

1. Consider the email spam classification problem of Murphy Problem 8.1. Suppose you intend to use a linear perceptron classifier on that data. In the parts below, unless stated otherwise, assume the dataset of  $N = 4601$  samples is split into  $N_{Tr} = 3000$  for training and  $N_{Test} = 1601$  for testing. Also, for the tolerance  $\delta$  in the VC generalization bound, use 0.1 (for a certainty of 0.9). The parts below have short answers.

**Hint:** You may use the relation that if  $\mathcal{H}$  is a linear perceptron classifier in  $D$  dimensions ( $D$  features),  $d_{VC}(\mathcal{H}) = D + 1$ . (This will be proved in Problem 2.)

- What is the VC dimension of the hypothesis set?
- Expressing the upper bound on the out-of-sample error as  $E_{out}(h_g) \leq E_{in}(h_g) + \epsilon_{vc}$   
For  $E_{in}(h_g)$  measured on the training data, use  $d_{vc}$  from part (a) to get a value for  $\epsilon_{vc}$ .
- To get a lower  $\epsilon_{vc}$ , suppose you reduce the number of features to  $D = 10$ , and also increase the training set size to 10,000. Now what is  $\epsilon_{vc}$ ?
- Suppose that you had control over the number of training samples  $N_{Tr}$  (by collecting more email data). How many training samples would ensure a generalization error of  $\epsilon_{vc} = 0.1$  again with probability 0.9 (the same tolerance  $\delta = 0.1$ ), and using the reduced feature set (10 features)?
- Instead suppose you use the test set to measure  $E_{in}(h_g)$ , so let's call it  $E_{test}(h_g)$ . What is the hypothesis set now? What is its cardinality?
- Continuing from part (e), use the bound:  
$$E_{out}(h_g) \leq E_{test}(h_g) + \epsilon$$
Use the original feature set and the original test set, so that  $N_{Test} = 1601$ . Give an appropriate expression for  $\epsilon$  and calculate it numerically.

2. AML Exercise 2.4 (page 52). In addition to the hints given in the book, you can solve the problem by following the steps outlined below (on the next page).

For part (a):

- i. Write a point  $\underline{x}_i$  as a  $d+1$  dimensional vector;
- ii. Construct the  $(d+1) \times (d+1)$  matrix suggested by the book;
- iii. Write  $\underline{h}(\underline{X})$ , the output of the perceptron, as function of  $\underline{X}$  and the weights  $\underline{w}$  (note that  $\underline{h}(\underline{X})$  is a  $d+1$  dimensional vector with elements  $+1$  and  $-1$ );
- iv. Using the nonsingularity of  $\underline{X}$ , justify how any  $\underline{h}(\underline{X})$  can be obtained.

For part (b):

- i. Write a point  $\underline{x}_k$  as a linear combination of the other  $d+1$  points;
- ii. Write  $h(\underline{x}_k)$  (output for the chosen point) and substitute the value of  $\underline{x}_k$  by the expression just found on the previous item (**Hint:** use the  $\text{sgn}\{\cdot\}$  function);
- iii. What part of your expression in (ii) determines the class assignment of each point  $\underline{x}_i$ , for  $i \neq k$ ?
- iv. You have just proven (part (a)) that  $\underline{h}(\underline{X})$  with  $\underline{X}_{(d+1) \times (d+1)}$  can be shattered.

When we add a  $(d+2)^{\text{th}}$  line to  $\underline{X}$  can it still be shattered? In other words, can you choose the value of  $h(\underline{x}_k)$ ? Justify your answer. **Hint:** you can choose the class label of the other  $(d+1)$  points.

3. AML Problem 2.24 (page 75), except

>> Replace part (a) with:

(a.1) For a single given dataset, give an expression for  $g^{(D)}(x)$ . (AML notation)

(a.2) Find  $\bar{g}(x)$  analytically; express your answer in simplest form.

>> For parts (b) and (c), obtain  $E_D\{E_{out}\}$  by direct numerical computation, not by adding bias and var.

>> For part (d), obtain  $\text{bias}(x)$ ,  $\text{var}(x)$ ,  $\text{bias}$ ,  $\text{var}$ , and  $E_D\{E_{out}\}$ , all by analytical (pencil and paper) techniques.

## Reading

AML Sections 4.0, 4.1 (pp. 119-126): *Overfitting*

AML Section 4.2 (pp. 126-137): *Regularization (AML perspective)*

**Problem based on the reading**

4. AML Exercise 4.3 (p. 125). part (a) - first question only; part (b) - first question only. [Think about the second question of each part if you like; there isn't necessarily a single answer to each, though.]
5. AML Exercise 4.5 (p. 131)