

Mehryar Mohri  
Foundations of Machine Learning 2014  
Courant Institute of Mathematical Sciences  
Homework assignment 3  
October 28, 2014  
Due: November 14, 2014

### A. Kernels

1. Show that the kernel  $K$  defined by

$$\forall (x, y) \in \mathbb{R}^N \times \mathbb{R}^N, \quad K(x, y) = \frac{1}{1 + \frac{\|x-y\|^2}{\sigma^2}}, \quad (1)$$

where  $\sigma > 0$  is a parameter, is PDS (*hint*: the function  $x \mapsto \int_0^{+\infty} e^{-sx} e^{-s} ds$  defined for all  $x \geq 0$  could be useful for the proof).

*Solution*: For  $x \geq 0$ , the integral  $\int_0^{+\infty} e^{-sx} e^{-s} ds$  is well defined and

$$\int_0^{+\infty} e^{-sx} e^{-s} ds = \int_0^{+\infty} e^{-s(1+x)} ds = \left[ -\frac{e^{-s(1+x)}}{1+x} \right]_0^{+\infty} = \frac{1}{1+x}. \quad (2)$$

Since the Gaussian kernel,  $(x, y) \mapsto e^{-\frac{\|x-y\|^2}{\sigma^2}}$  is PDS for any  $\sigma \neq 0$ , the kernel  $(x, y) \mapsto \int_0^{+\infty} e^{-\frac{s\|x-y\|^2}{\sigma^2}} e^{-s} ds$  is also PDS since for any  $x_1, \dots, x_m$  in  $\mathbb{R}^N$  and  $c_1, \dots, c_m$  in  $\mathbb{R}$ ,

$$\sum_{i,j=1}^m c_i c_j e^{-\frac{s\|x_i-x_j\|^2}{\sigma^2}} \geq 0 \Rightarrow \int_0^{+\infty} \sum_{i,j=1}^m c_i c_j e^{-\frac{s\|x_i-x_j\|^2}{\sigma^2}} e^{-s} ds \geq 0. \quad (3)$$

By (2),  $\int_0^{+\infty} e^{-\frac{s\|x-y\|^2}{\sigma^2}} e^{-s} ds = \frac{1}{1 + \frac{\|x-y\|^2}{\sigma^2}}$  for all  $x, y$  in  $\mathbb{R}^N$ , which concludes the proof.  $\square$

2. Show that the kernel  $K$  defined by

$$\forall (x, y) \in \mathbb{R}^N \times \mathbb{R}^N, \quad K(x, y) = \exp \left( \frac{\sum_{i=1}^N \min(|x_i|, |y_i|)}{\sigma^2} \right), \quad (4)$$

where  $\sigma > 0$  is a parameter, is PDS (*hint*: the function  $(x_0, y_0) \mapsto \int_0^{+\infty} 1_{t \in [0, |x_0|]} 1_{t \in [0, |y_0|]} dt$  defined over  $\mathbb{R} \times \mathbb{R}$  could be useful for the proof).

*Solution:* Observe that for all  $x, y \in \mathbb{R}$ ,

$$\int_0^{+\infty} 1_{t \in [0, |x|]} 1_{t \in [0, |y|]} dt = \min\{|x|, |y|\}. \quad (5)$$

Thus,  $(x, y) \mapsto \min(|x|, |y|)$  is a PDS kernel over  $\mathbb{R} \times \mathbb{R}$  since for any  $x_1, \dots, x_m$  in  $\mathbb{R}^N$  and  $c_1, \dots, c_m$  in  $\mathbb{R}$ ,

$$\begin{aligned} & \sum_{i,j=1}^m c_i c_j \min\{|x_i|, |x_j|\} \\ &= \int_0^{+\infty} \sum_{i,j=1}^m c_i c_j \sum_{k=1}^N 1_{t \in [0, |x_i|]} 1_{t \in [0, |x_j|]} dt \\ &= \int_0^{+\infty} \left\| \begin{bmatrix} c_1 1_{t \in [0, |x_1|]} \\ \vdots \\ c_m 1_{t \in [0, |x_m|]} \end{bmatrix} \right\|^2 dt \geq 0. \end{aligned}$$

Thus,  $(x, y) \mapsto \sum_{i=1}^N \min(|x_i|, |y_i|)$  is a PDS kernel over  $\mathbb{R}^N \times \mathbb{R}^N$  as a sum of PDS kernels. Its composition with  $x \mapsto e^{\frac{x}{\sigma^2}}$  with admits a power series with an infinite radius of convergence and non-negative coefficients is thus also PDS, which concludes the proof.  $\square$

## B. Support Vector Machines

1. Download and install the `libsvm` software library from:

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

and briefly consult the documentation to become more familiar with the tools.

