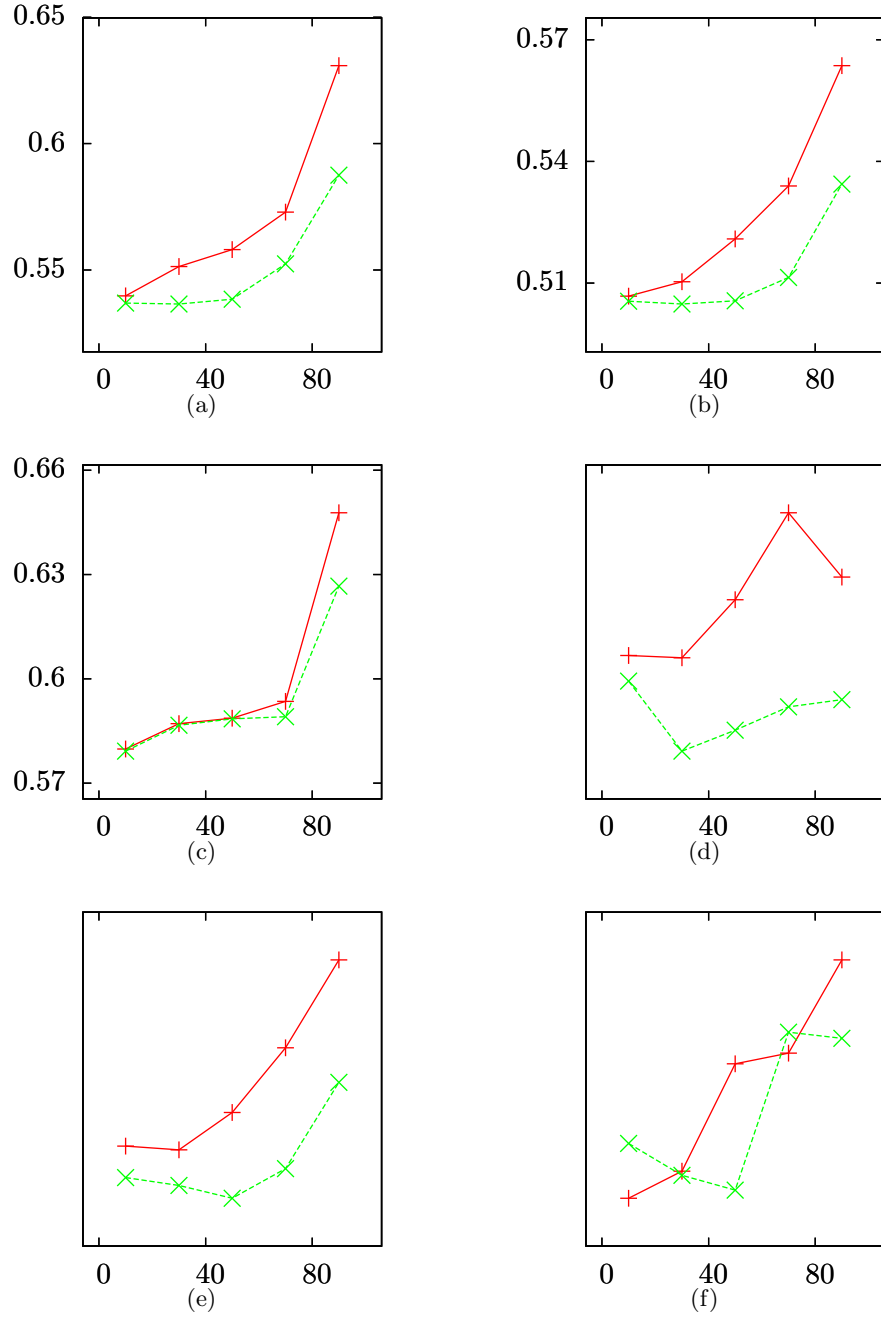


Figure 2.11: Comparison of two methods of removing support vectors for the test cases with ids 6-11 from Table 2.3, cont. The  $x$  axis - the percent of removed data,  $y$  axis - a percent difference in misclassified testing examples, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin

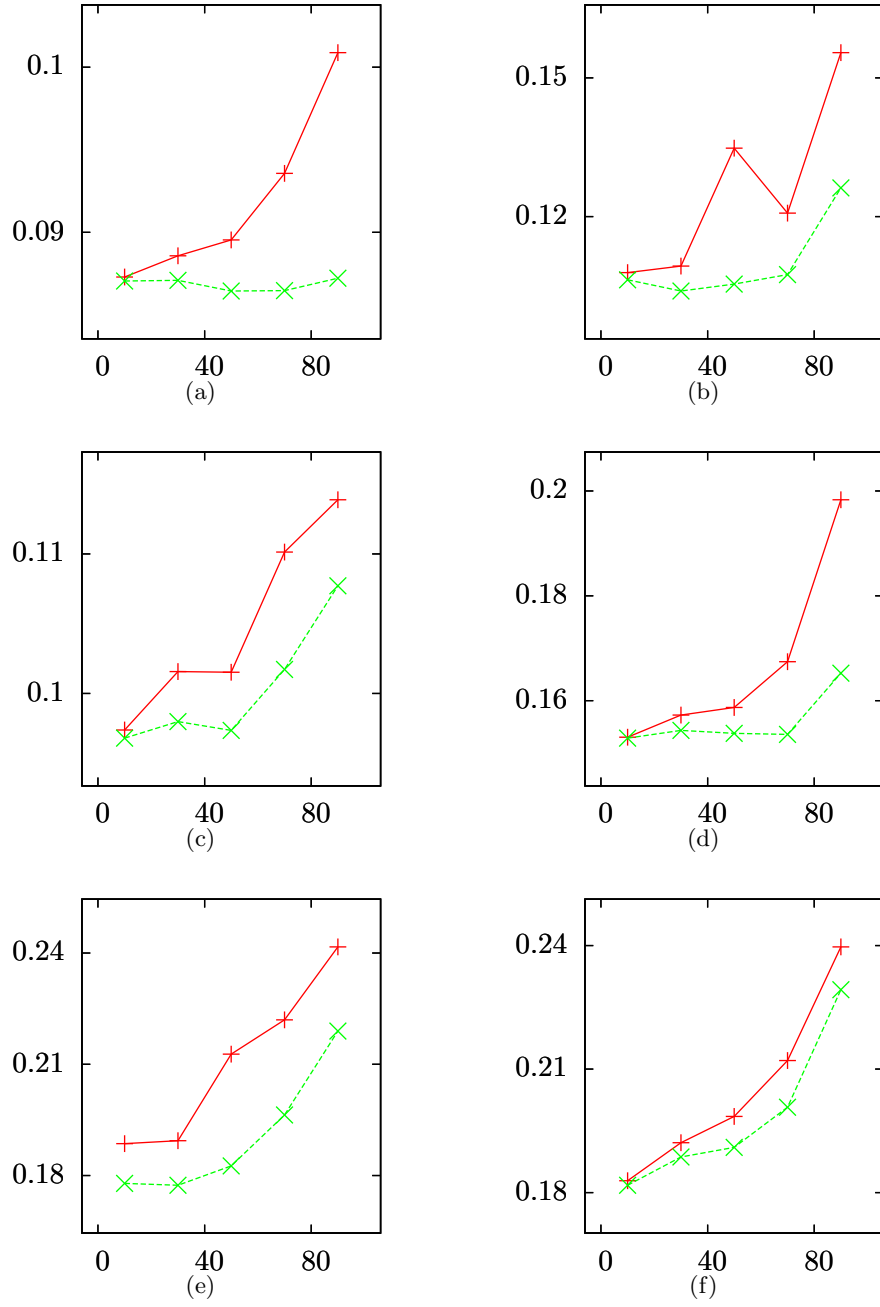


Figure 2.12: Comparison of two methods of removing support vectors for the test cases with ids 0-5 from Table 2.4. The  $x$  axis - the percent of removed data,  $y$  axis - a percent difference in MSE for testing data, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin

