

## EXAM 18.01.14 B

### Question 1

- 1). This dataset is linearly separable because exists a separation surface that splits our instance space into 3 regions such that different classified instances are separated.
- 2). You can use a kernel function  $K(x_i, x_j) = (Bx_i^T x_j + r)^d$ , for example with  $d=3$ . This is a polynomial kernel function.
- 3). We have to use also in this case the Kernel trick:

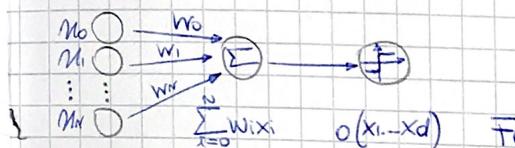
$y(x_2) = \text{sign}(w_0 + \sum_{n=1}^N \alpha_n x_n^T x_2)$  is possible to see an inner product.

$y(x_2) = \text{sign}(w_0 + \sum_{n=1}^N \alpha_n K(x_n, x_2))$  we have applied the Kernel trick.

$$w_0 = \frac{1}{|SV|} \sum_{x_i \in SV} (t_i - \sum_{x_j \in S} \alpha_j t_j K(x_i, x_j))$$

### Question 2

- 1). The perceptron model is based on the following structure:



$$\sigma(x_0, \dots, x_d) = \begin{cases} +1 & \text{if } w_0 + w_1 x_1 + \dots + w_d x_d > 0 \\ -1 & \text{otherwise} \end{cases} = \text{sign}(w^T x)$$

For the moment for simplicity we consider only  $\sigma(x) = w^T x$ .

We want to minimize the error function  $E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - \sigma(x_n))^2 = \frac{1}{2} \sum_{n=1}^N (t_n - w^T x_n)^2$

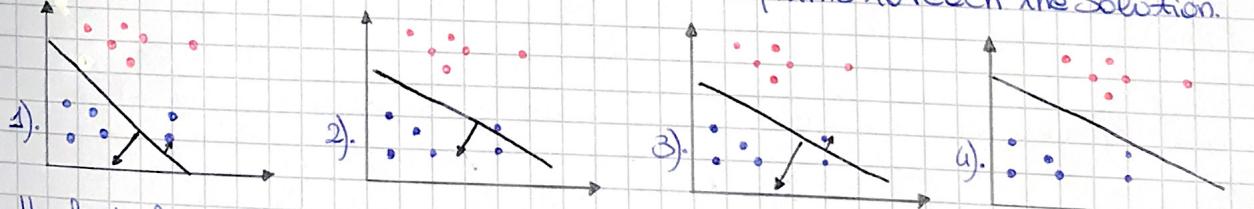
We compute the gradient:  $\frac{\partial E(w)}{\partial w_i} = \sum_{n=1}^N (t_n - w^T x_n)(-x_{i,n})$

- 2). In perceptron we try to find the optimal solution using a sequential method:

$\hat{w} \leftarrow \hat{w} + \Delta w_i \Rightarrow \hat{w} \leftarrow \hat{w} + \eta \sum_{n=1}^N (t_n - \hat{w}^T x_n)(x_{i,n})$ . Now we consider another time the sign function:  $\hat{w} \leftarrow \hat{w} + \eta \sum_{n=1}^N (t_n - \text{sign}(\hat{w}^T x_n))(x_{i,n})$ .

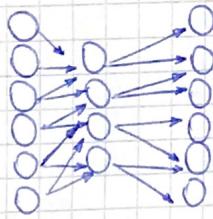
$\eta$  is a very important hyperparameter called learning rate. If  $\eta$  is too large we can diverge, if  $\eta$  is too small we can take a lot of time to reach the solution.

2).



In the first plot is possible to see that there are misclassified points. We sum the vector related to the separation surface with the feature vector of the misclassified point and we obtain the plot number 2. In the plot number 3 we do the same and we obtain the correct separation surface in the plot 4.

### Question 3



$$\# \text{parameters} = (7 \times 4) + (5 \times 6) = 58 = (6 \times 4 + 4) + (4 \times 5 + 6)$$

Backpropagation is not affected by local minima because it is not a learning algorithm but is an algorithm used to compute the gradient and to propagate the gradient through all the network. For this reason it is not affected by overfitting.

### Question 4

1). In linear regression we have a target function  $f: X \rightarrow Y$  with  $X \in \mathbb{R}^d$  and  $Y \in \mathbb{R}$  and a dataset  $\mathcal{D} = \{(x_n, t_n)\}_{n=1}^N$ . The goal is to find  $w^* = \underset{w}{\operatorname{argmin}} E_{\mathcal{D}}(w)$ . We have a model  $y(x, w) = \sum_{j=1}^N w_j \phi_j$

and we know that the target value is  $t = y(x, w) + \varepsilon$  with  $\varepsilon$  additive Gaussian noise. Using the information about the Gaussian  $p(\varepsilon | \beta) = N(\varepsilon | 0, \beta^{-1}) \Rightarrow p(t | x_1, \dots, x_N, w, \beta) = N(t | y(x, w), \beta)$ .

Now we assume the iid hypothesis:  $p(\{t_1, \dots, t_N\} | x_1, \dots, x_N, w, \beta) = \prod_{n=1}^N N(t_n | w^\top \phi(x_n), \beta^{-1}) = \prod_{n=1}^N \ln N(t_n | w^\top \phi(x_n), \beta^{-1}) = -\beta \cdot \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 / 2 \ln(2\pi\beta^{-1})$

We have found the error function  $E_{\mathcal{D}}(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$

2). For the batch mode we can use least squares:  $E_{\mathcal{D}}(w) = \frac{1}{2} (t - \phi w)^\top (t - \phi w)$  and  $w^* = \phi^+ t$ .

If we use sequential learning:  $\hat{w} \leftarrow \hat{w} + \eta [t_n - w^\top \phi(x_n)] \phi(x_n)$  with  $\eta$  learning rate.

### Question 5

1).  $h_{ML} = \underset{h \in H}{\operatorname{argmax}} \frac{p(D|h)p(h)}{p(D)} = \underset{h \in H}{\operatorname{argmax}} p(h|D) = \underset{h \in H}{\operatorname{argmax}} p(D|h)p(h)$

If we assume that  $p(h) = p(h_j)$  we can simplify and obtain the ML hypothesis:  $h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(D|h)$

2). given a target function  $P: X \rightarrow V$  knowing that each instance can be described as  $\langle o_1, \dots, o_N \rangle$ :

$$P(v|x, D) = P(v|o_1, \dots, o_N, D) = \frac{P(o_1, \dots, o_N|v, D)P(v|D)}{P(o_1, \dots, o_N|D)} = P(o_1, \dots, o_N|v, D)P(v|D)$$

$V_{NB} = \underset{v \in V}{\operatorname{argmax}} P(o_1, \dots, o_N|v, D)P(v|D)$ . Now we have to introduce a simplification. Each sample is mutually conditionally independent from the other  $\rightarrow P(o_1, \dots, o_N|v, D) = \prod_{i=1}^N P(o_i|v, D)$

$$\text{We obtain that: } V_{NB} = \underset{v \in V}{\operatorname{argmax}} P(v|D) \prod_{i=1}^N P(o_i|v, D)$$

3). NBC is not an optimal classifier differently from BOC because it is based on the conditional independence assumption, but can be used in all the cases.