

# CS 229, Autumn 2016

## Problem Set #2: Naive Bayes, SVMs, and Theory

---

**Due Wednesday, November 2 at 11:00 am on Gradescope.**

**Notes:** (1) These questions require thought, but do not require long answers. Please be as concise as possible. (2) If you have a question about this homework, we encourage you to post your question on our Piazza forum, at <https://piazza.com/stanford/autumn2016/cs229>. (3) If you missed the first lecture or are unfamiliar with the collaboration or honor code policy, please read the policy on Handout #1 (available from the course website) before starting work. (4) For problems that require programming, please include in your submission a printout of your code (with comments) and any figures that you are asked to plot.

If you are scanning your document by cellphone, please check the Piazza forum for recommended cellphone scanning apps and best practices.

### 1. [15 points] Constructing kernels

In class, we saw that by choosing a kernel  $K(x, z) = \phi(x)^T \phi(z)$ , we can implicitly map data to a high dimensional space, and have the SVM algorithm work in that space. One way to generate kernels is to explicitly define the mapping  $\phi$  to a higher dimensional space, and then work out the corresponding  $K$ .

However in this question we are interested in direct construction of kernels. I.e., suppose we have a function  $K(x, z)$  that we think gives an appropriate similarity measure for our learning problem, and we are considering plugging  $K$  into the SVM as the kernel function. However for  $K(x, z)$  to be a valid kernel, it must correspond to an inner product in some higher dimensional space resulting from some feature mapping  $\phi$ . Mercer's theorem tells us that  $K(x, z)$  is a (Mercer) kernel if and only if for any finite set  $\{x^{(1)}, \dots, x^{(m)}\}$ , the matrix  $K$  is symmetric and positive semidefinite, where the square matrix  $K \in \mathbb{R}^{m \times m}$  is given by  $K_{ij} = K(x^{(i)}, x^{(j)})$ .

Now here comes the question: Let  $K_1, K_2$  be kernels over  $\mathbb{R}^n \times \mathbb{R}^n$ , let  $a \in \mathbb{R}^+$  be a positive real number, let  $f : \mathbb{R}^n \mapsto \mathbb{R}$  be a real-valued function, let  $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^d$  be a function mapping from  $\mathbb{R}^n$  to  $\mathbb{R}^d$ , let  $K_3$  be a kernel over  $\mathbb{R}^d \times \mathbb{R}^d$ , and let  $p(x)$  a polynomial over  $x$  with *positive* coefficients.

For each of the functions  $K$  below, state whether it is necessarily a kernel. If you think it is, prove it; if you think it isn't, give a counter-example.

- (a) [1 points]  $K(x, z) = K_1(x, z) + K_2(x, z)$
- (b) [1 points]  $K(x, z) = K_1(x, z) - K_2(x, z)$
- (c) [1 points]  $K(x, z) = aK_1(x, z)$
- (d) [1 points]  $K(x, z) = -aK_1(x, z)$
- (e) [5 points]  $K(x, z) = K_1(x, z)K_2(x, z)$
- (f) [2 points]  $K(x, z) = f(x)f(z)$
- (g) [2 points]  $K(x, z) = K_3(\phi(x), \phi(z))$
- (h) [2 points]  $K(x, z) = p(K_1(x, z))$

[Hint: For part (e), the answer is that the  $K$  there *is* indeed a kernel. You still have to prove it, though. (This one may be harder than the rest.) This result may also be useful for another part of the problem.]

**Answer:** All 8 cases of proposed kernels  $K$  are trivially symmetric because  $K_1, K_2, K_3$  are symmetric; and because the product of 2 real numbers is commutative (for (1f)). Thanks to Mercer's theorem, it is sufficient to prove the corresponding properties for positive semidefinite matrices. To differentiate between matrix and kernel function, we'll use  $G_i$  to denote a kernel matrix (Gram matrix) corresponding to a kernel function  $K_i$ .