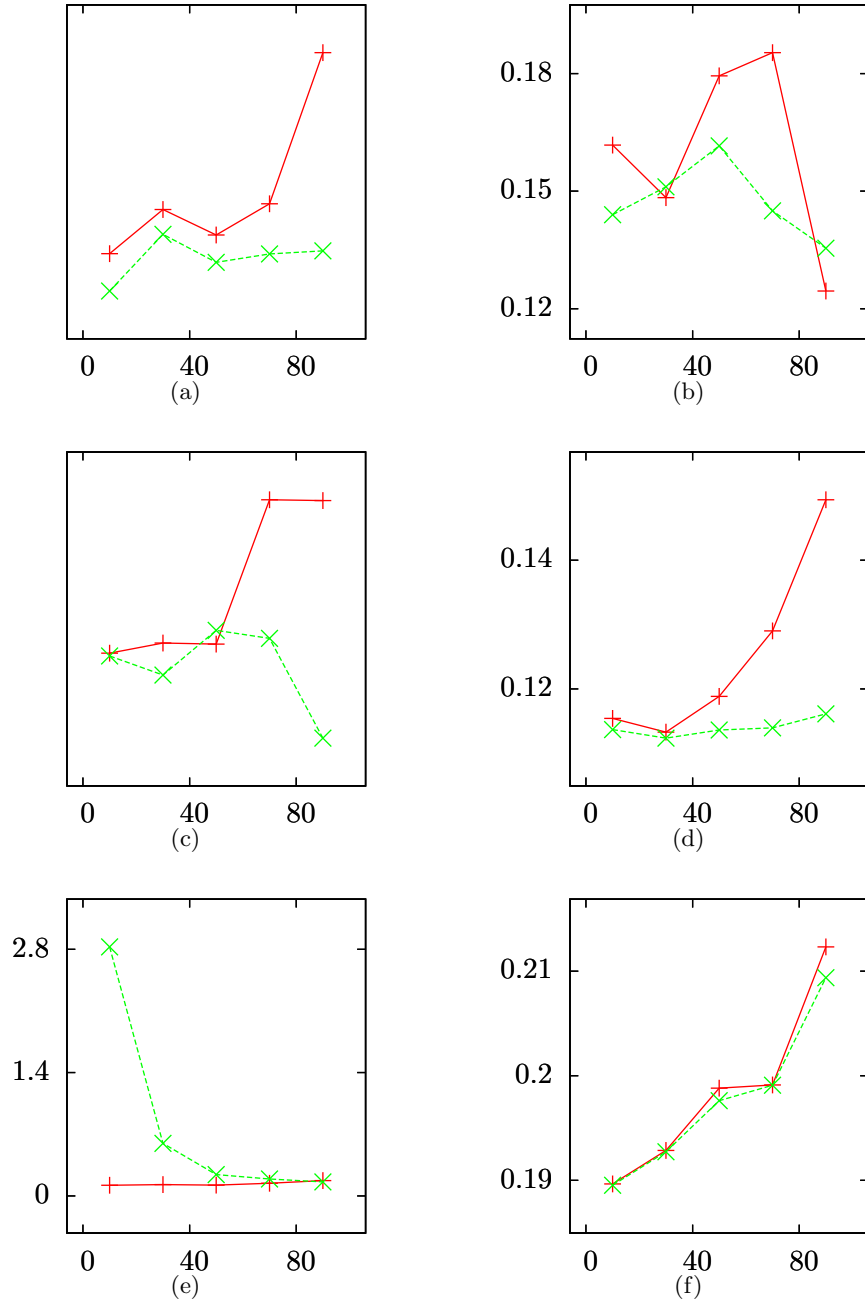


Figure 2.13: Comparison of two methods of removing support vectors for the test cases with ids 6-11 from Table 2.4, cont. The  $x$  axis - the percent of removed data,  $y$  axis - a percent difference in MSE for testing data, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin

Table 2.4: Performance of  $\varphi$ -SVC for reduced models for real world data, for regression. Column descriptions: *id* – id of the test, *dn* – a data set, *ker* – a kernel with a kernel parameter (for a polynomial kernel it is a dimension, for the RBF kernel it is  $\sigma$ ), *idRef* – a reference to the test, *te12M* – a percent average difference in MSE for testing data, if greater than 0 than a method with knowledge about a margin is better, *s1* – the average number of support vectors for a method without knowledge about a margin, *s2* – the average number of support vectors for a method with knowledge about a margin

(a)			(b)			
id	dn	ker	idRef	te12M	s1	s2
0	abalone	denseLinear 0.0	0	11.39	12	30
1	abalone	densePolynomial 5.0	1	34.86	14	48
2	abalone	denseRBF 0.125	2	11.11	18	53
3	cadata	denseLinear 0.0	3	14.08	25	47
4	cadata	densePolynomial 5.0	4	16.38	26	52
5	cadata	denseRBF 0.125	5	11.41	34	53
6	djia	denseLinear 0.0	6	5.81	9	26
7	djia	densePolynomial 5.0	7	34.16	13	48
8	djia	denseRBF 0.1	8	0.2	8	50
9	housing	denseLinear 0.0	9	19.47	17	29
10	housing	densePolynomial 5.0	10	18.39	19	53
11	housing	denseRBF 0.077	11	1.33	30	53

## 2.14 Summary

In this report, we analyzed the possibility to preserve knowledge about the original solution by using margin weights while creating reduced models. We also used knowledge about the margin of an example for incorporating the additional nonlinear constraint to the problem. We also showed that the method with knowledge about the margin of an example can achieve smaller generalization error for similar number of support vectors, for reduced models. A potential list of applications for margin weights is much broader. In future research, we plan to investigate possibility to create knowledge about a margin by experts and to apply it for time series data. Moreover, we plan to use it for combining SVM classifiers with each other in order to decrease generalization error.

