

## Statistical learning theory

### Exercise T8.1: Empirical Risk Minimization

(tutorial)

It is known from the literature (e.g., Anthony & Bartlett book from 1999) that the number of *affinely*<sup>1</sup> separable assignments of  $p$  data points in  $N$  dimensions is given by

$$\tilde{C}_{(p,N)} := 2 \sum_{k=0}^N \binom{p-1}{k}.$$

- (a) How is the Binomial coefficient  $\binom{n}{k}$  defined for  $n, k \in \mathbb{N}_0$  with  $n \geq k$ ?
- (b) Use the recursive formula for binomial coefficients

$$\binom{n}{k} + \binom{n}{k-1} = \binom{n+1}{k} \quad \text{to show that} \quad \tilde{C}_{(p,N)} + \tilde{C}_{(p,N-1)} = \tilde{C}_{(p+1,N)}.$$

- (c) Explain how the number of affinely separable assignments is related to the *Vapnik-Chervonenkis (VC) dimension* and overfitting.

### Exercise H8.1: Vapnik-Chervonenkis dimension

(homework, 2 points)

Use the definition of  $\tilde{C}_{(p,N)}$  and the recursion property from above together with the binomial formula,

$$(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k,$$

for  $x, y \in \mathbb{R}$ ,  $n \in \mathbb{N}$ , to show that a linear classifier,  $y(\mathbf{x}; \mathbf{w}) = \text{sign}(w_0 + \sum_{i=1}^N x_i w_i)$ , has a Vapnik-Chervonenkis dimension of  $d_{VC} = N + 1$ .

*Hint:* It suffices to show that  $\tilde{C}_{(N+1,N)} = 2^{N+1}$  and  $\tilde{C}_{(N+2,N)} < 2^{N+2}$ .

---

<sup>1</sup>Note that in the lecture's formula for  $C_{(p,N)}$  the sum runs to  $N-1$  only (instead of  $N$ ). This is due to the fact that the lecture considers *linear* separability, i.e., no bias term, whereas here we consider the linear neuron with bias, that is,  $y(\mathbf{x}) = \text{sign}(w_0 + \sum_{i=1}^N w_i x_i)$ . Therefore we define the tilde-variant  $\tilde{C}_{(p,N)}$  here.

**Exercise H8.2: Variability of classification (homework, 5 points)**

Assume data  $\underline{x}^{(\alpha)} \in \mathbb{R}^2$  drawn from two clusters,  $C_1$  and  $C_2$  and distributed according to the (multivariate) Normal distributions  $\mathcal{N}(\underline{\mu}_i, 2\mathbf{I})$ ,  $i = 1, 2$  with  $\underline{\mu}_1 = (0, 1)^T$ ,  $\underline{\mu}_2 = (1, 0)^T$ , and identity matrix  $\mathbf{I}$ . This task examines, how well a linear connectionist neuron can separate these two classes for increasing amounts  $N$  of available training data. Proceed as follows:

1. Generate a sample of  $N/2$  data points  $\underline{x}^{(\alpha)}$  from each of the two clusters (i.e.,  $N$  data points in total). Let  $y^{(\alpha)} = 1$  for  $\underline{x}^{(\alpha)}$  from  $C_1$  and  $y^{(\alpha)} = -1$  for  $\underline{x}^{(\alpha)}$  from  $C_2$ .
2. Find the weights of a linear connectionist neuron with output  $y(\underline{x}) = \text{sign}(w_0 + \sum_{i=1}^N w_i x_i)$  minimizing the squared error according to the known analytical formula for (see, e.g., problem T4.2c).
3. Find the predictions of this classifier for  $N_{\text{test}} = 1000$  new data drawn from the same distributions (again, 50% of the data points for each class).
4. Calculate the accuracies (percentage of correct classifications) for the training ( $r_{\text{train}}$ ) and test samples ( $r_{\text{test}}$ ).

Repeat these steps 50 times for each  $N \in \{2, 4, 6, 8, 10, 20, 40, 100\}$ , and save the resulting parameters as well as the accuracies for training and testing.

- (a) (2 point) Plot the mean and standard deviation of  $r_{\text{train}}$ , and  $r_{\text{test}}$  against the number of samples  $N$  in an errorbar-plot (plotting  $N$  on the x-axis and the corresponding statistic on the y-axis).
- (b) (1 point) Plot the means and standard deviation of  $w_1$ ,  $w_2$  and  $w_0$  against  $N$ .
- (c) (2 points) Interpret your results. How do these estimates depend on  $N$ ?

**Exercise H8.3: The Binomial distribution (homework, 3 points)**

This exercise examines the relation between the following 3 distributions that are used in the statistical learning theory:

$$f(k; n, p) = \binom{n}{k} p^k (1-p)^{n-k} \quad (\text{Binomial distribution})$$

$$f(x; \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (\text{Normal distribution})$$

$$f(k; \lambda) = \frac{\lambda^k}{k!} e^{-\lambda} \quad (\text{Poisson distribution})$$

- (a) (1 point) Visualize the probability mass function  $f(k; n, p)$  of the binomial distribution for a few different values of  $k, n, p$  that demonstrate the different shapes that function can have.
- (b) (1 point) The normal distribution is sometimes used as an approximation to the binomial distribution. Under which conditions is this reasonable? Under which conditions is it problematic? Visualize one example where the Normal approximation is good and one where it is not. Give at least one reason why this distribution is so widely used.
- (c) (1 point) The Poisson distribution is often used as an alternative approximation to the binomial distribution. Under which conditions is it a good approximation? Visualize one example parametrization where the Poisson approximation is good and one where it is not.

**Total 10 points.**