



## SOLUTION

### General Information:

Exercises (1 SWS): Tue 12:15 – 13:45 (0.154-115) and Fri 08:15 – 09:45 (0.151-115)

Certificate: Oral exam at the end of the semester

Contact: peter.fischer@fau.de, shiyang.hu@fau.de

## Support Vector Regression

**Exercise 1** In the lecture, you learn how an SVM can be used for classification. In this exercise, we consider *Support Vector Regression* (SVR). Let  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$  be a set of observations. The task for regression is to predict  $y_i$  from  $\mathbf{x}_i$  according to the linear regression function:

$$F(\mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{x} + \alpha_0, \quad (1)$$

for a weight vector  $\boldsymbol{\alpha} \in \mathbb{R}^d$  and the bias  $\alpha_0 \in \mathbb{R}$ . The intuition behind SVR is to penalize only deviations that are larger than  $\epsilon$ .

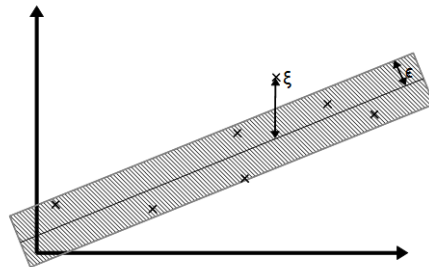


Figure 1:  $\epsilon$ -tube of the SVR

The primal optimization problem for SVR is given by the following inequality-constraint minimization:

$$\begin{aligned} \boldsymbol{\alpha}^* &= \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \frac{1}{2} \|\boldsymbol{\alpha}\|^2 + C \sum_i (\xi_i + \hat{\xi}_i), \text{ s.t.} \\ y_i &\leq (\boldsymbol{\alpha}^T \mathbf{x}_i + \alpha_0) + \epsilon + \xi_i \\ y_i &\geq (\boldsymbol{\alpha}^T \mathbf{x}_i + \alpha_0) - \epsilon - \hat{\xi}_i \\ \xi_i, \hat{\xi}_i &\geq 0 \end{aligned}$$

Here,  $\xi_i, \hat{\xi}_i$  are slack variables (see also SVM classification) and  $\epsilon$  specifies uncertainty of the regression function.

- (a) Write down the Lagrangian  $L$  of the primal optimization problem using Lagrange multipliers  $\lambda_i, \hat{\lambda}_i, \mu_i, \hat{\mu}_i$ .  
 Hint: bring the constraints to the standard form  $f_i(\mathbf{x}) \leq 0$   
 First, reformulate the inequality constraints to standard form  $f_i(\boldsymbol{\alpha}) \leq 0$

$$\begin{aligned} y_i - (\boldsymbol{\alpha}^T \mathbf{x}_i + \alpha_0) - \epsilon - \xi_i &\leq 0 \rightarrow f_i(\boldsymbol{\alpha}) \leq 0 \\ -y_i + (\boldsymbol{\alpha}^T \mathbf{x}_i + \alpha_0) - \epsilon - \hat{\xi}_i &\leq 0 \rightarrow \hat{f}_i(\boldsymbol{\alpha}) \leq 0 \\ -\xi_i &\leq 0 \\ -\hat{\xi}_i &\leq 0 \end{aligned}$$

no equality constraints

Summarize into one equation:

$$\begin{aligned} L(\boldsymbol{\alpha}, \alpha_0, \boldsymbol{\xi}, \hat{\boldsymbol{\xi}}, \boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}, \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}) = & \frac{1}{2} \|\boldsymbol{\alpha}\|^2 + C \sum_i (\xi_i + \hat{\xi}_i) + \\ & \sum_i \lambda_i (y_i - \boldsymbol{\alpha}^T \mathbf{x}_i - \alpha_0 - \epsilon - \xi_i) + \\ & \sum_i \hat{\lambda}_i (-y_i + \boldsymbol{\alpha}^T \mathbf{x}_i + \alpha_0 - \epsilon - \hat{\xi}_i) + \\ & \sum_i (-\mu_i \xi_i - \hat{\mu}_i \hat{\xi}_i) \end{aligned}$$

- (b) Write down the Karush-Kuhn-Tucker (KKT) conditions for the primal optimization problem given above.

(a) Primal constraints: see above

(b) Dual constraints:  $\forall i : \lambda_i, \hat{\lambda}_i, \mu_i, \hat{\mu}_i \geq 0$

(c) Complementary slackness:

$$\forall i : \lambda_i f_i(\boldsymbol{\alpha}) = 0, \hat{\lambda}_i \hat{f}_i(\boldsymbol{\alpha}) = 0, \mu_i \xi_i = 0, \hat{\mu}_i \hat{\xi}_i = 0$$

(d) Gradient of Lagrangian is zero:  $\nabla L(\boldsymbol{\alpha}, \alpha_0, \boldsymbol{\xi}, \hat{\boldsymbol{\xi}}, \boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}, \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}) = 0$

