

## Machine Learning Exercise Sheet 09

## SVM and Kernels

## In-Class

## 1 SVM

$$y_i = f(x_i) + \xi_i \quad y_i (w^T x_i + b) \geq 1 - \xi_i \quad \text{deal with noise}$$

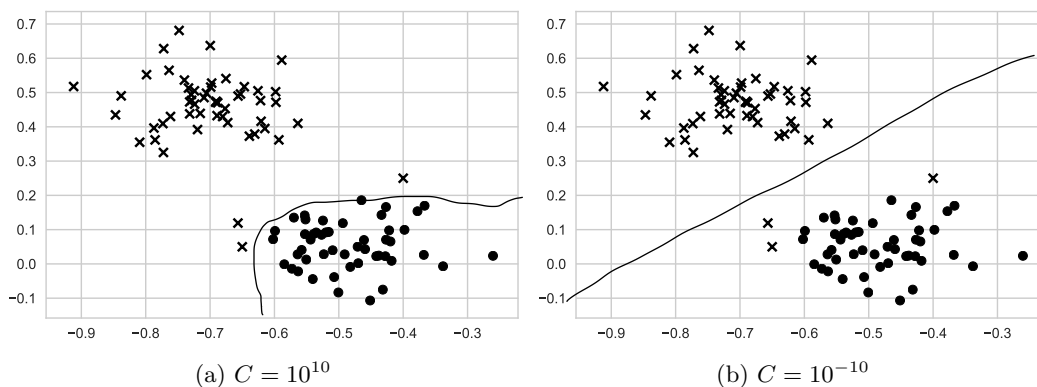
**Problem 1:** What is the connection between soft-margin SVM and logistic regression?

**Problem 2:** Consider a soft-margin SVM fitted to a linearly separable dataset  $\mathcal{D}$  using the Hinge loss formulation of the optimization task.

$$\text{learn of} \quad \text{minimize}_{w,b} \quad \frac{1}{2} w^T w + C \sum_{i=1}^N \max(0, 1 - y_i(w^T x_i + b))$$

- a) *Yes, linearly separable dataset* Is it guaranteed that all training samples in  $\mathcal{D}$  will be assigned the correct label by the model?
- b) Prove that if for some  $C_0 \geq 0$  the resulting model classifies all training samples correctly and all samples lie outside of the margin then it will also be the case if we train the model with any larger  $C > C_0$ .

**Problem 3:** Sketch the decision boundary of an SVM with a quadratic kernel (polynomial with degree 2) for the data in the figure below, for two specified values of the penalty parameter  $C$ . The two classes are denoted as  $\bullet$ 's and  $\times$ 's.



1. Soft margin.

$$\begin{aligned} \min \quad & f_0(w, b, \xi) = \frac{1}{2} w^T w + c \sum_{i=1}^N \xi_i \\ \text{sub to} \quad & y_i (w^T x_i + b) - 1 + \xi_i \geq 0 \quad i = 1, \dots, N \\ & \xi_i \geq 0 \quad i = 1, \dots, N \end{aligned}$$

Loss function

$$\bar{E}(w, b, c) = \frac{1}{2c} w^T w + \sum_{i=1}^N \text{Hinge}(y_i \cdot (w^T x_i + b))$$

Loss Regression

$$p(y_i = 1 | w, x) = \sigma(w^T x_i + b)$$

$$p(y_i = -1 | w, x) = \sigma(-(w^T x_i + b)) = 1 - \sigma(w^T x_i + b)$$

$$p(y | w, x) = \prod_{i=1}^N p(y_i | w, x_i) = \frac{1}{1 + e^{-(w^T x_i + b)}} \cdot \frac{1}{1 + e^{-(w^T x_i + b)}}$$

$$\begin{aligned} E(w, b) &= -\ln p(y | w, x) = - \sum_{i=1}^N \ln (1 + e^{-y_i (w^T x_i + b)}) \\ &= \sum_{i=1}^N \ln (1 + e^{-y_i (w^T x_i + b)}) \end{aligned}$$

