

# Selected Exercises of Elements of Statistical Learning

Mike Craig

Last Updated September 3, 2017

## **1 Introduction**

This is a collection of certain exercise solutions I did as I was reading through this book. As I am not a mathematician, the "proofs" listed here are rough and probably not rigorous. They only serve as a means to explain the concepts in a way that aids my own understanding.

## 2 Overview of Supervised Learning

### Exercises

**Exercise 2.1.** Suppose each of  $K$ -classes has an associated target  $t_k$ , which is a vector of all zeros, except a one in the  $k$ th position. Show that classifying to the largest element of  $\hat{y}$  amounts to choosing the closest target,  $\min_k ||t_k - \hat{y}||$ , if the elements of  $\hat{y}$  sum to one.

We are trying to show that

$$\arg \max_k \hat{y}_k = \arg \min_k ||t_k - \hat{y}||$$

To do this, we can show that for any  $k^* \neq \arg \max_k \hat{y}_k$  and  $k = \arg \max_k \hat{y}_k$ ,

$$||t_{k^*} - \hat{y}|| > ||t_k - \hat{y}||$$

Note that we can equivalently consider the  $||\cdot||^2$  instead of  $||\cdot||$ , because they are both monotonic for  $\geq 0$ .

Using the definition of the Euclidean norm,

$$||t_{k^*} - \hat{y}||^2 = ||\hat{y}||^2 + ||t_{k^*}||^2 - 2t_{k^*}\hat{y} \quad (1)$$

$$= \hat{y}^2 \quad (2)$$

$$(3)$$

Similarly,

$$||t_k - \hat{y}||^2 = ||\hat{y}||^2 + ||t_k||^2 - 2t_k\hat{y} \quad (4)$$

$$= \hat{y}^2 + 1 - 2\hat{y} \quad (5)$$

Therefore

$$||t_{k^*} - \hat{y}||^2 - ||t_k - \hat{y}||^2 = \hat{y}^2 - (\hat{y}^2 + 1 - 2\hat{y}) \quad (6)$$

$$= -1 + 2\hat{y} \quad (7)$$

$$\geq 0 \quad (8)$$

since  $\hat{y}$  sums to one. So,  $||t_k - \hat{y}||$  is minimized when  $k = \arg \max_k \hat{y}_k$ .

**Exercise 2.2.** Show how to compute the Bayes decision boundary for the simulation example in Figure 2.5.

**Exercise 2.3.** Consider  $N$  data points uniformly distributed in a  $p$ -dimensional unit ball centered at the origin. Suppose we consider a nearest-neighbor estimate at the origin. Show that the median distance from the origin to the closest data point is given by:  $d(p, N) = (1 - \frac{1}{2}^{1/N})^{1/p}$ .

Let  $m$  be the median distance from the origin to the closest point. This means that the probability that all data points are further than  $m$  is 0.5. "Further" simply means a greater norm. Since samples  $x_i$  are i.i.d, we can more formally state this as:

$$\prod_{i=1}^N P(\|x_i\| > m) = \frac{1}{2}$$

Note that we can flip this around to be  $\prod_{i=1}^N P(\|x_i\| \leq m) = \frac{1}{2}$ . Now we can use the cumulative function of the uniform distribution as follows:

$$\prod_{i=1}^N P(\|x_i\| \leq m) = \prod_{i=1}^N (1 - m^p) \quad (9)$$

$$= \prod_{i=1}^N (1 - m^p) \quad (10)$$

$$= (1 - m^p)^N = \frac{1}{2} \quad (11)$$

Now we can solve for  $m$ :

$$\frac{1}{2} = (1 - m^p)^N \quad (12)$$

$$\frac{1}{2}^{1/N} = 1 - m^p \quad (13)$$

$$m^p = 1 - \frac{1}{2}^{1/N} \quad (14)$$

$$m = (1 - \frac{1}{2}^{1/N})^{1/p} \quad (15)$$

*Exercise 2.4.* The edge effect problem discussed on page 23 is not peculiar to uniform sampling from bounded domains. Consider inputs drawn from a spherical multinormal distribution  $XN(0, I_p)$ . The squared distance from any sample point to the origin has a  $\chi_p^2$  distribution with mean  $p$ . Consider a prediction point  $x_0$  drawn from this distribution, and let  $a = x_0/\|x_0\|$  be an associated unit vector. Let  $z_i = a^T x_i$  be the projection of each of the training points on this direction. Show that the  $z_i$  are distributed  $N(0, 1)$  with expected squared distance from the origin 1, while the target point has expected squared distance  $p$  from the origin. Hence for  $p = 10$ , a randomly drawn test point is about 3.1 standard deviations from the origin, while all the training points are on average one standard deviation along direction  $a$ . So most prediction points see themselves as lying on the edge of the training set.