2. Consider the `splice` data set

<http://www.cs.toronto.edu/~dave/data/splice/desc.html>.

Download the already formatted training and test files of a noisy version of that dataset from

[http://www.cs.nyu.edu/~mohri/ml14/splice\\_noise\\_train.txt](http://www.cs.nyu.edu/~mohri/ml14/splice_noise_train.txt)  
[http://www.cs.nyu.edu/~mohri/ml14/splice\\_noise\\_test.txt](http://www.cs.nyu.edu/~mohri/ml14/splice_noise_test.txt).

Use the `libsvm` scaling tool to scale the features of all the data. The scaling parameters should be computed only on the training data and then applied to the test data.

3. Consider the corresponding binary classification which consists of distinguishing two types of splice junctions in DNA sequences using about 60 features. Use SVMs combined with polynomial kernels to tackle this problem.

To do that, randomly split the training data into ten equal-sized disjoint sets. For each value of the polynomial degree,  $d = 1, 2, 3, 4$ , plot the average cross-validation error plus or minus one standard deviation as a function of  $C$  (let other parameters of polynomial kernels in `libsvm` be equal to their default values), varying  $C$  in powers of 5, starting from a small value  $C = 5^{-k}$  to  $C = 5^k$ , for some value of  $k$ .  $k$  should be chosen so that you see a significant variation in training error, starting from a very high training error to a low training error. Expect longer training times with `libsvm` as the value of  $C$  increases.

*Solution:*

Figure 1 shows the average cross-validation performance as a function of the regularization parameter  $C$ . Note that the algorithm starts to exhibit some over-fitting as  $C$  becomes very large. The performance for several choices of  $d$  and  $C$  are essentially indistinguishable; one suitable choice of optimal parameters is  $C^* = 5^1$  and  $d^* = 3$ .  $\square$

4. Let  $(C^*, d^*)$  be the best pair found previously. Fix  $C$  to be  $C^*$ . Plot the ten-fold cross-validation error and the test errors for the hypotheses obtained as a function of  $d$ . Plot the average number of support vectors obtained as a function of  $d$ . How many of the support vectors lie on the marginal hyperplanes? Plot the soft margin of the solution as a function of  $d$ .

*Solution:*

The first plot in Figure 4 compares CV and test errors for  $C^*$  as a function of  $g$ . As expected test error is slightly higher than CV error. The second plot shows that the total number of support vectors and the number of support vectors on the marginal hyperplanes. The last plot shows the margin as a function of  $g$ .  $\square$

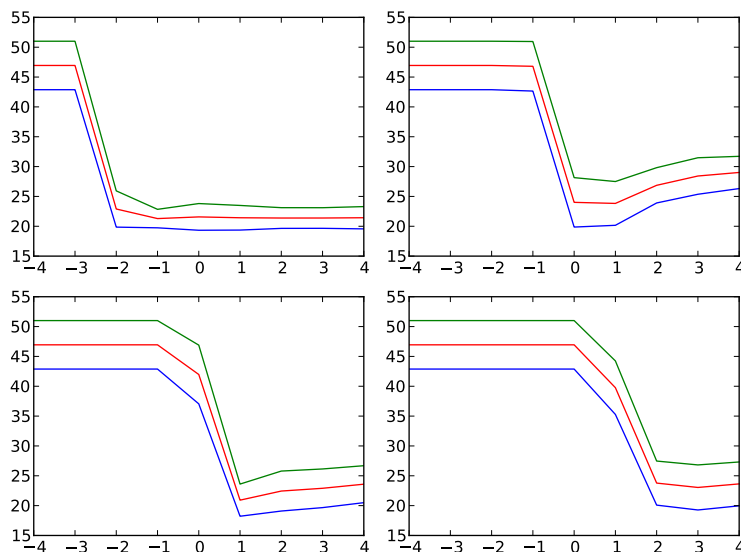


Figure 1: Average error (red) according to 10-fold cross-validation, with error-bars (green and blue) indicating one standard deviation. Top left, top right, bottom left and bottom right correspond to  $d = 1, 2, 3, 4$  respectively. On the  $x$ -axis we have  $\log_5 C$ .

5. Now, combine SVMs with Gaussian kernels to tackle the same task. Use cross-validation as before to determine the best value of  $C$  and  $\sigma$ , varying  $C$  in powers of 5, and  $\sigma$  in powers of 2 for a reasonable range so that you see a significant variation in training error, as before. Fix  $C$  and  $\sigma$  to the best values found via cross-validation. How does the test error of the solution compare to the best result obtained using polynomial kernels? What is the value of the soft margin?