- (a) Kernel. The sum of 2 positive semidefinite matrices is a positive semidefinite matrix:  $\forall z \ z^T G_1 z \geq 0, z^T G_2 z \geq 0$  since  $K_1, K_2$  are kernels. This implies  $\forall z \ z^T G z = z^T G_1 z + z^T G_2 z \geq 0$ .
- (b) Not a kernel. Counterexample: let  $K_2 = 2K_1$  (we are using (1c) here to claim  $K_2$  is a kernel). Then we have  $\forall z \ z^T G z = z^T (G_1 - 2G_1) z = -z^T G_1 z \leq 0$ .
- (c) Kernel.  $\forall z \ z^T G_1 z \geq 0$ , which implies  $\forall z \ a z^T G_1 z \geq 0$ .
- (d) Not a kernel. Counterexample:  $a = 1$ . Then we have  $\forall z \ -z^T G_1 z \leq 0$ .
- (e) Kernel.  $K_1$  is a kernel, thus  $\exists \phi^{(1)} \ K_1(x, z) = \phi^{(1)}(x)^T \phi^{(1)}(z) = \sum_i \phi_i^{(1)}(x) \phi_i^{(1)}(z)$ . Similarly,  $K_2$  is a kernel, thus  $\exists \phi^{(2)} \ K_2(x, z) = \phi^{(2)}(x)^T \phi^{(2)}(z) = \sum_j \phi_j^{(2)}(x) \phi_j^{(2)}(z)$ .

$$K(x, z) = K_1(x, z) K_2(x, z) \quad (1)$$

$$= \sum_i \phi_i^{(1)}(x) \phi_i^{(1)}(z) \sum_i \phi_i^{(2)}(x) \phi_i^{(2)}(z) \quad (2)$$

$$= \sum_i \sum_j \phi_i^{(1)}(x) \phi_i^{(1)}(z) \phi_j^{(2)}(x) \phi_j^{(2)}(z) \quad (3)$$

$$= \sum_i \sum_j (\phi_i^{(1)}(x) \phi_j^{(2)}(x)) (\phi_i^{(1)}(z) \phi_j^{(2)}(z)) \quad (4)$$

$$= \sum_{(i,j)} \psi_{i,j}(x) \psi_{i,j}(z) \quad (5)$$

Where the last equality holds because that's how we define  $\psi$ . We see  $K$  can be written in the form  $K(x, z) = \psi(x)^T \psi(z)$  so it is a kernel.

Here is an alternate super-slick linear-algebraic proof. If  $G$  is the Gram matrix for the product  $K_1 \times K_2$ , then  $G$  is a principal submatrix of the Kronecker product  $G_1 \otimes G_2$ , where  $G_i$  is the Gram matrix for  $K_i$ . As the Kronecker product is positive semi-definite, so are its principal submatrices.

- (f) Kernel. Just let  $\psi(x) = f(x)$ , and since  $f(x)$  is a scalar, we have  $K(x, z) = \phi(x)^T \phi(z)$  and we are done.
- (g) Kernel. Since  $K_3$  is a kernel, the matrix  $G_3$  obtained for *any* finite set  $\{x^{(1)}, \dots, x^{(m)}\}$  is positive semidefinite, and so it is also positive semidefinite for the sets  $\{\phi(x^{(1)}), \dots, \phi(x^{(m)})\}$ .
- (h) Kernel. By combining (1a) sum, (1c) scalar product, (1e) powers, (1f) constant term, we see that any polynomial of a kernel  $K_1$  will again be a kernel.

## 2. [10 points] Kernelizing the Perceptron

Let there be a binary classification problem with  $y \in \{-1, 1\}$ . The perceptron uses hypotheses of the form  $h_\theta(x) = g(\theta^T x)$ , where  $g(z) = \text{sign}(z) = 1$  if  $z \geq 0$ ,  $-1$  otherwise.

In this problem we will consider a stochastic gradient descent-like implementation of the perceptron algorithm where each update to the parameters  $\theta$  is made using only one training example. However, unlike stochastic gradient descent, the perceptron algorithm will only make one pass through the entire training set. The update rule for this version of the perceptron algorithm is given by

$$\theta^{(i+1)} := \begin{cases} \theta^{(i)} + \alpha y^{(i+1)} x^{(i+1)} & \text{if } h_{\theta^{(i)}}(x^{(i+1)}) y^{(i+1)} < 0 \\ \theta^{(i)} & \text{otherwise,} \end{cases}$$

where  $\theta^{(i)}$  is the value of the parameters after the algorithm has seen the first  $i$  training examples. Prior to seeing any training examples,  $\theta^{(0)}$  is initialized to  $\vec{0}$ .