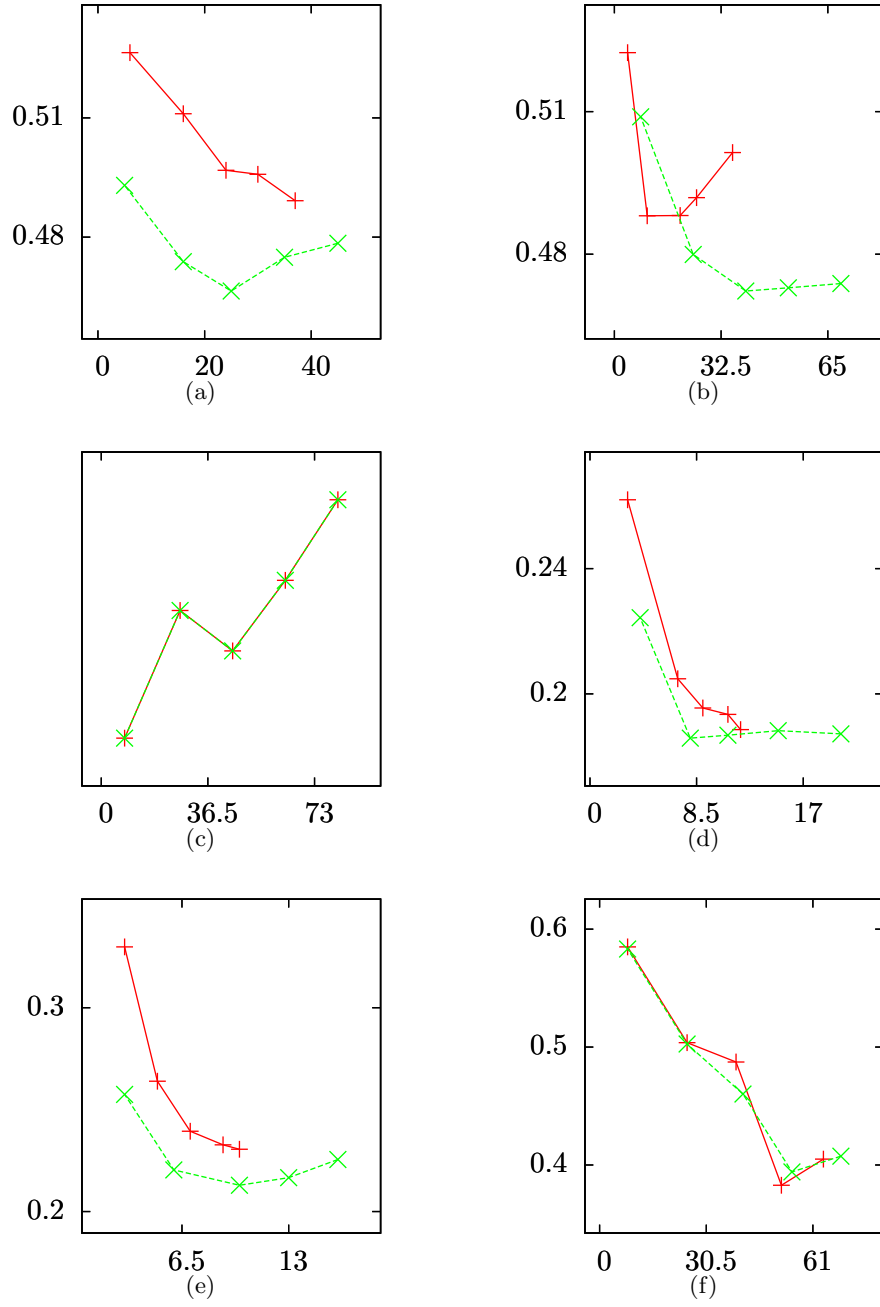


Figure 2.14: Comparison of two methods of removing support vectors for the test cases with ids 0-5 from Table 2.3. The  $x$  axis - the number of support vectors,  $y$  axis - a percent difference in misclassified testing examples, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin

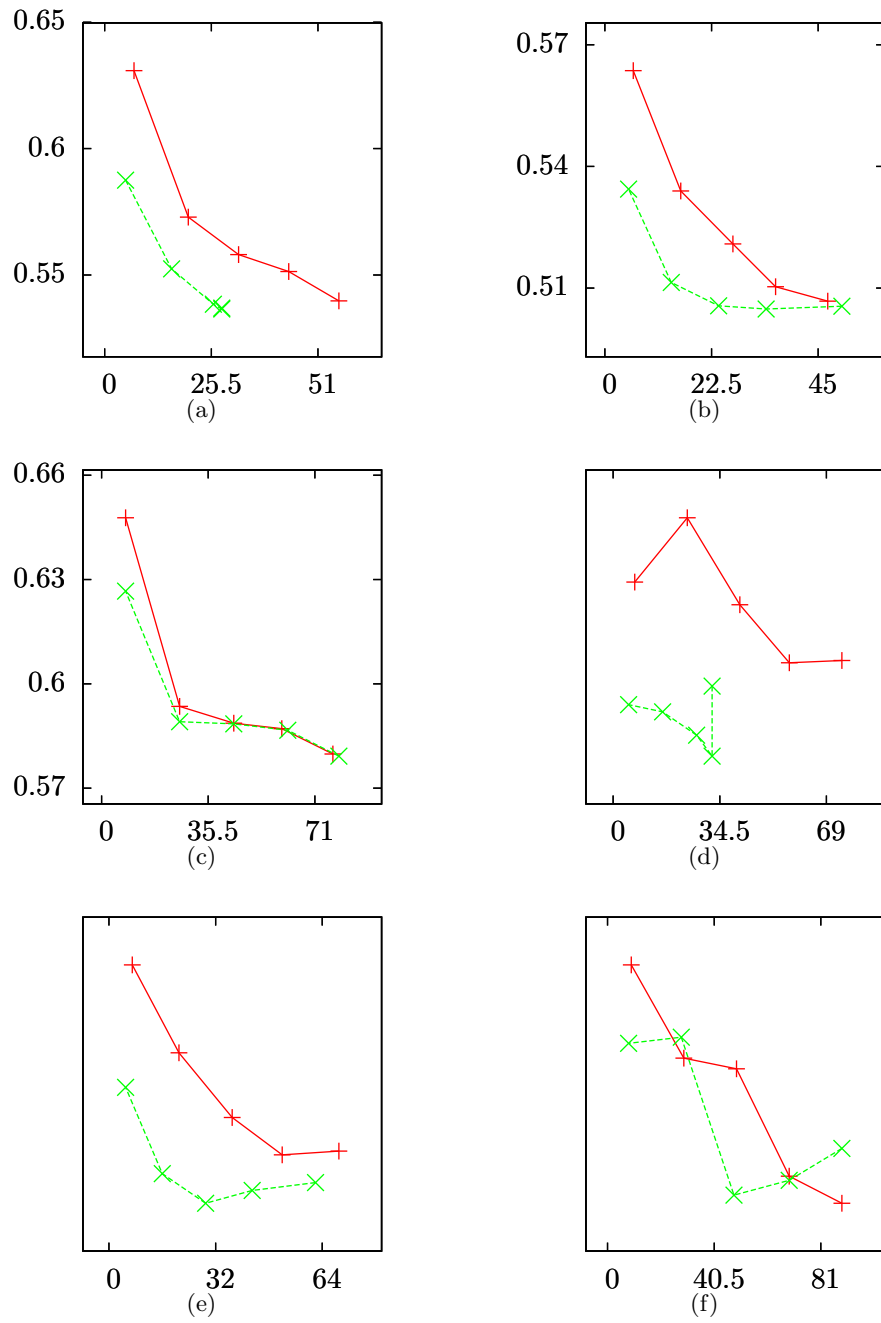


Figure 2.15: Comparison of two methods of removing support vectors for the test cases with ids 6-11 from Table 2.3, cont. The  $x$  axis - the number of support vectors,  $y$  axis - a percent difference in misclassified testing examples, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin



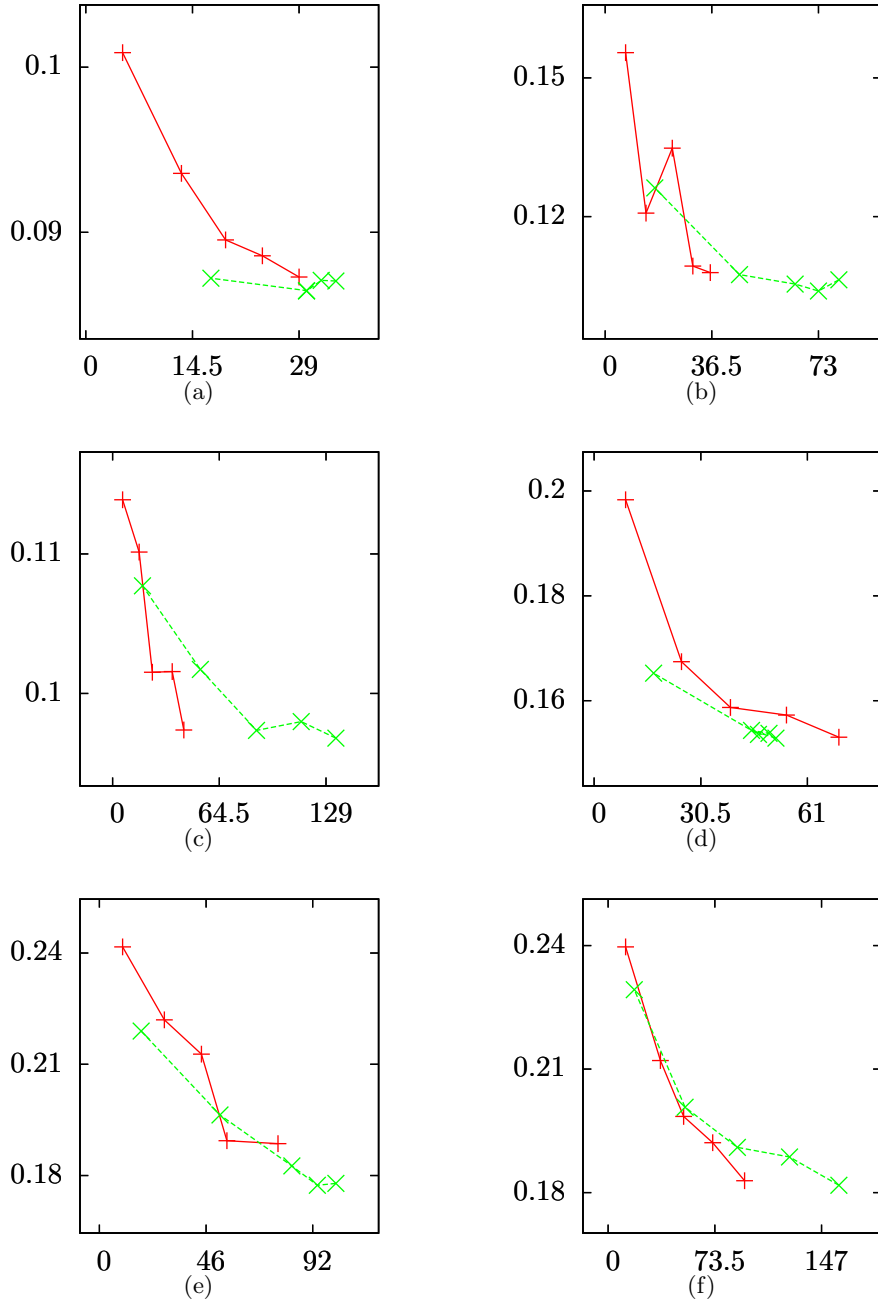


Figure 2.16: Comparison of two methods of removing support vectors for the test cases with ids 0-5 from Table 2.4. The  $x$  axis - the number of support vectors,  $y$  axis - a percent difference in MSE for testing data, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin

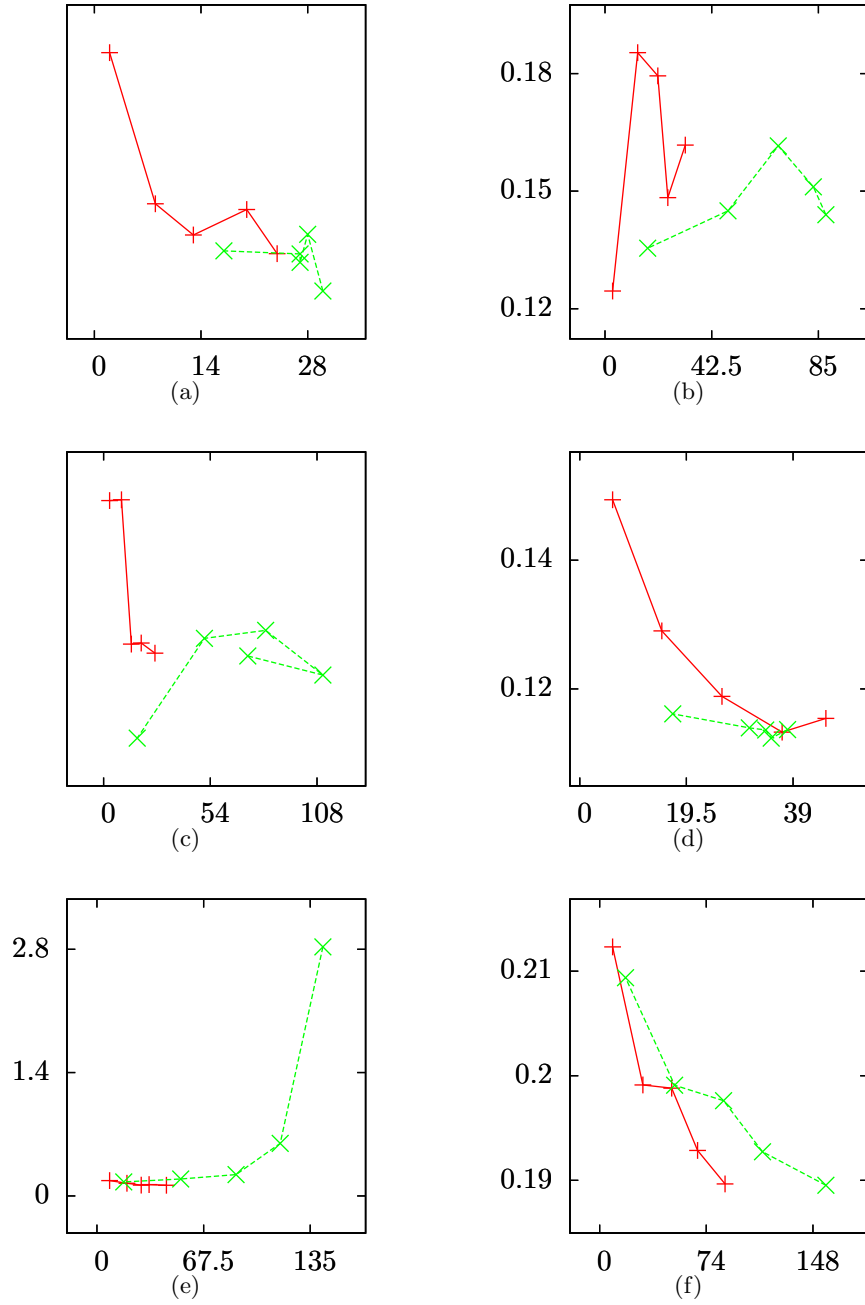


Figure 2.17: Comparison of two methods of removing support vectors for the test cases with ids 6-11 from Table 2.4, cont. The  $x$  axis - the number of support vectors,  $y$  axis - a percent difference in MSE for testing data, the line with '+' points - a random removing method, the line with 'x' points - proposed removing method with knowledge about a margin

# Appendix A

## A.1 Derivation of the Dual Form of OP 11

OP 17.

$$\max_{\vec{\alpha}, \vec{r}} d(\vec{\alpha}, \vec{r}) \quad (\text{A.1})$$

where

$$\begin{aligned} d(\vec{\alpha}, \vec{r}) &= \min_{\vec{w}, b, \vec{\xi}} t(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{r}) \\ t(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{r}) &= \frac{1}{2} \|\vec{w}\|^2 + \sum_{i=1}^n C_i \xi_i - \\ &\quad - \sum_{i=1}^n \alpha_i (y_c^i h(\vec{x}_i) - 1 + \xi_i - \varphi_i) - \sum_{i=1}^n r_i \xi_i \end{aligned}$$

subject to

$$\begin{aligned} \alpha_i &\geq 0 \\ r_i &\geq 0 \end{aligned}$$

for  $i \in \{1, \dots, n\}$ .