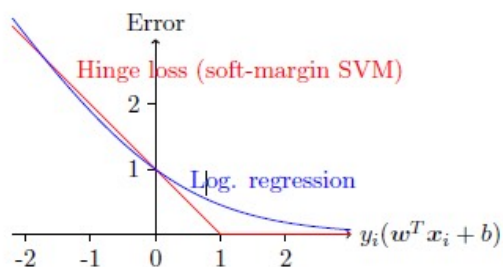
$$\textcircled{2} \quad \bar{E}(w, b, \lambda) = \lambda w^T w + \sum_{i=1}^N \ln (1 + e^{-y_i (w^T x_i + b)})$$

$$\textcircled{1} \quad \bar{E}(w, b, c) = \frac{1}{2c} w^T w + \sum_{i=1}^N \text{Hinge}(y_i \cdot (w^T x_i + b))$$

Hence, we see the close relationship between the soft-margin SVM (Eq. 1) and logistic regression (Eq.

2). While soft-margin SVM uses the hinge function for its loss, logistic regression uses  $\ln(1 + e^{-x})$ .

For better comparison, we can plot these two (with logistic regression rescaled by  $\frac{1}{\ln 2}$ ):



2  
a) No, because of misclassifying a ten point nearby by the decision boundary will lead to significantly increasing margin.

range  $C \rightarrow$  Hard SUM  $\rightarrow$  less behaviour

b) let  $h(w, b) = \sum_{i=1}^N \max(0, 1 - y_i (w^T x_i + b))$

when right classified  $h(w_b) = 0$

$$(w^*, b^*) = \arg \min_{w, b} \frac{1}{2} w^T w + c_0 h(w, b)$$

$$(v^*, d^*) = \arg \min_{v, d} \frac{1}{2} w^T w + c \cdot h(w, b)$$

$$\frac{1}{2} \|w^*\|^2 + c_0 h(w^*, b^*) \leq \frac{1}{2} \|v^*\|^2 + c_0 h(v^*, d^*)$$

$$\frac{1}{2} \|v^*\|^2 + C h(v^*, d^*) \leq \frac{1}{2} \|w^*\|^2 + C h(w^*, b^*)$$

$$(c_0 - c) (h(w^*, b^*) - h(v^*, d^*)) \leq 0$$

$$\therefore c > c_0 \quad \therefore h(w^*, b^*) \leq h(v^*, d^*)$$

That means the overall Hinge loss decreases for the optimal solution if we solve the SVM fitting problem with a larger value of  $C$ . Since we could separate the data perfectly (with no points within the margin) with the solution  $(w^*, b^*)$  for  $C_0$  we have  $h(w^*, b^*) = 0$ . As shown above the Hinge loss can not get worse for larger  $C$ , hence also  $h(v^*, d^*)$  has to be 0 and the solution  $(v^*, d^*)$  also corresponds to a decision boundary with no error on training data and no points lying inside the margin.

## 2 Kernels

**Problem 4:** Consider the Gaussian kernel

$$k_G(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|^2}{2\sigma^2}\right), \text{ with } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

$$\begin{aligned} 4. \textcircled{1} \quad k_G(\theta(\mathbf{x}_1), \theta(\mathbf{x}_2)) &= \exp\left(-\frac{\|\theta(\mathbf{x}_1) - \theta(\mathbf{x}_2)\|^2}{2\sigma^2}\right) \\ &= \exp\left(-\frac{\theta(\mathbf{x}_1)^T \theta(\mathbf{x}_1) - 2\theta(\mathbf{x}_1)^T \theta(\mathbf{x}_2) + \theta(\mathbf{x}_2)^T \theta(\mathbf{x}_2)}{2\sigma^2}\right) \\ &= \exp\left(-\frac{k(\mathbf{x}_1, \mathbf{x}_1) - 2k(\mathbf{x}_1, \mathbf{x}_2) + k(\mathbf{x}_2, \mathbf{x}_2)}{2\sigma^2}\right) \end{aligned}$$

- a) Suppose you have found a feature map  $\theta : \mathbb{R}^d \rightarrow \mathbb{R}^{d_1}$  that transforms your data into a feature space in which a SVM with a Gaussian kernel works well. However, computing  $\theta(\mathbf{x})$  is computationally expensive and luckily you discover an efficient method to compute the scalar product

$$k(\mathbf{x}_1, \mathbf{x}_2) = \theta(\mathbf{x}_1)^T \theta(\mathbf{x}_2)$$

in your feature space without having to compute  $\theta(\mathbf{x}_1)$  and  $\theta(\mathbf{x}_2)$  explicitly. Show how you can use the scalar product  $k(\mathbf{x}_1, \mathbf{x}_2)$  to efficiently compute the Gaussian kernel  $k_G(\theta(\mathbf{x}_1), \theta(\mathbf{x}_2))$  in your feature space.

