

## Support Vector Regression

### Exercise T10.1: Regression with SVM

(tutorial)

In regression problems, we are given a training data set

$$\{(\underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)})\}, \quad \alpha \in \{1, \dots, p\}, \quad \underline{\mathbf{x}} \in \mathbb{R}^N, \quad y_T \in \mathbb{R},$$

and want to fit the linear regression function

$$y(\underline{\mathbf{x}}) = \underline{\mathbf{w}}^\top \underline{\mathbf{x}} + b.$$

- (a) What is the  $\varepsilon$ -insensitive cost function for regression?
- (b) Derive the primal problem of the  $\varepsilon$ -support vector regression ( $\varepsilon$ -SVR).
- (c) The optimal  $\varepsilon$ -parameter depends linearly on the noise level in the data, which is unknown. Derive the primal problem for the  $\nu$ -SVR, which adjusts  $\varepsilon$  as a primal parameter.
- (d) Derive the Lagrangian of the  $\nu$ -SVR.

Solution:

- (a) The  $\varepsilon$ -insensitive cost function

$$e(\underline{\mathbf{x}}, y_T) = \max(0, |y(\underline{\mathbf{x}}) - y_T| - \varepsilon).$$

- (b)  $\varepsilon$ -SVR has the following primal problem:

$$\min_{\underline{\mathbf{w}}, b, \varphi_\alpha, \varphi_\alpha^*} \frac{1}{2} \|\underline{\mathbf{w}}\|^2 + C \left( \frac{1}{p} \sum_{\alpha=1}^p (\varphi_\alpha + \varphi_\alpha^*) \right)$$

s.t.

$$\begin{aligned} (\underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} + b) - y_T^{(\alpha)} &\leq \varepsilon + \varphi_\alpha \\ y_T^{(\alpha)} - (\underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} + b) &\leq \varepsilon + \varphi_\alpha^* \\ \varphi_\alpha, \varphi_\alpha^* &\geq 0, \end{aligned}$$

where  $\varphi_\alpha, \varphi_\alpha^*$  are slack variables,  $\underline{\mathbf{w}}, b, \varphi_\alpha, \varphi_\alpha^*$  are called the primal variables and where the constant  $C > 0$  determines the trade-off between the 'flatness' of  $y$  and the amount up to which deviations larger than  $\varepsilon$  are tolerated.

- (c) The optimal  $\varepsilon$ -parameter linearly depends on the noise level in the data, which is unknown. There exists, however, a method to automatically adjust  $\varepsilon$ , and at the same time have a predetermined fraction of support vectors: The so-called  $\nu$ -SVR allows the  $\varepsilon$ -tube width to automatically adapt to the data. In contrast to  $\varepsilon$ -support vector regression,  $\varepsilon$  becomes a variable of the primal optimization problem, which now includes an extra term which

attempts to minimize  $\varepsilon$ . Introducing a fixed parameter  $\nu \geq 0$  (which was shown to provide a lower bound on the fraction of support vectors), the **primal problem** of the  $\nu$ -SVR is:

$$\min_{\underline{\mathbf{w}}, b, \varphi_\alpha, \varphi_\alpha^*, \varepsilon} \frac{1}{2} \|\underline{\mathbf{w}}\|^2 + C \left( \nu \varepsilon + \frac{1}{p} \sum_{\alpha=1}^p (\varphi_\alpha + \varphi_\alpha^*) \right)$$

s.t.  $\forall \alpha \in \{1, \dots, p\} :$

$$\begin{aligned} (\underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} + b) - y_T^{(\alpha)} &\leq \varepsilon + \varphi_\alpha \\ y_T^{(\alpha)} - (\underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} + b) &\leq \varepsilon + \varphi_\alpha^* \\ \varphi_\alpha, \varphi_\alpha^*, \varepsilon &\geq 0 \end{aligned}$$

where  $\varphi_\alpha, \varphi_\alpha^*$  are slack variables and  $\underline{\mathbf{w}}, b, \varphi_\alpha, \varphi_\alpha^*, \varepsilon$  are the primal variables.