A partial derivative with respect to  $w_i$  is

$$\frac{\partial t(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{r})}{\partial w_i} = w_i - \sum_{j=1}^n \alpha_j y_c^j x_{ji} = 0 \quad (\text{A.2})$$

for  $i \in \{1, \dots, m\}$ . A partial derivative with respect to  $b$  is

$$\frac{\partial t(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{r})}{\partial b} = \sum_{i=1}^n \alpha_i y_c^i = 0 \quad (\text{A.3})$$

A partial derivative with respect to  $\xi_i$  is

$$\frac{\partial t(\vec{w}, b, \vec{\alpha}, \vec{\xi}, \vec{r})}{\partial \xi_i} = C_i - r_i - \alpha_i = 0 \quad (\text{A.4})$$

After substitution of above equations to  $d(\vec{\alpha}, \vec{r})$  we get

$$\begin{aligned} d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=1}^n \alpha_j y_c^j x_{ji} \right) \left( \sum_{k=1}^n \alpha_k y_c^k x_{ki} \right) \\ &- \sum_{i=1}^n \alpha_i y_c^i \left( \sum_{j=1}^m w_j x_{ij} + b \right) + \sum_{i=1}^n \alpha_i (1 + \varphi_i) + C_i \sum_{i=1}^n \xi_i \\ &- \sum_{i=1}^n \alpha_i \xi_i - \sum_{i=1}^n r_i \xi_i \end{aligned} \quad (\text{A.5})$$

$$\begin{aligned} d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^n \alpha_k \alpha_j y_c^k y_c^j x_{ki} x_{ji} - \sum_{i=1}^n \alpha_i y_c^i \sum_{j=1}^m w_j x_{ij} \\ &- b \sum_{i=1}^n \alpha_i y_c^i + \sum_{i=1}^n \alpha_i (1 + \varphi_i) \end{aligned} \quad (\text{A.6})$$

$$\begin{aligned} d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_k \alpha_j y_c^k y_c^j \sum_{i=1}^m x_{ji} x_{ki} \\ &- \sum_{i=1}^n \alpha_i y_c^i \sum_{j=1}^m x_{ij} \sum_{k=1}^n \alpha_k y_c^k x_{kj} + \sum_{i=1}^n \alpha_i (1 + \varphi_i) \end{aligned} \quad (\text{A.7})$$

$$\begin{aligned} d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_k \alpha_j y_c^k y_c^j \sum_{i=1}^m x_{ji} x_{ki} \\ &- \sum_{i=1}^n \sum_{k=1}^n \alpha_k \alpha_i y_c^k y_c^i \sum_{j=1}^m x_{ij} x_{kj} + \sum_{i=1}^n \alpha_i (1 + \varphi_i) \end{aligned} \quad (\text{A.8})$$

$$d(\vec{\alpha}, \vec{r}) = -\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_k \alpha_i y_c^k y_c^i \sum_{j=1}^m x_{ij} x_{kj} + \sum_{i=1}^n \alpha_i (1 + \varphi_i) . \quad (\text{A.9})$$

The dual form is

**OP 18.**

$$\max_{\vec{\alpha}, \vec{r}} d(\vec{\alpha}, \vec{r}) = \sum_{i=1}^n \alpha_i (1 + \varphi_i) - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_k \alpha_i y_c^k y_c^i \sum_{j=1}^m x_{ij} x_{kj} \quad (\text{A.10})$$

subject to

$$\sum_{i=1}^n \alpha_i y_c^i = 0 \quad (\text{A.11})$$

$$C_i = r_i + \alpha_i \quad (\text{A.12})$$

$$\alpha_i \geq 0 \quad (\text{A.13})$$

$$r_i \geq 0 \quad (\text{A.14})$$

for  $i \in \{1, \dots, n\}$ .

## A.2 Derivation of $\varepsilon$ -SVR Reformulation as $\varphi$ -SVC

We present derivation of  $\varepsilon$ -SVR in the form OP 8 as  $\varphi$ -SVC

**OP 19.**

$$\min_{\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*} f(\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*) = \frac{1}{2} \|\vec{w}_r\|^2 + C_r \sum_{i=1}^n (\xi_r^i + \xi_r^{*i}) \quad (\text{A.15})$$

subject to

$$y_r^i - g(\vec{x}_i) \leq \varepsilon_u^i + \xi_r^i \quad (\text{A.16})$$

$$g(\vec{x}_i) - y_r^i \leq \varepsilon_d^i + \xi_r^{*i} \quad (\text{A.17})$$

$$\vec{\xi}_r \geq 0 \quad (\text{A.18})$$

$$\vec{\xi}_r^* \geq 0 \quad (\text{A.19})$$

for  $i \in \{1, \dots, n\}$ , where

$$g(\vec{x}_i) = \vec{w}_r \cdot \vec{x}_i + b_r \quad (\text{A.20})$$

**OP 20.**

$$\min_{\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*} f(\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*) = \frac{1}{2} \|\vec{w}_r\|^2 + C_r \sum_{i=1}^{2n} \xi_r^i \quad (\text{A.21})$$

subject to

$$g(\vec{x}_i) \geq y_r^i - \varepsilon_u^i - \xi_r^i \quad (\text{A.22})$$

for  $i \in \{1, \dots, n\}$ ,

$$g(\vec{x}_i) \leq \varepsilon_d^i + y_r^i + \xi_r^i \quad (\text{A.23})$$

for  $i \in \{n+1, \dots, 2n\}$ ,

$$\vec{\xi}_r \geq 0 \quad (\text{A.24})$$

where

$$g(\vec{x}_i) = \vec{w}_r \cdot \vec{x}_i + b_r \quad (\text{A.25})$$

**OP 21.**

$$\min_{\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*} f(\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*) = \frac{1}{2} \|\vec{w}_r\|^2 + C_r \sum_{i=1}^{2n} \xi_r^i \quad (\text{A.26})$$

subject to

$$y_c^i g(\vec{x}_i) \geq y_r^i - \varepsilon_u^i - \xi_r^i \quad (y_c^i = 1) \quad (\text{A.27})$$

for  $i \in \{1, \dots, n\}$ ,

$$y_c^i g(\vec{x}_i) \geq -\varepsilon_d^i - y_r^i - \xi_r^i \quad (y_c^i = -1) \quad (\text{A.28})$$

for  $i \in \{n+1, \dots, 2n\}$ ,

$$\vec{\xi}_r \geq 0 \quad (\text{A.29})$$

where

$$g(\vec{x}_i) = \vec{w}_r \cdot \vec{x}_i + b_r \quad (\text{A.30})$$

**OP 22.**

$$\min_{\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*} f(\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*) = \frac{1}{2} \|\vec{w}_r\|^2 + C_r \sum_{i=1}^{2n} \xi_r^i \quad (\text{A.31})$$

subject to

$$y_c^i g(\vec{x}_i) \geq y_c^i y_r^i - \varepsilon_i - \xi_r^i \quad (\text{A.32})$$

$$\vec{\xi}_r \geq 0 \quad (\text{A.33})$$

for  $i \in \{1, \dots, 2n\}$ , where

$$\varepsilon_i = \varepsilon_u^i \text{ for } i \in \{1, \dots, n\} \quad , \quad (\text{A.34})$$

$$\varepsilon_i = \varepsilon_d^i \text{ for } i \in \{n+1, \dots, 2n\} \quad , \quad (\text{A.35})$$

$$g(\vec{x}_i) = \vec{w}_r \cdot \vec{x}_i + b_r \quad . \quad (\text{A.36})$$

**OP 23.**

$$\min_{\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*} f(\vec{w}_r, b_r, \vec{\xi}_r, \vec{\xi}_r^*) = \frac{1}{2} \|\vec{w}_r\|^2 + C_r \sum_{i=1}^{2n} \xi_r^i \quad (\text{A.37})$$

subject to

$$y_c^i g(\vec{x}_i) \geq 1 + y_c^i y_r^i - \varepsilon_i - \xi_r^i - 1 \quad (\text{A.38})$$

$$\vec{\xi}_r \geq 0 \quad (\text{A.39})$$

for  $i \in \{1, \dots, 2n\}$ , where

$$\varepsilon_i = \varepsilon_u^i \text{ for } i \in \{1, \dots, n\} \quad , \quad (\text{A.40})$$

$$\varepsilon_i = \varepsilon_d^i \text{ for } i \in \{n+1, \dots, 2n\} \quad , \quad (\text{A.41})$$

$$g(\vec{x}_i) = \vec{w}_r \cdot \vec{x}_i + b_r \quad . \quad (\text{A.42})$$

And we have

$$w_i = \sum_{j=1}^{2n} y_c^j \alpha_j x_{ij} = \sum_{j=1}^n \alpha_j x_{ij} - \sum_{j=n+1}^{2n} \alpha_j^* x_{ij} = \sum_{j=1}^n (\alpha_j - \alpha_j^*) x_{ij} \quad . \quad (\text{A.43})$$