- b) One of the nice things about kernels is that new kernels can be constructed out of already given ones. Use the five kernel construction rules from the lecture to prove that  $k_G$  is a kernel.

*Hint: Use the Taylor expansion of the exponential function to prove that  $\exp \circ k_1$  is a kernel if  $k_1$  is a kernel. Also, consider  $k_2(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2))$  with the linear kernel  $k_2(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^T \mathbf{x}_2$  and a feature map  $\phi$  with only one feature.*

- c) Can any finite set of points be linearly separated in the feature space of the Gaussian kernel if  $\sigma$  can be chosen freely?

**Problem 5:** Let  $\mathcal{M} = \bigcup_{n \in \mathbb{N}} \bigcup_{m \in \mathbb{N}} \mathbb{R}^{n \times m}$  denote the set of all real-valued matrices of arbitrary size. Prove that the function  $k : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$  is a valid kernel, where

$$k(\mathbf{X}, \mathbf{Y}) = \min(\text{rank}(\mathbf{X}), \text{rank}(\mathbf{Y})).$$

(construct a function

$$\phi(\mathbf{x}) = \{1, 1, 1, \dots, 0, 0, 0\}$$

$$\phi^T(\mathbf{x}) \phi(\mathbf{y}) = \sum_{j=1}^{\infty} \phi_j(\mathbf{x}) \phi_j(\mathbf{y}) = \min(\text{rank}(\mathbf{X}), \text{rank}(\mathbf{Y})) \in \mathbb{R}$$

4. ② prove  $\exp \circ k_1$  is a kernel if  $k(x_1, x_2)$  is a kernel

Taylor expansion:

$$\exp(k(x_1, x_2)) = 1 + \sum_{n=1}^{\infty} \frac{1}{n!} k(x_1, x_2)^n$$

$\nwarrow$  rule 3  $k_1(x_1, x_2) \cdot k_2(x_1, x_2)$   
 $\nwarrow$  rule 2  $c \cdot k_1(x_1, x_2)$   
 $\nwarrow$  rule 1  $k_1(x_1, x_2) + k_2(x_1, x_2)$   
 $\nwarrow$  rule 4  $k(\phi(x_1), \phi(x_2))$   
 $\phi(z) = 1$

$$\begin{aligned} \exp(k_G(x_1, x_2)) &= \exp\left(-\frac{\|x_1 - x_2\|^2}{2\sigma}\right) = \exp\left(-\frac{x_1^T x_1 - 2x_1^T x_2 + x_2^T x_2}{2\sigma}\right) \\ &= \exp\left(-\frac{x_1^T x_1}{2\sigma}\right) \cdot \exp\left(\frac{x_1^T x_2}{\sigma}\right) \exp\left(-\frac{x_2^T x_2}{2\sigma}\right) \\ &\quad \uparrow \text{use rule} \end{aligned}$$

Consider the last term first. The scalar product  $x_1^T x_2$  is the linear kernel and by rule 2 the product  $x_1^T x_2 / \sigma^2$  is a kernel because  $\sigma^2$  is positive. As proved above the term  $\exp(x_1^T x_2 / \sigma^2)$  is then a kernel as well.

The product of the first two terms  $\exp(-x_1^T x_1 / 2\sigma^2) \exp(-x_2^T x_2 / 2\sigma^2)$  is a kernel by rule 4 with  $k_3(x_1, x_2) = x_1 x_2$  and the feature map  $\phi(z) = \exp(-z^T z / 2\sigma^2)$ .

Finally, by rule 3 the product of the first two terms with the third term is a kernel.

$$\textcircled{3} \quad k_G(x_1, x_2, \sigma) \xrightarrow{\sigma=0} k(x_1, x_2) = \begin{cases} 1 & x_1 = x_2 \\ 0 & x_1 \neq x_2 \end{cases}$$

All training sample are correctly classified if  $y_i (w^T \phi_G(x_i) + b) > 0$

$$y_i \left( \sum_j y_j \alpha_j k_G(x_i, x_j, \sigma) + b \right) > 0 \xrightarrow{\sigma=0} y_i^2 \alpha_i + y_i b$$

for small  $\sigma$ . By choosing  $b = 0$  we see that the resulting condition is fulfilled for all training samples, since  $y_i^2 = 1$  for all  $i$  and we can simply set all  $\alpha_i > 0$  (note, that this means that every sample is a support vector in this case). Hence all finite sets of points can be linearly separated using the Gaussian kernel if the variance  $\sigma$  is chosen small enough.