*Solution:*

Figure 3 shows the average cross-validation performance as a function of the regularization parameter  $C$ .  $C^* = 1$  and  $g^* = 0.05$ .  $\square$

### C. Boosting

As discussed in class, AdaBoost can be viewed as coordinate descent applied to an exponential objective function. Here, we consider an alternative en-

semble method algorithm, HingeBoost, that consists of applying coordinate descent to an objective function based on the hinge loss. Using the same notation as in class, consider the function  $F$  defined for all  $\alpha \in \mathbb{R}^N$  by

$$F(\alpha) = \sum_{i=1}^m \max \left( 0, 1 - y_i \sum_{j=1}^N \alpha_j h_j(x_i) \right), \quad (6)$$

where the  $h_j$ s are base classifiers belonging to a hypothesis set  $H$  of functions taking values  $-1$  or  $+1$ .

1. Show that  $F$  is convex and admits a right- and left-derivative along any direction.

*Solution:*

Since the hinge loss is convex, its composition with affine function of  $\alpha$  is also convex and  $F$  is convex as a sum of convex functions.

For the existence of one-sided directional derivatives, one can use the fact that any convex function has one-sided directional derivatives or alternatively, that our specific function is the sum of piecewise affine functions, which are also known to have one-sided directional derivatives (think of one-dimensional hinge loss).

□

2. For any  $j \in [1, N]$ , let  $\mathbf{e}_j$  denote the direction corresponding to the base hypothesis  $h_j$ . Let  $\alpha_t$  denote the vector of coefficients  $\alpha_{t,j}$ ,  $j \in [1, N]$  obtained after  $t \geq 0$  iterations of coordinate descent and  $f_t = \sum_{j=1}^N \alpha_{t,j} h_j$  the predictor obtained after  $t$  iterations.

Give the expression of the right-derivative  $F'_+(\alpha_{t-1}, \mathbf{e}_j)$  and the left-derivative  $F'_-(\alpha_{t-1}, \mathbf{e}_j)$  after  $t-1$  iterations in terms of  $f_{t-1}$ .

*Solution:*

Distinguishing different cases depending on the value of  $y_i f_{t-1}(x_i) = 1$ , it is straightforward to derive the following expressions for all  $j \in [1, N]$ :

$$F'_+(\alpha_{t-1}, \mathbf{e}_j) = \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f_{t-1}(x_i) < 1} + 1_{(y_i h_j(x_i) < 0) \wedge (y_i f_{t-1}(x_i) = 1)}]$$

$$F'_-(\alpha_{t-1}, \mathbf{e}_j) = \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f_{t-1}(x_i) < 1} + 1_{(y_i h_j(x_i) > 0) \wedge (y_i f_{t-1}(x_i) = 1)}].$$

The key here is that when  $y_i f_{t-1}(x_i) \neq 1$ , each term in the sum will be either 0 or the affine function independent of  $y_i h_j(x_i)$ . On the other hand, when  $y_i f_{t-1}(x_i) = 1$ , the sign of  $y_i h_j(x_i)$  determines whether the finite differences will extend into the 0 portion of the affine portion of the term.  $\square$

3. For any  $j \in [1, N]$ , define the maximum directional derivative  $\delta F(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j)$  at  $\boldsymbol{\alpha}_{t-1}$  as follows:

$$\delta F(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) = \begin{cases} 0 & \text{if } F'_-(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) \leq 0 \leq F'_+(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) \\ F'_+(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) & \text{if } F'_-(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) \leq F'_+(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) \leq 0 \\ F'_-(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) & \text{if } 0 \leq F'_-(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j) \leq F'_+(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j). \end{cases}$$

The direction  $\mathbf{e}_j$  considered by the coordinate descent considered here is the one maximizing  $|\delta F(\boldsymbol{\alpha}_{t-1}, \mathbf{e}_j)|$ . Once the best direction  $j$  is selected, the step  $\eta$  can be determined by minimizing  $F(\boldsymbol{\alpha}_{t-1} + \eta \mathbf{e}_j)$  using a grid search. Give the pseudocode of HingeBoost.

*Solution:* The pseudocode of the HingeBoost algorithm is given in Figure 5.  $\square$

4. Bonus question: implement HingeBoost and AdaBoost using as base classifiers boosting stumps and compare their performance on the data set of Problem B. The number of rounds of boosting  $T$  can be determined via cross-validation varying  $T$  in powers of 10. Report the test errors of each algorithm and compare them with those obtained for SVMs in problem B.