The gradient of the Lagrangian is:

$$\frac{\partial L}{\partial \boldsymbol{\alpha}} = \boldsymbol{\alpha} - \sum_i \lambda_i \mathbf{x}_i + \sum_i \hat{\lambda}_i \mathbf{x}_i = 0 \Rightarrow \boldsymbol{\alpha} = \sum_i (\lambda_i - \hat{\lambda}_i) \mathbf{x}_i \quad (2)$$

$$\frac{\partial L}{\partial \alpha_0} = -\sum_i \lambda_i + \sum_i \hat{\lambda}_i = 0 \Rightarrow \sum_i (\lambda_i - \hat{\lambda}_i) = 0 \quad (3)$$

$$\frac{\partial L}{\partial \xi_i} = C - \lambda_i - \mu_i = 0 \Rightarrow \lambda_i + \mu_i = C \quad (4)$$

$$\frac{\partial L}{\partial \hat{\xi}_i} = C - \hat{\lambda}_i - \hat{\mu}_i = 0 \Rightarrow \hat{\lambda}_i + \hat{\mu}_i = C \quad (5)$$

- (c) Derive the dual optimization problem. To derive the dual optimization problem, you have to eliminate  $\boldsymbol{\alpha}$ ,  $\boldsymbol{\xi}$ , and  $\hat{\boldsymbol{\xi}}$  from  $L$  using the gradient of  $L$ .

Preliminary solution:

$$L(\boldsymbol{\alpha}, \alpha_0, \boldsymbol{\xi}, \hat{\boldsymbol{\xi}}, \boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}, \boldsymbol{\mu}, \hat{\boldsymbol{\mu}}) = \frac{1}{2} \|\boldsymbol{\alpha}\|^2 + C \sum_i (\xi_i + \hat{\xi}_i) + \sum_i \left( -\mu_i \xi_i - \hat{\mu}_i \hat{\xi}_i \right) +$$

$$\sum_i \lambda_i (y_i - \boldsymbol{\alpha}^T \mathbf{x}_i - \alpha_0 - \epsilon - \xi_i) +$$

$$\sum_i \hat{\lambda}_i \left( -y_i + \boldsymbol{\alpha}^T \mathbf{x}_i + \alpha_0 - \epsilon - \hat{\xi}_i \right)$$

For  $\boldsymbol{\alpha}$ , a direct replacement was found in equation (2).  $\xi_i$  are eliminated after collecting all the terms in the Lagrangian, because  $C \sum_i \xi_i - \sum_i (\lambda_i + \mu_i) \xi_i = 0$  using (4). The same is true for  $\hat{\xi}_i$ .  $\alpha_0$  is eliminated, because  $\sum_i (\hat{\lambda}_i - \lambda_i) \alpha_0 = 0$  using (3). The resulting dual optimization problem is:

$$\tilde{L}(\boldsymbol{\lambda}, \hat{\boldsymbol{\lambda}}) = \overbrace{\frac{1}{2} \boldsymbol{\alpha}^T \boldsymbol{\alpha}}^{\frac{1}{2} \boldsymbol{\alpha}^T \boldsymbol{\alpha}} = \frac{1}{2} \sum_i \sum_j (\lambda_i - \hat{\lambda}_i) (\lambda_j - \hat{\lambda}_j) \mathbf{x}_i^T \mathbf{x}_j -$$

$$\sum_i \lambda_i \left[ \sum_j (\lambda_j - \hat{\lambda}_j) \mathbf{x}_j^T \right] \mathbf{x}_i + \sum_i \hat{\lambda}_i \left[ \sum_j (\lambda_j - \hat{\lambda}_j) \mathbf{x}_j^T \right] \mathbf{x}_i -$$

$$\epsilon \sum_i (\lambda_i + \hat{\lambda}_i) + \sum_i (\lambda_i - \hat{\lambda}_i) y_i$$

$$= -\frac{1}{2} \sum_i \sum_j (\lambda_i - \hat{\lambda}_i) (\lambda_j - \hat{\lambda}_j) \mathbf{x}_i^T \mathbf{x}_j -$$

$$\epsilon \sum_i (\lambda_i + \hat{\lambda}_i) + \sum_i (\lambda_i - \hat{\lambda}_i) y_i \quad (6)$$

The dual optimization is also constrained.  $\lambda_i$  and  $\hat{\lambda}_i$  are Lagrangian multipliers and therefore non-negative. In addition, they are constrained by (4) and (5). Equation (3) is also a constraint.

$$0 \leq \lambda_i, \hat{\lambda}_i \leq C \quad (7)$$

$$\sum_i (\lambda_i - \hat{\lambda}_i) = 0 \quad (8)$$

(d) Which property must be fulfilled for *support vectors* in SVR?

Hint: replace  $\boldsymbol{\alpha}$  in Equation (1).

Support vectors are the only ones necessary to compute the regression. Therefore, they are included in (1).

$$F(\mathbf{x}) = \boldsymbol{\alpha}^T \mathbf{x} + \alpha_0$$

$$= \sum_{i=1}^N (\lambda_i - \hat{\lambda}_i) \mathbf{x}_i^T \mathbf{x} + \alpha_0$$

For support vectors, the term in the sum is not 0  $\Rightarrow \lambda_i - \hat{\lambda}_i \neq 0$ . Due to the complementary slackness condition, this is the case only for points that are on the boundary or outside of the  $\epsilon$ -tube. This is a sparse solution, because most points are close to the estimated output of the regression.