

Support Vector Machines

The winter holidays start next week.

Homework of this exercise sheet is due next year, on 05.01.2017!

Exercise T8.1: Structural Risk Minimization

(tutorial)

- (a) Discuss the concept of the *margin* for the linear connectionist neuron: What is the effect of a small vs. a big margin on generalization?
- (b) Write down and explain the *primal optimization problem* of model selection through structural risk minimization (SRM).
- (c) Write down the Lagrangian of the primal problem and explain the intuition behind the theorem of Kuhn and Tucker. Why can we expect sparse dual variables?
- (d) Discuss SVM classification of non-separable classes. How can this be regularized? Write down the primal problem of the C-SVM.
- (e) What is the kernel-trick and how can we exploit it?

Exercise H8.1: Deriving the C-SVM optimization problem (homework, 3 points)

- (a) (1 point) Linear connectionist neurons have a degree of freedom that is not used in classification. By setting the constraint

$$\min_{\alpha=1,\dots,p} \left| \underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} + b \right| \stackrel{!}{=} 1$$

this degree is eliminated. Show that under this constraint the Euclidean distance $d(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{w}}, b)$ of sample $\underline{\mathbf{x}}^{(\alpha)}$ to the closest point of the decision boundary $\{x|y(x) = 0\}$ is bounded by

$$d(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{w}}, b) \geq \frac{1}{\|\underline{\mathbf{w}}\|}, \quad \forall \alpha \in \{1, \dots, p\}.$$

- (b) (2 points) Write down the Lagrangian of the primal optimization problem of the C-SVM and derive the dual optimization problem of the C-SVM:

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

with $0 \leq \lambda_\alpha \leq \frac{C}{p}$ and $\sum_{\alpha=1}^p \lambda_\alpha y_T^{(\alpha)} = 0.$

Exercise H8.2: C-SVM with standard parameters (homework, 3 points)

In this exercise, we use C-SVMs to solve the “XOR”-classification problem from exercise sheet 6. To this end (1) first create a *training set* of 80 data as described in exercise H6.1 and (2) create a *test set* of 80 data from the same distribution.

You can use existing software: `libsvm`¹ implements optimization routines (Matlab & Python) for SVMs. Alternatively, you can use the corresponding `scikit.learn` class². For R, the package `e1071` implements SVM-optimization.

- Download, install, and familiarize yourself with `LIBSVM` or one of the other packages.
- Read the *Practical Guide to Support Vector Classification*³ especially section 3.2 on *Cross-Validation*.

Next, use your chosen SVM implementation to train a C-SVM with RBF kernel and the software’s standard parameters. Classify the test data and report the classification error quantified by the 0/1 loss function (percentage of wrong predictions). Visualize the results as in exercise H6.2: plot the training patterns and the decision boundary (e.g. with a contour plot) in input space.

Exercise H8.3: C-SVM parameter optimization (homework, 4 points)

- (2 points) Use cross-validation and grid-search to determine good values for C and the kernel parameter γ . Follow the procedure described in the *guide*: Define the grid using exponentially growing sequences of C and γ , e.g. $C \in \{2^{-6}, 2^{-4}, \dots, 2^{10}\}$, $\gamma \in \{2^{-5}, 2^{-3}, \dots, 2^9\}$. Make sure you only use the training data in this step. Plot the mean training-set classification rate and cross-validation performance as a function of C and γ (e.g. using contour plots as in figure 2 of the *guide*).
- (1 point) Find the best combination of C and γ and train the RBF C-SVM on the *entire* training data, this time using these “optimal” parameters. Plot the results in the same way as in exercise H8.2.
- (1 point) Compare the results with those obtained in H8.2, both in terms of statistics (e.g. classification performance, number of support vectors) and visually (e.g. signs of over- and under-fitting). What happens when you divide C or γ by 4?

Total 10 points.

¹ <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

² <http://tinyurl.com/lrpxw9k>

³ <http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>