*Solution:* The test errors of each algorithm should be approximately:

- SVM with Gaussian kernel: 18%
- SVM with Polynomial kernel: 20%
- AdaBoost: 17%
- HingeBoost: 19%

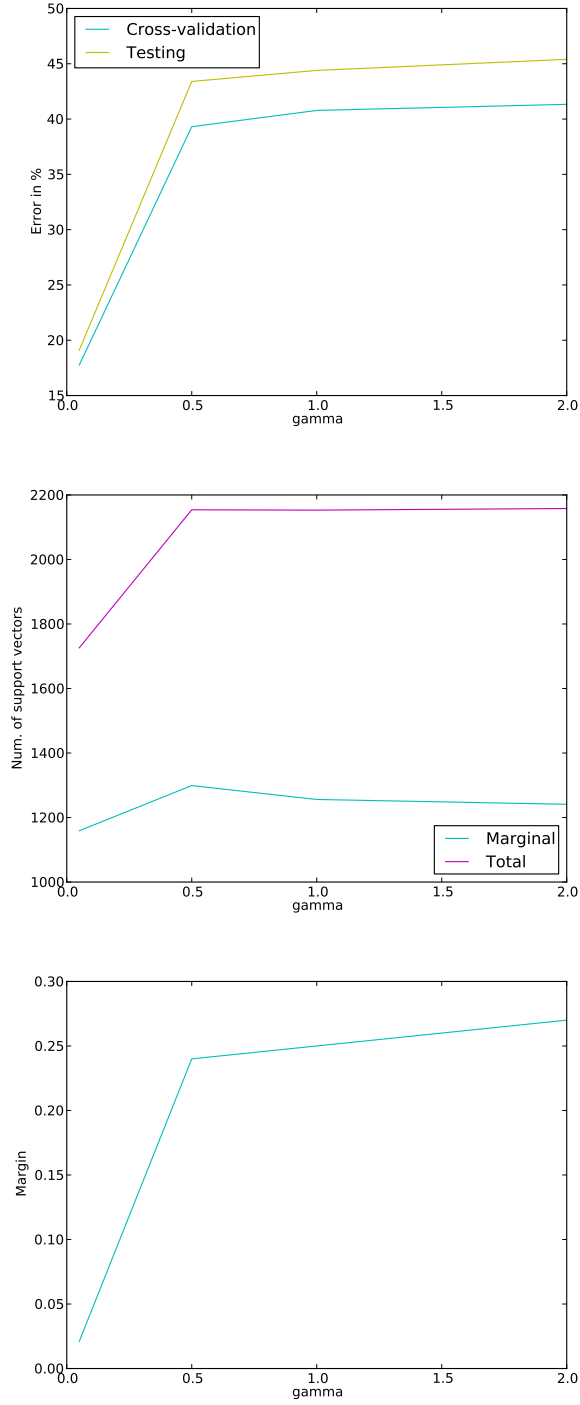


Figure 2: The test and validation error as a function of  $g$  (top panel), the number of total and marginal support vectors (second panel) and the margin as function of  $g$ .

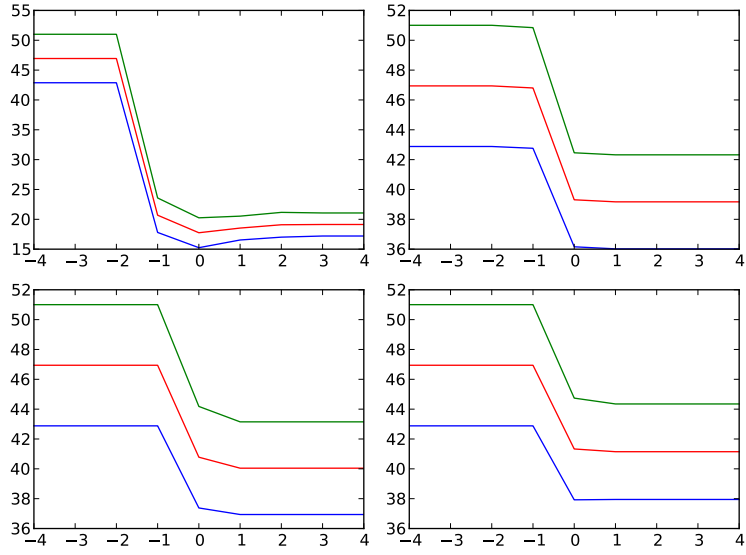


Figure 3: Average error (red) according to 10-fold cross-validation, with error-bars (green and blue) indicating one standard deviation. Top left, top right, bottom left and bottom right correspond to  $g = 0.05, 0.5, 1, 2$  respectively. Note that  $g = 1/2\sigma^2$ . On the  $x$ -axis we have  $\log_5 C$ .