**Homework**

Creed boundary decision = hyperplane  
 SVM want to maximize the margin

**3 SVM**

**Problem 6:** Explain the similarities and differences between the SVM and perceptron algorithms. How do they perform classification? In what way do they differ?

**Problem 7:** Recall that the dual function in the setting of the SVM training task (Slide 17) can be written as

$$g(\alpha) = \frac{1}{2} \alpha^T Q \alpha + \alpha^T \mathbf{1}_N.$$

- (a) Write down the matrix  $Q$  using the vector of labels  $\mathbf{y}$  and feature matrix  $\mathbf{X}$ . Denote the element-wise product between two matrices (in case you want to use it) by  $\odot$  (also known as Hadamard product or Schur product).
- (b) Prove that we can search for a *local* maximizer of  $g$  to find its *global* maximum (don't forget to prove properties of  $Q$  that you decide to use in this task).

**Problem 8:** Consider training a standard hard-margin SVM on a linearly separable training set of  $N$  samples. Let  $s$  denote the number of support vectors we would obtain if we would train on the entire dataset. Furthermore, let  $\varepsilon$  denote the leave-one-out cross validation (LOOCV) misclassification rate. Prove that the following relation holds:

$$\varepsilon \leq \frac{s}{N}.$$

**Problem 9:** Load the notebook `exercise_09_notebook.ipynb` from Moodle. Fill in the missing code and run the notebook. Convert the evaluated notebook to pdf and add it to the printout of your homework.

*Note: We suggest that you use Anaconda for installing Python and Jupyter, as well as for managing packages. We recommend that you use Python 3.*

*For more information on Jupyter notebooks and how to convert them to other formats, consult the Jupyter documentation and nbconvert documentation.*

**4 Kernels**

**Problem 10:** Show that for  $N \in \mathbb{N}$  and  $a_i \geq 0$  for  $i = 0, \dots, N$  the following function  $k$  is a valid kernel.

$$k(\mathbf{x}_1, \mathbf{x}_2) = \sum_{i=1}^N a_i (\mathbf{x}_1^T \mathbf{x}_2)^i + a_0, \text{ with } \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

**Problem 11:** Find the feature transformation  $\phi(x)$  corresponding to the kernel

$$k(x_1, x_2) = \frac{1}{1 - x_1 x_2}, \text{ with } x_1, x_2 \in (0, 1).$$

*Hint: Consider an infinite-dimensional feature space and infinite series.*

10

$$k(x_1, x_2) = \sum_{i=1}^N a_i (x_1^T x_2)^2 + a_0, \text{ with } x_1, x_2 \in \mathbb{R}^d.$$

Annotations:

- rule 3 points to  $(x_1^T x_2)^2$
- rule 3 points to  $a_i$
- rule 1 points to  $a_0$
- rule 2 points to  $a_i \geq 0$
- rule points to the summation symbol  $\sum$

$k_1 = x_1^T x_2$  is kernel because of scalar product of vector

$k_c(x_1, x_2)$  is kernel, because we can define feature map  $\phi(x) = \sqrt{a_0}$   
 $\phi^T(x_1) \phi(x_2) = a_0$

---

11 geometric transformation

$$k(x_1, x_2) = \frac{1}{1 - x_1 x_2} = \sum_{i=0}^{\infty} x_1^i x_2^i = \phi^T(x_1) \phi(x_2)$$

$$\phi(x) = (1, x, x^2, \dots, x^n)^T$$

$$\begin{aligned}
 g(\alpha) &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i y_j x_i^T x_j \alpha_j \\
 Q &= - y_i y_j x_i^T x_j \\
 A &= y_i y_j \quad B = x_i^T x_j \quad y \in (N \times 1) \\
 A &= \begin{pmatrix} y_1 y_1 & y_1 y_2 & \dots & y_1 y_N \\ \vdots & \vdots & \ddots & \vdots \\ y_N y_1 & y_N y_2 & \dots & y_N y_N \end{pmatrix} = y y^T \quad B = \begin{pmatrix} x_1^T x_1 & \dots & x_1^T x_N \\ \vdots & \ddots & \vdots \\ x_N^T x_1 & \dots & x_N^T x_N \end{pmatrix} = X \cdot X^T \\
 Q &= - y y^T \odot X \cdot X^T
 \end{aligned}$$