Let  $K$  be a Mercer kernel corresponding to some very high-dimensional feature mapping  $\phi$ . Suppose  $\phi$  is so high-dimensional (say,  $\infty$ -dimensional) that it's infeasible to ever represent  $\phi(x)$  explicitly. Describe how you would apply the “kernel trick” to the perceptron to make it work in the high-dimensional feature space  $\phi$ , but without ever explicitly computing  $\phi(x)$ . [Note: You don't have to worry about the intercept term. If you like, think of  $\phi$  as having the property that  $\phi_0(x) = 1$  so that this is taken care of.] Your description should specify

- How you will (implicitly) represent the high-dimensional parameter vector  $\theta^{(i)}$ , including how the initial value  $\theta^{(0)} = \vec{0}$  is represented (note that  $\theta^{(i)}$  is now a vector whose dimension is the same as the feature vectors  $\phi(x)$ );
- How you will efficiently make a prediction on a new input  $x^{(i+1)}$ . I.e., how you will compute  $h_{\theta^{(i)}}(x^{(i+1)}) = g(\theta^{(i)T} \phi(x^{(i+1)}))$ , using your representation of  $\theta^{(i)}$ ; and
- How you will modify the update rule given above to perform an update to  $\theta$  on a new training example  $(x^{(i+1)}, y^{(i+1)})$ ; i.e., using the update rule corresponding to the feature mapping  $\phi$ :

$$\theta^{(i+1)} := \theta^{(i)} + \alpha \mathbf{1}\{\theta^{(i)T} \phi(x^{(i+1)}) y^{(i+1)} < 0\} y^{(i+1)} \phi(x^{(i+1)}).$$

[Hint: our discussion of the representer theorem may be useful.]

### Answer:

In the high-dimensional space we update  $\theta$  as follows:

$$\theta := \theta + \alpha(y^{(i)} - h_{\theta}(\phi(x^{(i)})))\phi(x^{(i)})$$

So (assuming we initialize  $\theta^{(0)} = \vec{0}$ )  $\theta$  will always be a linear combination of the  $\phi(x^{(i)})$ , i.e.,  $\exists \beta_l$  such that  $\theta^{(i)} = \sum_{l=1}^i \beta_l \phi(x^{(l)})$  after having incorporated  $i$  training points. Thus  $\theta^{(i)}$  can be compactly represented by the coefficients  $\beta_l$  of this linear combination, i.e.,  $i$  real numbers after having incorporated  $i$  training points  $x^{(i)}$ . The initial value  $\theta^{(0)}$  simply corresponds to the case where the summation has no terms (i.e., an empty list of coefficients  $\beta_l$ ).

We do not work explicitly in the high-dimensional space, but use the fact that  $g(\theta^{(i)T} \phi(x^{(i+1)})) = g(\sum_{l=1}^i \beta_l \cdot \phi(x^{(l)})^T \phi(x^{(i+1)})) = g(\sum_{l=1}^i \beta_l K(x^{(l)}, x^{(i+1)}))$ , which can be computed efficiently.

We can efficiently update  $\theta$ . We just need to compute  $\beta_i = \alpha(y^{(i)} - g(\theta^{(i-1)T} \phi(x^{(i)})))$  at iteration  $i$ . This can be computed efficiently, if we compute  $\theta^{(i-1)T} \phi(x^{(i)})$  efficiently as described above.

In an alternative approach, one can observe that, unless a sample  $\phi(x^{(i)})$  is misclassified,  $y^{(i)} - h_{\theta^{(i)}}(\phi(x^{(i)}))$  will be zero; otherwise, it will be  $\pm 1$  (or  $\pm 2$ , if the convention  $y, h \in \{-1, 1\}$  is taken). The vector  $\theta$ , then, can be represented as the sum  $\sum_{\{i: y^{(i)} \neq h_{\theta^{(i)}}(\phi(x^{(i)}))\}} \alpha(2y^{(i)} - 1)\phi(x^{(i)})$  under the  $y, h \in \{0, 1\}$  convention, and containing  $(2y^{(i)})$  under the other convention. This can then be expressed as  $\theta^{(i)} = \sum_{i \in \text{Misclassified}} \beta_i \phi(x^{(i)})$  to be in more obvious congruence with the above. The efficient representation can now be said to be a list which stores only those indices that were misclassified, as the  $\beta_i$ s can be recomputed from the  $y^{(i)}$ s and  $\alpha$  on demand. The derivation for (b) is then only cosmetically different, and in (c) the update rule is to add  $(i + 1)$  to the list if  $\phi(x^{(i+1)})$  is misclassified.