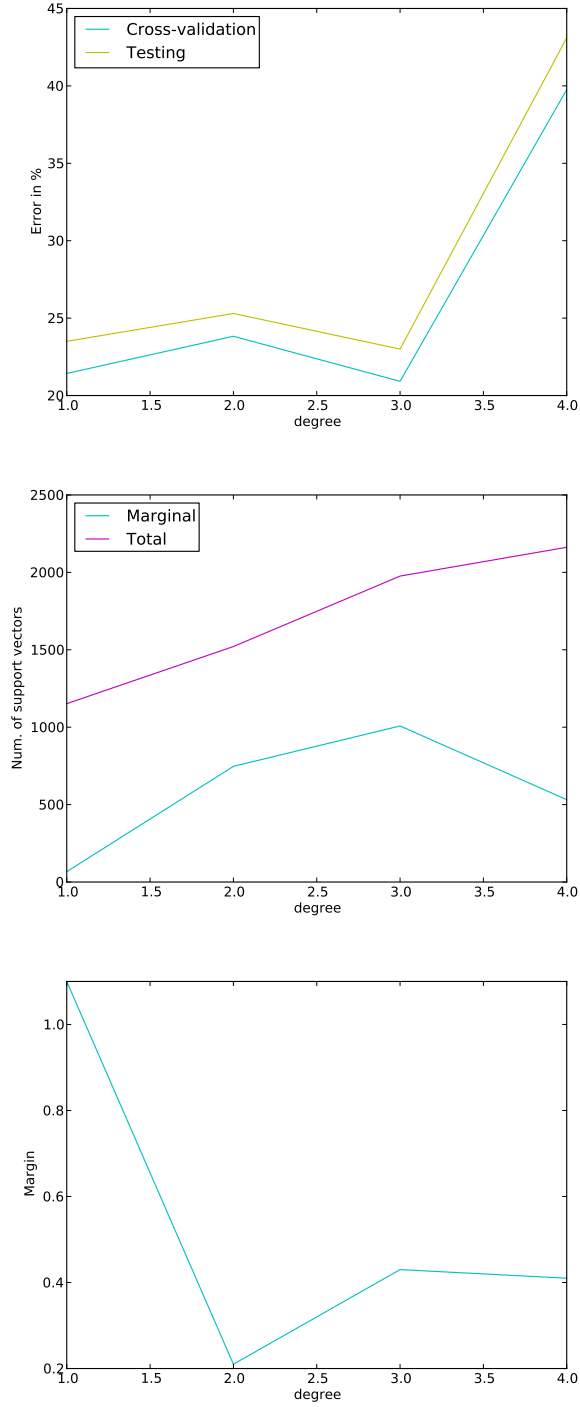


Figure 4: The test and validation error as a function of degree (top panel), the number of total and marginal support vectors (second panel) and the margin as function of  $d$ .

---

```

HINGEBOOST( $S = ((x_1, y_1), \dots, (x_m, y_m))$ )
1   $f \leftarrow 0$ 
2  for  $j \leftarrow 1$  to  $N$  do
3       $r \leftarrow \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f(x_i) < 1} + 1_{(y_i h_j(x_i) < 0) \wedge (y_i f(x_i) = 1)}]$ 
4       $l \leftarrow \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f(x_i) < 1} + 1_{(y_i h_j(x_i) > 0) \wedge (y_i f(x_i) = 1)}]$ 
5      if  $(l \leq 0) \wedge (r \geq 0)$  then
6           $d[j] \leftarrow 0$ 
7      elseif  $(l \leq r)$  then
8           $d[j] \leftarrow r$ 
9      else  $d[j] \leftarrow l$ 
10 for  $t \leftarrow 1$  to  $T$  do
11      $k \leftarrow \operatorname{argmin}_{j \in [1, N]} |d[j]|$ 
12      $\eta \leftarrow \operatorname{argmin}_{\eta \geq 0} G(f + \eta h_k) \triangleright \text{line search}$ 
13      $f \leftarrow f + \eta h_k$ 
14 return  $f$ 

```

---

Figure 5: Pseudocode of the HingeBoost algorithm. The function  $G$  is defined for any  $f$  by  $G(f) = \sum_{i=1}^m \max(0, f(x_i))$ .