$$\begin{aligned}
 \text{b)} \quad x^T (A \odot B) x &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i y_j x_i^T x_j \alpha_j \\
 &= \left( \sum_{i=1}^N \alpha_i y_i x_i \right)^T \left( \sum_{j=1}^N \alpha_j y_j x_j \right) \\
 &= \left\| \sum_{i=1}^N \alpha_i y_i x_i \right\|^2 \geq 0
 \end{aligned}$$

$\therefore A \odot B$  is PSD  $\rightarrow Q = -A \odot B$  is negative semi-definite  
 $\Rightarrow$  Function  $g$  is concave, every local maximum is global maximum

$\delta$  for support vector  $\alpha_i^* f_i(\theta^*) = 0$

$$\alpha_i^* \neq 0$$

for training set

$$\alpha_i^* = 0 \Rightarrow w^* \text{ and } b^* \text{ will not change}$$

直观地讲，这个结果来自于下面的主张（我们在下面证明）：如果在整个训练集上进行训练时， $x_i$ 不是支持向量，那么当把 $x_i$ 从训练集中排除时，最优的 $w^*$ 和 $b^*$ 不会改变。训练集之外的 $x_i$ 不会改变。  
 由于原始数据是线性可分离的，而且我们使用的是硬边际分类器，所以由原始 $w^*$ 和 $b^*$ 给出的假设不会在 $x_i$ 上出错，因此，在第 $i$ 步中也不会出错。因此，在LOOCV的第 $i$ 步中不会出现错误。等价地，LOOCV程序中唯一可能的错误是在整个训练集上训练时对支持向量的 $x_i$ 产生的，因此 $\epsilon \leq \frac{s}{N}$ 。

Intuitively, the result follows from the following claim (which we prove below): if  $x_i$  is not a support vector when training on the entire training set, then the optimal  $w^*$  and  $b^*$  do not change when leaving  $x_i$  out of the training set.

Since the original data are linearly separable and since we are using a hard-margin classifier, the hypothesis given by the original  $w^*$  and  $b^*$  will not make an error on  $x_i$ , and hence, no error will be made in the  $i$ -th step of the LOOCV. Equivalently, the only possible errors in the LOOCV procedure are made on  $x_i$ 's which are support vectors when training on the entire training set, and hence  $\epsilon \leq \frac{s}{N}$ .

(The following details are not required for full points).

Formally, let  $(w_D^*, b_D^*)$  and  $\alpha_D^*$  denote the optimal primal and dual solutions for the SVM when training on the entire  $D$ . Also let,  $D_i = D \setminus \{(x_i, y_i)\}$  be the set of training examples when omitting the  $i$ -th example, and let  $(w_{D_i}^*, b_{D_i}^*)$  and  $\alpha_{D_i}^*$  be the optimal primal and dual variables of the optimization problem when training on  $D_i$ .

$\alpha_{D_i}$  consists of only  $n - 1$  variables, namely  $\alpha_{D_i,1}, \dots, \alpha_{D_i,i-1}, \alpha_{D_i,i+1}, \dots, \alpha_{D_i,N}$ . Now consider setting the dual variables as follows  $\alpha_{D_i,j} = \alpha_{D,j}^*$  for  $j \neq i$ . Note that, if  $x_i$  is not a support vector when training on  $D$  then  $\alpha_{D,i}^* = 0$ . We can verify that  $(w_{D_i}^*, b_{D_i}^*)$  and  $\alpha_{D_i}^*$  satisfy the KKT conditions for the SVM optimization problem when training on  $D_i$  (e.g. the condition regarding the derivatives of the Lagrangian with respect to the primal variables is guaranteed to hold by our construction). From this, and the fact that  $w_{D_i}^*, b_{D_i}^*$  are unique since the objective function is strictly convex we can conclude that  $w^*, b^*$  do not change when omitting  $\{(x_i, y_i)\}$  as desired.