### A.3 Derivation of the Dual Form of OP 16

**OP 24.**

$$\max_{\vec{\alpha}, \vec{r}} d(\vec{\alpha}, \vec{r}) \quad (\text{A.44})$$

where

$$\begin{aligned} d(\vec{\alpha}, \vec{r}) &= \min_{\vec{w}} t(\vec{w}, \vec{\alpha}, \vec{\xi}, \vec{r}) \\ t(\vec{w}, \vec{\alpha}, \vec{\xi}, \vec{r}) &= \frac{1}{2} \|\vec{w}\|^2 + \sum_{i=1}^n C_i \xi_i - \\ &\quad - \sum_{i=1}^n \alpha_i (y_i \vec{w} \cdot \vec{x}_i - 1 + \xi_i - \varphi_i) - \sum_{i=1}^n r_i \xi_i \end{aligned}$$

subject to

$$\begin{aligned}\alpha_i &\geq 0 \\ r_i &\geq 0\end{aligned}$$

for  $i \in \{1, \dots, n\}$ .

A partial derivative with respect to  $w_i$  is

$$\frac{\partial h(\vec{w}, \vec{\alpha}, \vec{\xi}, \vec{r})}{\partial w_i} = w_i - \sum_{j=1}^n \alpha_j y_j x_{ji} = 0 \quad (\text{A.45})$$

for  $i \in \{1, \dots, m\}$ . A partial derivative with respect to  $\xi_i$  is

$$\frac{\partial h(\vec{w}, \vec{\alpha}, \vec{\xi}, \vec{r})}{\partial \xi_i} = C_i - r_i - \alpha_i = 0 \quad (\text{A.46})$$

After substitution of above equations to  $d(\vec{\alpha}, \vec{r})$  we get

$$\begin{aligned}d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{i=1}^m \left( \sum_{j=1}^n \alpha_j y_j x_{ji} \right) \left( \sum_{k=1}^n \alpha_k y_k x_{ki} \right) \\ &\quad - \sum_{i=1}^n \alpha_i y_i \left( \sum_{j=1}^m w_j x_{ij} \right) + \sum_{i=1}^n \alpha_i (1 + \varphi_i) + C_i \sum_{i=1}^n \xi_i \\ &\quad - \sum_{i=1}^n \alpha_i \xi_i - \sum_{i=1}^n r_i \xi_i\end{aligned} \quad (\text{A.47})$$

$$\begin{aligned}d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^n \alpha_k \alpha_j y_k y_j x_{ki} x_{ji} - \sum_{i=1}^n \alpha_i y_i \sum_{j=1}^m w_j x_{ij} \\ &\quad + \sum_{i=1}^n \alpha_i (1 + \varphi_i)\end{aligned} \quad (\text{A.48})$$

$$\begin{aligned}d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_k \alpha_j y_k y_j \sum_{i=1}^m x_{ji} x_{ki} \\ &\quad - \sum_{i=1}^n \alpha_i y_i \sum_{j=1}^m x_{ij} \sum_{k=1}^n \alpha_k y_k x_{kj} + \sum_{i=1}^n \alpha_i (1 + \varphi_i)\end{aligned} \quad (\text{A.49})$$

$$\begin{aligned}d(\vec{\alpha}, \vec{r}) &= \frac{1}{2} \sum_{j=1}^n \sum_{k=1}^n \alpha_k \alpha_j y_k y_j \sum_{i=1}^m x_{ji} x_{ki} \\ &\quad - \sum_{i=1}^n \sum_{k=1}^n \alpha_k \alpha_i y_k y_i \sum_{j=1}^m x_{ij} x_{kj} + \sum_{i=1}^n \alpha_i (1 + \varphi_i)\end{aligned} \quad (\text{A.50})$$

$$d(\vec{\alpha}, \vec{r}) = -\frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_k \alpha_i y_k y_i \sum_{j=1}^m x_{ij} x_{kj} + \sum_{i=1}^n \alpha_i (1 + \varphi_i) \quad (\text{A.51})$$

The dual form is

**OP 25.**

$$\max_{\vec{\alpha}, \vec{r}} d(\vec{\alpha}, \vec{r}) = \sum_{i=1}^n \alpha_i (1 + \varphi_i) - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_k \alpha_i y_k y_i \sum_{j=1}^m x_{ij} x_{kj} \quad (\text{A.52})$$

subject to

$$C_i = r_i + \alpha_i \quad (\text{A.53})$$

$$\alpha_i \geq 0 \quad (\text{A.54})$$

$$r_i \geq 0 \quad (\text{A.55})$$

for  $i \in \{1, \dots, n\}$ .

## A.4 Incorporation of the Linear Dependency to $\varphi$ -SVC

We will incorporate (2.44) to OP 11. After reformulation

$$\sum_{i=1}^s s_i \vec{w}_c \cdot \vec{d}_i + b \sum_{i=1}^s s_i = e \quad (\text{A.56})$$

$$b = \frac{1}{\sum_{i=1}^s s_i} \left( e - \sum_{i=1}^s s_i \vec{w}_c \cdot \vec{d}_i \right) . \quad (\text{A.57})$$

Now we can substitute  $b$  to  $h(\vec{x})$  and we get

$$h(\vec{x}) = \vec{w}_c \cdot \vec{x} + \frac{1}{\sum_{i=1}^s s_i} \left( e - \sum_{i=1}^s s_i \vec{w}_c \cdot \vec{d}_i \right) \quad (\text{A.58})$$

after reformulation

$$h(\vec{x}) = \vec{w}_c \cdot \vec{x} - \frac{1}{\sum_{i=1}^s s_i} \vec{w}_c \cdot \sum_{i=1}^s s_i \vec{d}_i + \frac{e}{\sum_{i=1}^s s_i} . \quad (\text{A.59})$$

After substituting above to OP 11, we get the  $\varphi$ -SVC problem without the offset with the new kernel in the form of transformation of the input vectors

$$\vec{x} \rightarrow \vec{x} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i \quad (\text{A.60})$$

$$K(\vec{x}, \vec{y}) = \left( \vec{x} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i \right) \cdot \left( \vec{y} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i \right) \quad (\text{A.61})$$

and  $\varphi_i$  weights set as

$$\varphi_i = \varphi_{old} - y_i \frac{e}{\sum_{i=1}^s s_i} . \quad (\text{A.62})$$



We get nonlinear solutions by using the following way of further kernelization

$$K(\vec{x}, \vec{y}) = K_o(\vec{x}, \vec{y}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(\vec{x}, \vec{d}_i) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(\vec{y}, \vec{d}_i) \quad (\text{A.63})$$

$$+ \frac{1}{(\sum_{i=1}^s s_i)^2} \sum_{i=1}^s \sum_{j=1}^s s_i s_j K_o(\vec{d}_i, \vec{d}_j) \quad (\text{A.64})$$