- (a) Let two hypothesis classes  $H_1$  and  $H_2$  satisfy  $H_1 \subseteq H_2$ . Prove or disprove:  $VC(H_1) \leq VC(H_2)$ .
- (b) Let  $H_1 = H_2 \cup \{h_1, \dots, h_k\}$ . (I.e.,  $H_1$  is the union of  $H_2$  and some set of  $k$  additional hypotheses.) Prove or disprove:  $VC(H_1) \leq VC(H_2) + k$ . [Hint: You might want to start by considering the case of  $k = 1$ .]
- (c) Let  $H_1 = H_2 \cup H_3$ . Prove or disprove:  $VC(H_1) \leq VC(H_2) + VC(H_3)$ .

**Answer:**

- (a) True. Suppose that  $VC(H_1) = d$ . Then there exists a set of  $d$  points that is shattered by  $H_1$  (i.e., for each possible labeling of the  $d$  points, there exists a hypothesis  $h \in H_1$  which realizes that labeling). Now, since  $H_2$  contains all hypotheses in  $H_1$ , then  $H_2$  shatters the same set, and thus we have  $VC(H_2) \geq d = VC(H_1)$ .
- (b) True. If we can prove the result for  $k = 1$ , then the result stated in the problem set follows immediately by applying the same logic inductively, one hypothesis at a time. So, let us prove that if  $H_1 = H_2 \cup \{h\}$ , then  $VC(H_1) \leq VC(H_2) + 1$ . Suppose that  $VC(H_1) = d$ , and let  $S_1$  be a set of  $d$  points that is shattered by  $H_1$ . Now, pick an arbitrary  $x \in S_1$ . Since  $H_1$  shatters  $S_1$ , there must be some  $\bar{h} \in H_1$  such that  $h$  and  $\bar{h}$  agree on labelings for all points in  $S_1$  except  $x$ . This means that  $H' := H_1 \setminus \{h\}$  achieves all possible labelings on  $S' := S_1 \setminus \{x\}$  (i.e.  $H'$  shatters  $S'$ ), so  $VC(H') \geq |S'| = d - 1$ . But  $H' \subseteq H_2$ , so from part (a),  $VC(H') \leq VC(H_2)$ . It follows that  $VC(H_2) \geq d - 1$ , or equivalently,  $VC(H_1) \leq VC(H_2) + 1$ , as desired.

For this problem, there were a number of possible correct proof methods; generally, to get full credit, you needed to argue formally that there exists no set of  $(VC(H_2) + 2)$  points shattered by  $H_1$ , or equivalently, that there always exists a set of  $(VC(H_1) - 1)$  points shattered by  $H_2$ . Here are a couple of the more common errors:

- Some submitted solutions stated that adding a single hypothesis to  $H_2$  increases the VC dimension by at most one, since the new hypothesis can only realize a single labeling. While this statement is vaguely true, it is neither sufficiently precise, nor is its correctness immediately obvious.
  - Some solutions made arguments relating to the cardinality of the sets  $H_1$  and  $H_2$ . However, generally when we speak about VC dimension, the sets  $H_1$  and  $H_2$  often have infinite cardinality (e.g., the set of all linear classifiers in  $\mathbb{R}^2$ ).
- (c) False. Counterexample: let  $H_1 = \{h_1\}$ ,  $H_2 = \{h_2\}$ , and  $\forall x, h_1(x) = 0, h_2(x) = 1$ . Then we have  $VC(H_1) = VC(H_2) = 0$ , but  $VC(H_1 \cup H_2) = 1$ .

## 5. [20 points] Training and testing on different distributions

In the discussion in class about learning theory, a key assumption was that we trained and tested our learning algorithms on the same distribution  $\mathcal{D}$ . In this problem, we'll investigate one special case of training and testing on different distributions. Specifically, we will consider what happens when the training labels are *noisy*, but the test labels are not.