(d) This corresponds to the following Lagrangian

$$\begin{aligned} L(\underbrace{\underline{\mathbf{w}}, b, \{\varphi_\alpha\}, \{\varphi_\alpha^*\}, \varepsilon}_{\text{primal variables}}, \underbrace{\{\lambda_\alpha\}, \{\lambda_\alpha^*\}, \{\eta_\alpha\}, \{\eta_\alpha^*\}, \delta}_{\text{dual variables (Lagrange multipliers)}}) \\ = \frac{1}{2} \|\underline{\mathbf{w}}\|^2 + C \left( \nu \varepsilon + \frac{1}{p} \sum_{\alpha=1}^p (\varphi_\alpha + \varphi_\alpha^*) \right) \\ - \sum_{\alpha=1}^p \lambda_\alpha \{ \varphi_\alpha + \varepsilon + y_T^{(\alpha)} - \underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} - b \} \\ - \sum_{\alpha=1}^p \lambda_\alpha^* \{ \varphi_\alpha^* + \varepsilon - y_T^{(\alpha)} + \underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} + b \} \\ - \sum_{\alpha=1}^p \eta_\alpha \varphi_\alpha - \sum_{\alpha=1}^p \eta_\alpha^* \varphi_\alpha^* - \delta \varepsilon \end{aligned}$$

The Lagrange multipliers  $\lambda_\alpha, \lambda_\alpha^*, \eta_\alpha, \eta_\alpha^*, \delta$  must all be  $\geq 0$  (since they correspond to inequality constraints) and are called the dual variables.

### Exercise H10.1: The dual problem of the $\nu$ -SVR

(homework, 5 points)

In this exercise you will derive the dual problem of the  $\nu$ -SVR.

- (a) (2 points) Calculate the derivatives of the Lagrangian with respect to the primal variables.
- (b) (3 points) By setting the derivatives from (a) to zero and using the results to eliminate the primal variables from the Lagrangian show that the dual problem takes the following form:

$$\max_{\lambda_\alpha, \lambda_\alpha^*} -\frac{1}{2} \sum_{\alpha, \beta=1}^p (\lambda_\alpha^* - \lambda_\alpha)(\lambda_\beta^* - \lambda_\beta) (\underline{\mathbf{x}}^{(\alpha)})^\top \underline{\mathbf{x}}^{(\beta)} + \sum_{\alpha=1}^p (\lambda_\alpha^* - \lambda_\alpha) y_T^{(\alpha)}$$

s.t.  $\forall \alpha \in \{1, \dots, p\} :$

$$0 \leq \lambda_\alpha \leq \frac{C}{p}, \quad 0 \leq \lambda_\alpha^* \leq \frac{C}{p}, \quad \sum_{\alpha=1}^p (\lambda_\alpha - \lambda_\alpha^*) = 0, \quad \sum_{\alpha=1}^p (\lambda_\alpha + \lambda_\alpha^*) \leq \nu C.$$

**Exercise H10.2: Regression with the  $\nu$ -SVR (homework, 5 points)**

In this exercise you will apply  $\nu$ -SVR from a software package of your choice (e.g. `scikit-learn`, `libsvm`) to the same dataset used in exercise sheet 5. The training set `TrainingRidge.csv` and the validation set `ValidationRidge.csv` can be found on ISIS. Do **not** center, whiten or expand the data before training (otherwise the proposed hyperparameter ranges become inadequate).

- (a) (2 points) Train the  $\nu$ -SVR on the training set with the standard parameters of your library (“out of the box”).

Deliverables:

1. Plot the model prediction for the validation set as an image plot (where colors represent the output values, the axes represent the two coordinates:  $x_1$  and  $x_2$ ). Add the data points from the training set by highlighting their locations (e.g. colored rectangles) in the same plot.
  2. Compute the mean squared error (MSE) between model prediction and true labels of the validation set. Make a second plot over  $x_1$  and  $x_2$  with a heat map of the MSE.
- (b) (2 points) Perform a 10-fold cross-validation with a  $\nu$ -SVR with parameters  $\nu = 0.5$  and  $C \in 2^i$ ,  $i \in \{-2, \dots, 12\}$ . Use a Gaussian RBF kernel with  $\gamma \in 2^j$ ,  $j \in \{-12, \dots, 0\}$ .

Deliverables:

Plot the resulting mean (test set) MSE over the folds as an image plot. Note that the RBF kernel is parametrized as in the previous sheet (with parameter  $\gamma$  instead of  $\sigma$ ).

- (c) (1 point) Extract the best parameter combination  $C$  and  $\gamma$ . Use the entire training set to train a new  $\nu$ -SVR with these parameters.

Deliverables:

1. Plot the model prediction for the validation set as an image plot. Compare the plot with the true labels and the results from (a).
2. Visualize the mean squared error for the validation set as a heat map for comparison.

**Total 10 points.**