This kernel is used only for solving the optimization problem. Now we derive the solution, because

$$\vec{w}_c = \sum_{j=1}^n \alpha_j y_c^j \left( \vec{x}_j - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i \right) \quad (\text{A.65})$$

we get

$$h(\vec{x}) = \vec{w}_c \cdot \vec{x} + b = \sum_{j=1}^n \alpha_j y_c^j \left( \vec{x}_j - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \vec{d}_i \right) \cdot \vec{x} + b \quad (\text{A.66})$$

$$= \sum_{j=1}^n \alpha_j y_c^j \vec{x}_j \cdot \vec{x} - \frac{1}{\sum_{i=1}^s s_i} \sum_{j=1}^n \sum_{i=1}^s \alpha_j y_c^j s_i \vec{d}_i \cdot \vec{x} + b \quad (\text{A.67})$$

and after adding internal kernels we get

$$h(\vec{x}) = \sum_{j=1}^n \alpha_j y_c^j K_o(\vec{x}_j, \vec{x}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{j=1}^n \sum_{i=1}^s \alpha_j y_c^j s_i K_o(\vec{d}_i, \vec{x}) + b \quad (\text{A.68})$$

when  $b$  is computed as

$$b = \frac{1}{\sum_{i=1}^s s_i} \left( e - \sum_{i=1}^s s_i \sum_{j=1}^n \alpha_j y_c^j \left( \vec{x}_j - \frac{1}{\sum_{i=1}^s s_i} \sum_{k=1}^s s_k \vec{d}_k \right) \cdot \vec{d}_i \right) \quad (\text{A.69})$$

$$= \frac{1}{\sum_{i=1}^s s_i} \left( e - \sum_{i=1}^s s_i \sum_{j=1}^n \alpha_j y_c^j K_o(\vec{x}_j, \vec{d}_i) + \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \sum_{j=1}^n \sum_{k=1}^s \alpha_j y_c^j s_k K_o(\vec{d}_k, \vec{d}_i) \right) \quad (\text{A.70})$$

## A.5 Incorporation of the Linear Dependency to $\delta$ -SVR

We will incorporate (2.43) to  $\delta$ -SVR. We assume that we use special kernels developed for  $\delta$ -SVR, (1.49). For  $\delta$ -SVR the condition (2.43) becomes

$$\sum_{i=1}^s s_i \frac{-\vec{w}_{\text{red}} \cdot \vec{d}_{i,\text{red}} - b_c}{w_c^{m+1}} = e, \quad (\text{A.71})$$

After transformation

$$ew_c^{m+1} + \sum_{i=1}^s s_i w_{\text{red}}^{\vec{}} \cdot d_{i,\text{red}}^{\vec{}} + \sum_{i=1}^s s_i b_c = 0 \quad (\text{A.72})$$

$$b_c = \frac{1}{\sum_{i=1}^s s_i} \left( -ew_c^{m+1} - \sum_{i=1}^s s_i w_{\text{red}}^{\vec{}} \cdot d_{i,\text{red}}^{\vec{}} \right) \quad (\text{A.73})$$

Substituting it to  $h(\vec{x})$  we get

$$h(\vec{x}) = \vec{w}_c \cdot \vec{x} + \frac{1}{\sum_{i=1}^s s_i} \left( -ew_c^{m+1} - \sum_{i=1}^s s_i w_{\text{red}}^{\vec{}} \cdot d_{i,\text{red}}^{\vec{}} \right) \quad (\text{A.74})$$

$$h(\vec{x}) = w_{\text{red}}^{\vec{}} \cdot x_{\text{red}}^{\vec{}} + w_c^{m+1} x_{m+1} + \frac{1}{\sum_{i=1}^s s_i} \left( -ew_c^{m+1} - \sum_{i=1}^s s_i w_{\text{red}}^{\vec{}} \cdot d_{i,\text{red}}^{\vec{}} \right) \quad (\text{A.75})$$

$$h(\vec{x}) = w_{\text{red}}^{\vec{}} \cdot \left( x_{\text{red}}^{\vec{}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}}^{\vec{}} \right) + w_c^{m+1} \left( x_{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) \quad (\text{A.76})$$

After substituting above to OP 11, we get the  $\varphi$ -SVC problem without the offset with the new kernel in the form

$$K(\vec{x}, \vec{y}) = \left( x_{\text{red}}^{\vec{}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}}^{\vec{}} \right) \cdot \left( y_{\text{red}}^{\vec{}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}}^{\vec{}} \right) \quad (\text{A.77})$$

$$+ \left( x_{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) \left( y_{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) \quad (\text{A.78})$$

Then

$$K(\vec{x}, \vec{y}) = x_{\text{red}}^{\vec{}} \cdot y_{\text{red}}^{\vec{}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}}^{\vec{}} \cdot x_{\text{red}}^{\vec{}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}}^{\vec{}} \cdot y_{\text{red}}^{\vec{}} \quad (\text{A.79})$$

$$+ \frac{1}{(\sum_{i=1}^s s_i)^2} \sum_{i=1}^s \sum_{j=1}^s s_i s_j d_{i,\text{red}}^{\vec{}} \cdot d_{j,\text{red}}^{\vec{}} + x_{m+1} y_{m+1} - x_{m+1} \frac{e}{\sum_{i=1}^s s_i} - y_{m+1} \frac{e}{\sum_{i=1}^s s_i} \quad (\text{A.80})$$

$$+ \frac{e^2}{(\sum_{i=1}^s s_i)^2} \quad (\text{A.81})$$

We get nonlinear solutions by using the following way of further kernelization

$$K(\vec{x}, \vec{y}) = K_o(x_{\text{red}}^{\vec{}}, y_{\text{red}}^{\vec{}}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(d_{i,\text{red}}^{\vec{}}, x_{\text{red}}^{\vec{}}) \quad (\text{A.82})$$

$$- \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(d_{i,\text{red}}^{\vec{}}, y_{\text{red}}^{\vec{}}) + \frac{1}{(\sum_{i=1}^s s_i)^2} \sum_{i=1}^s \sum_{j=1}^s s_i s_j K_o(d_{i,\text{red}}^{\vec{}}, d_{j,\text{red}}^{\vec{}}) \quad (\text{A.83})$$

$$+ x_{m+1}y_{m+1} - x_{m+1}\frac{e}{\sum_{i=1}^s s_i} - y_{m+1}\frac{e}{\sum_{i=1}^s s_i} + \frac{e^2}{(\sum_{i=1}^s s_i)^2} . \quad (\text{A.84})$$

This kernel is used only for solving the optimization problem. So solving  $\delta$ -SVR with the additional constraint leads to SVC optimization problem without the offset and with a special kernel presented above. Now we derive the solution

$$h(\vec{x}) = \vec{w}_c \cdot \vec{x} + b_c = \sum_{j=1}^{2n} \alpha_j y_c^j \left( x_{j,\text{red}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i d_{i,\text{red}} \right) \cdot \vec{x}_{\text{red}} \quad (\text{A.85})$$

$$+ \sum_{j=1}^{2n} \alpha_j y_c^j \left( x_j^{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) x_{m+1} + b_c \quad (\text{A.86})$$

with internal kernels

$$h(\vec{x}) = \sum_{j=1}^{2n} \alpha_j y_c^j \left( K_o(x_{j,\text{red}}, x_{\text{red}}) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i K_o(d_{i,\text{red}}, x_{\text{red}}) \right) \quad (\text{A.87})$$

$$+ \sum_{j=1}^{2n} \alpha_j y_c^j \left( x_j^{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) x_{m+1} + b_c \quad (\text{A.88})$$

where  $b_c$  is computed as

$$b_c = \frac{1}{\sum_{i=1}^s s_i} \left( -e \sum_{j=1}^{2n} \alpha_j y_c^j \left( x_j^{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) \right) \quad (\text{A.89})$$

$$- \frac{1}{\sum_{i=1}^s s_i} \left( \sum_{i=1}^s s_i \sum_{j=1}^{2n} \alpha_j y_c^j \left( x_{j,\text{red}} - \frac{1}{\sum_{i=1}^s s_i} \sum_{k=1}^s s_k d_{k,\text{red}} \right) \cdot d_{i,\text{red}} \right) \quad (\text{A.90})$$

$$= - \frac{1}{\sum_{i=1}^s s_i} e \sum_{j=1}^{2n} \alpha_j y_c^j \left( x_j^{m+1} - \frac{e}{\sum_{i=1}^s s_i} \right) - \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \sum_{j=1}^{2n} \alpha_j y_c^j K_o(x_{j,\text{red}}, d_{i,\text{red}}) \quad (\text{A.91})$$

$$+ \frac{1}{\sum_{i=1}^s s_i} \sum_{i=1}^s s_i \sum_{j=1}^{2n} \alpha_j y_c^j \frac{1}{\sum_{i=1}^s s_i} \sum_{k=1}^s s_k K_o(d_{k,\text{red}}, d_{i,\text{red}}) \quad (\text{A.92})$$