Consider a binary classification problem with labels  $y \in \{0, 1\}$ , and let  $\mathcal{D}$  be a distribution over  $(x, y)$ , that we'll think of as the original, "clean" or "uncorrupted" distribution. Define  $\mathcal{D}_\tau$  to be a "corrupted" distribution over  $(x, y)$  which is the same as  $\mathcal{D}$ , except that the labels  $y$  have some probability  $0 \leq \tau < 0.5$  of being flipped. Thus, to sample from  $\mathcal{D}_\tau$ ,

we would first sample  $(x, y)$  from  $\mathcal{D}$ , and then with probability  $\tau$  (independently of the observed  $x$  and  $y$ ) replace  $y$  with  $1 - y$ . Note that  $\mathcal{D}_0 = \mathcal{D}$ .

The distribution  $\mathcal{D}_\tau$  models a setting in which an unreliable human (or other source) is labeling your training data for you, and on each example he/she has a probability  $\tau$  of mislabeling it. Even though our training data is corrupted, we are still interested in evaluating our hypotheses with respect to the original, uncorrupted distribution  $\mathcal{D}$ .

We define the generalization error *with respect to*  $\mathcal{D}_\tau$  to be

$$\varepsilon_\tau(h) = P_{(x,y) \sim \mathcal{D}_\tau}[h(x) \neq y].$$

Note that  $\varepsilon_0(h)$  is the generalization error with respect to the “clean” distribution; it is with respect to  $\varepsilon_0$  that we wish to evaluate our hypotheses.

- (a) For any hypothesis  $h$ , the quantity  $\varepsilon_0(h)$  can be calculated as a function of  $\varepsilon_\tau(h)$  and  $\tau$ . Write down a formula for  $\varepsilon_0(h)$  in terms of  $\varepsilon_\tau(h)$  and  $\tau$ , and justify your answer.
- (b) Let  $|H|$  be finite, and suppose our training set  $S = \{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$  is obtained by drawing  $m$  examples IID from the corrupted distribution  $\mathcal{D}_\tau$ . Suppose we pick  $h \in H$  using empirical risk minimization:  $\hat{h} = \arg \min_{h \in H} \hat{\varepsilon}_S(h)$ . Also, let  $h^* = \arg \min_{h \in H} \varepsilon_0(h)$ .

Let any  $\delta, \gamma > 0$  be given. Prove that for

$$\varepsilon_0(\hat{h}) \leq \varepsilon_0(h^*) + 2\gamma$$

to hold with probability  $1 - \delta$ , it suffices that

$$m \geq \frac{1}{2(1-2\tau)^2\gamma^2} \log \frac{2|H|}{\delta}.$$

**Remark.** This result suggests that, roughly,  $m$  examples that have been corrupted at noise level  $\tau$  are worth about as much as  $(1 - 2\tau)^2 m$  uncorrupted training examples. This is a useful rule-of-thumb to know if you ever need to decide whether/how much to pay for a more reliable source of training data. (If you’ve taken a class in information theory, you may also have heard that  $(1 - \mathcal{H}(\tau))m$  is a good estimate of the information in the  $m$  corrupted examples, where  $\mathcal{H}(\tau) = -(\tau \log_2 \tau + (1 - \tau) \log_2 (1 - \tau))$  is the “binary entropy” function. And indeed, the functions  $(1 - 2\tau)^2$  and  $1 - \mathcal{H}(\tau)$  are quite close to each other.)

- (c) Comment **briefly** on what happens as  $\tau$  approaches 0.5.

**Answer:**

- (a) We compute  $\varepsilon_\tau$  as a function of  $\varepsilon_0$  and then invert the obtained expression. An error occurs on the corrupted distribution, if and only if, an error occurred for the original distribution and the point that was not corrupted, or no error occurred for the original distribution but the point was corrupted. So we have

$$\varepsilon_\tau = \varepsilon_0(1 - \tau) + (1 - \varepsilon_0)\tau$$

Solving for  $\varepsilon_0$  gives

$$\varepsilon_0 = \frac{\varepsilon_\tau - \tau}{1 - 2\tau}$$

