

Question 1

- 1). If $C = C_1 \text{ AND } B = b_1$ THEN NO ①
- 2). If $C = C_1 \text{ AND } B = b_2$ THEN YES ②
- 3). If $C = C_2 \text{ AND } A = a_1$ THEN YES ③
- 4). If $C = C_2 \text{ AND } A = a_2 \text{ AND } B = b_1$ THEN YES ④
- 5). If $C = C_2 \text{ AND } A = a_2 \text{ AND } B = b_2$ THEN NO ⑤
- 6). If $C = C_2 \text{ AND } A = a_3$ THEN NO ⑥
- 7). If $C = C_3$ THEN NO ⑦

- 2). T is consistent with S_1 because of ①.
 T " " " S_2 " " " ④.
 T is not consistent with S_3 because of ⑦.
" " " S_4 " " " ⑤.

Question 2

- 1). The maximum a posteriori hypothesis is $h_{MAP} = \underset{h \in H}{\operatorname{argmax}} \frac{p(D|h)p(h)}{p(D)} = \underset{h \in H}{\operatorname{argmax}} p(D|h)p(h)$

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

- 2). Boc is an optimal classifier, it returns always the most probable prediction.

$$V_{ML} = \underset{v \in V}{\operatorname{argmax}} \sum_{h \in H} p(v|x, D, h) p(h|x, D) = \sum_{h \in H} p(v|a, h) p(h|D)$$

- 3). The problem of Boc is that we can use it only if we have analytical solutions or if the hypotheses space is not large, otherwise is not practical.

Question 3

- 1). In linear regression we have a target function $f: X \rightarrow Y$ with $X \subseteq \mathbb{R}^d$ and $Y = \mathbb{R}$. The model is $y(x, w) = \sum_{j=1}^n w_j \phi_j(x)$ that is still linear in the parameters w . We want to

find $w^* = \underset{w}{\operatorname{argmin}} E_S(w)$ with $E_S(w)$ error function but we have to discover $E_S(w)$.

We have a dataset $D = \{(x_n, t_n)\}_{n=1}^N$ and we know that: $t = y(x, w) + \epsilon$ with ϵ additive Gaussian noise. $p(\epsilon | \beta) = N(\epsilon | 0, \beta^{-1}) \Rightarrow p(t | n_1, \dots, n_N, w, \beta) = N(t | y(x, w), \beta^{-1})$

$$\text{If we use the iid hypothesis: } p(\{t_1, \dots, t_N\} | n_1, \dots, n_N, w, \beta) = \prod_{n=1}^N p(t_n | w^\top \phi(x_n), \beta^{-1}) = \prod_{n=1}^N \ln N(t_n | w^\top \phi(x_n), \beta^{-1}) = -\beta \cdot \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 - \frac{N}{2} \ln(2\pi\beta^{-1})$$

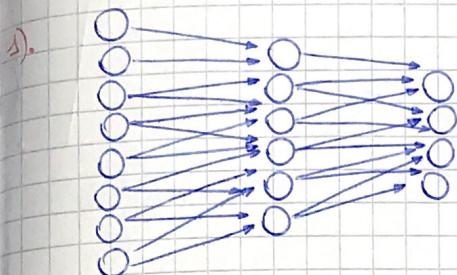
$$E_S(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$$

2). We can find w^* in batch or in a sequential way:

if we use least squares $E_\Delta(w) = \frac{1}{2} (\mathbf{t} - \mathbf{\Phi}w)^T (\mathbf{t} - \mathbf{\Phi}w)$ and $w^* = \mathbf{\Phi}^T \mathbf{t}$.

if we use sequential learning: $\hat{w} \leftarrow \hat{w} + \eta [t_n - w^T \phi(x_n)] \phi(x_n)$

Question 4



2). Backpropagation is an algorithm that is used for compute the gradient and the gradient will be propagated throw all the network. So BackProp is not affected by overfitting or by local minima.

Question 5

1). The dataset is linearly separable because exists a separation surface that splits our instance space into two regions such that different classified instances are separated.

2). I can use a kernel function of the polynomial type $K(x, x') = (Bx^T x' + \gamma)^d$

Using for example $d=3$.

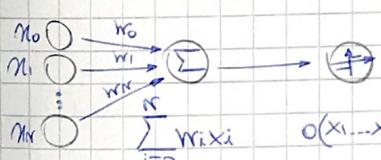
3). $y(x_i, x) = \text{Sign}(w_0 + \sum_{n=1}^N w_n x_n^T x)$ if apply the kernel trick

$$y(x_i, x) = \text{Sign}(w_0 + \sum_{n=1}^N w_n K(x_n, x))$$

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

Question 6

1). Perceptron is based on the following structure:



$$o(x_1, \dots, x_d) = \begin{cases} +1 & \text{if } w_0 + w_1 x_1 + \dots + w_d x_d > 0 \\ -1 & \text{otherwise} \end{cases}$$

For the moment we will consider $o(x) = w^T x$

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

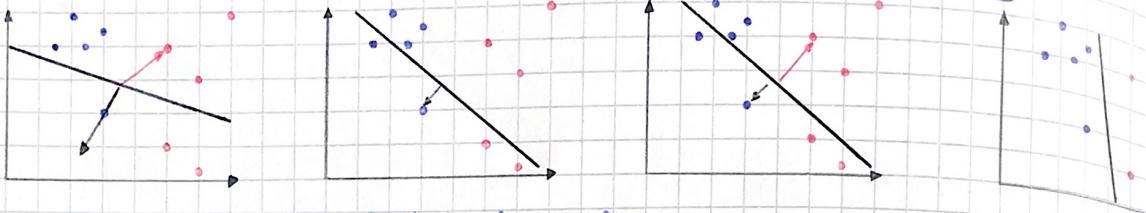
$\frac{\partial E(w)}{\partial w_i} = \sum_{n=1}^N (t_n - w^T x_n)(-x_{in})$. We find w^* in a sequential way.

$$\hat{w} \leftarrow \hat{w} + \Delta w_i \text{ with } \Delta w_i = -\eta \sum_{n=1}^N (t_n - w^\top x_n)(x_{i,n}).$$

$$\text{Now we consider another time the sign function: } \Delta w_i = -\eta \sum_{n=1}^N (\text{sign}(w^\top x_n))(x_{i,n})$$

η is a very important hyperparameter called learning rate, if η is too small we can take a lot of time to reach the solution, if η is too big we can diverge.

2).



w is plotted as a vector (black). In the first plot is possible to see that we have an error, so we add the vector w to the feature vector of the misclassified point. In the plot 3 is possible to see that we have other misclassified points so we have to sum another time the vector w to the feature vectors of the misclassified point.