# References

- [1] Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: a library for support vector machines*, 2001. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>. 15, 29
- [2] Yuh jye Lee and Olvi L. Mangasarian. Rsvm: Reduced support vector machines. In *Data Mining Institute, Computer Sciences Department, University of Wisconsin*, pages 00–07, 2001. 27
- [3] Masayuki Karasuyama, Ichiro Takeuchi, and Ryohei Nakano. Reducing svr support vectors by using backward deletion. In *KES '08: Proceedings of the 12th international conference on Knowledge-Based Intelligent Information and Engineering Systems, Part III*, pages 76–83, Berlin, Heidelberg, 2008. Springer-Verlag. ISBN 978-3-540-85566-8. 27
- [4] S. Sathiya Keerthi, Olivier Chapelle, and Dennis DeCoste. Building support vector machines with reduced classifier complexity. *J. Mach. Learn. Res.*, 7:1493–1515, December 2006. ISSN 1532-4435. 27
- [5] LibSVM data sets. Libsvm data sets. [www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/](http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/), 06 2011. 29
- [6] Olvi L. Mangasarian, Jude W. Shavlik, and Edward W. Wild. Knowledge-based kernel approximation. *Journal of Machine Learning Research*, 5:1127–1141, 2004. 27
- [7] Kuan ming Lin and Chih jen Lin. A study on reduced support vector machines. *IEEE TRANSACTIONS ON NEURAL NETWORKS*, 14:1449–1459, 2003. 27
- [8] Marcin Orchel. Incorporating detractors into svm classification. In Krzysztof Cyran, Stanislaw Kozielski, James Peters, Urszula Stańczyk, and Alicja Wakulicz-Deja, editors, *Man-Machine Interactions*, volume 59 of *Advances in Intelligent and Soft Computing*, pages 361–369. Springer Berlin / Heidelberg, 2009. ISBN 978-3-642-00562-6. doi: 10.1007/978-3-642-00563-3\_38. URL [http://dx.doi.org/10.1007/978-3-642-00563-3\\_38](http://dx.doi.org/10.1007/978-3-642-00563-3_38). 11

- [9] Marcin Orchel. Incorporating a priori knowledge from detractor points into support vector classification. In Andrej Dobnikar, Uroš Lotric, and Branko Šter, editors, *Adaptive and Natural Computing Algorithms*, volume 6594 of *Lecture Notes in Computer Science*, pages 332–341. Springer Berlin / Heidelberg, 2011. ISBN 978-3-642-20266-7. doi: 10.1007/978-3-642-20267-4\_35. URL [http://dx.doi.org/10.1007/978-3-642-20267-4\\_35](http://dx.doi.org/10.1007/978-3-642-20267-4_35). 11, 27
- [10] Marcin Orchel. Regression based on support vector classification. In Andrej Dobnikar, Uroš Lotric, and Branko Šter, editors, *Adaptive and Natural Computing Algorithms*, volume 6594 of *Lecture Notes in Computer Science*, pages 353–362. Springer Berlin / Heidelberg, 2011. ISBN 978-3-642-20266-7. doi: 10.1007/978-3-642-20267-4\_37. URL [http://dx.doi.org/10.1007/978-3-642-20267-4\\_37](http://dx.doi.org/10.1007/978-3-642-20267-4_37). 9
- [11] Marcin Orchel. Support vector regression as a classification problem with a priori knowledge in the form of detractors. In Tadeusz Czachorski, Stanislaw Kozielski, and Urszula Stańczyk, editors, *Man-Machine Interactions 2*, volume 103 of *Advances in Intelligent and Soft Computing*, pages 353–362. Springer Berlin / Heidelberg, 2011. ISBN 978-3-642-23168-1. doi: 10.1007/978-3-642-23169-8\_38. URL [http://dx.doi.org/10.1007/978-3-642-23169-8\\_38](http://dx.doi.org/10.1007/978-3-642-23169-8_38). 11, 15, 27
- [12] John C. Platt. *Fast training of support vector machines using sequential minimal optimization*, pages 185–208. MIT Press, Cambridge, MA, USA, 1999. ISBN 0-262-19416-3. 14
- [13] Ingo Steinwart, Don R. Hush, and Clint Scovel. Training svms without offset. *Journal of Machine Learning Research*, 12:141–202, 2011. 5
- [14] Francis Eng Hock Tay and Lijuan Cao. Modified support vector machines in financial time series forecasting. *Neurocomputing*, 48(1-4):847–861, 2002. 7
- [15] Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley-Interscience, September 1998. ISBN 0471030031. 7
- [16] Mingrui Wu, Bernhard Schölkopf, and Gökhan Bakır. A direct method for building sparse kernel learning algorithms. *JOURNAL OF MACHINE LEARNING RESEARCH*, 7:603–624, 2006. 27
- [17] Xiaoyun Wu and Rohini Srihari. Incorporating prior knowledge with weighted margin support vector machines. In *KDD '04: Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 326–333, New York, NY, USA, 2004. ACM. ISBN 1-58113-888-1. 11

# Notation and Symbols

## Miscellaneous

- $|A|$  the cardinality of a finite set  $A$ , i.e., the number of elements in the set  $A$ .
- Dot product of two vectors, sometimes it is written with additional parentheses, for example for two vectors:  $\vec{u}$  and  $\vec{v}$ , the dot product is  $\vec{u} \cdot \vec{v}$  or  $(\vec{u} \cdot \vec{v})$ .
- $\vec{v} \geq \vec{w}$  For two  $n$  dimensional vectors  $\vec{v}$  and  $\vec{w}$ , it means that for all  $i = 1 \dots n$   $v_i \geq w_i$ .
- $\vec{v} \gg \vec{w}$  For two  $n$  dimensional vectors  $\vec{v}$  and  $\vec{w}$ , it means that for all  $i = 1 \dots n$   $v_i > w_i$ .
- $\rho(A)$  the rank of a matrix  $A$ .
- $w_{\mathbf{r}}^i$  When a vector has an index in the subscript, the coefficient index is placed in the superscript, the example means the  $i$ -th coefficient of the  $\vec{w}_{\mathbf{r}}$ .

## Optimization theory

\* an asterisk as a superscript in optimization theory denotes a solution of the optimization problem.

# Abbreviations

**$\delta$ -SVR**  $\delta$  support vector regression.

**$\nu$ -SVC**  $\nu$  support vector classification.

**$\varepsilon$ -SVR**  $\varepsilon$ -insensitive support vector regression.

**$\varphi$ -SVC**  $\varphi$  support vector classification.

**C-SVC** C support vector classification.

**DJIA** Dow Jones Industrial Average.

**KKT** Karush-Kuhn-Tucker.

**RBF** radial basis function.

**RSVM** reduced support vector machines.

**SMO** sequential minimal optimization.

**SVC** support vector classification.

**SVM** support vector machines.

**VC** Vapnik-Chervonenkis.