

Question 1

$$\begin{array}{lll} 1). \quad p(YH) = 0.25 & p(YW) = 0.45 & p(O) = 0.3 \\ p(S|YH) = 0.3 & p(T|YH) = 0.5 & p(R|YH) = 0.2 \\ p(S|YW) = 0.5 & p(T|YW) = 0.3 & p(R|YW) = 0.2 \\ p(S|O) = 0.3 & p(T|O) = 0.3 & p(R|O) = 0.4 \end{array}$$

$$v^* = \operatorname{argmax}_{c \in C} \{p(T|c)p(c)\} = \operatorname{argmax}_{c \in C} \{0.5 \cdot 0.25, 0.3 \cdot 0.45, 0.3 \cdot 0.3\} =$$

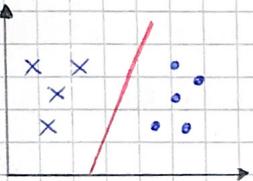
$$= \operatorname{argmax} \{0.125, 0.135, 0.09\} = YW$$

$$2). L = \sum_{c \in C} p(T|c)p(c) = 0.125 + 0.135 + 0.09 = 0.35 \text{ (likelihood)}$$

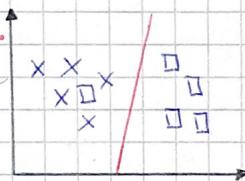
Question 2

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

2).



3).



4). The first dataset is linearly separable. It's possible to see that this solution is computed with SVM because SVM aims at maximum margin with the better accuracy.

Let x_n be the closest point to the separation surface $\bar{h} = \bar{w}_0 + \bar{w}^T x = 0$, we define

margin $\frac{|y(x_n)|}{\|\bar{w}\|}$. To find the margin we have to find $\min_{n=1 \dots N} \frac{\bar{t}_n (\bar{w}^T x_n + \bar{w}_0)}{\|\bar{w}\|}$. But

we want to find the maximum margin: $w^*, w^* = \operatorname{argmax}_{w_0, w} \frac{1}{\|\bar{w}\|} \min_{n=1 \dots N} \bar{t}_n (\bar{w}^T x_n + \bar{w}_0)$

Scaling we are not affecting our solution: $t_n y(x_n) = 1$ and $t_n y(x_n) \geq 1 \quad \forall n = 1 \dots N$

In this mode $w^*, w_0^* = \operatorname{argmin}_{w_0, w} \frac{1}{2} \|w\|^2$. Now we move to a Lagrangian domain.

$w^* = \sum_{n=1}^N \alpha_n^* t_n x_n$. We know that, using the KKT hypothesis, if $t_n y(x_n) \geq 1 \Rightarrow \alpha_n^* = 0$

We obtain that $w^* = w_0^* + \sum_{x_j \in SV} \alpha_j^* t_j x_j^T x_j = 0$ with $SV = \{x_n \in D | t_n y(x_n) = 1\}$

Also the second solution is computed with SVM. In fact, we can use SVM if the dataset is not perfectly linearly separable introducing slack variables $\xi_n \geq 0$.

$y_n(x_n) = 1 - \xi_n \quad \forall n = 1 \dots N$ and the modified optimal value w^* :

$$w^* = \underset{w}{\operatorname{argmin}} \frac{1}{2} \|w\|^2 + C \sum_{n=1}^N \xi_n$$

Question 3

1). In PCA we want to maximize the data variance after the projection to some direction U_1 . The projected points are $U_1^T x_n$ subject to $U_1^T U_1 = 1$.

The data variance is: $\frac{1}{N} \sum_{n=1}^N (U_1^T x_n - U_1^T \bar{x})^2 = U_1^T S U_1$ with $\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$

$$\text{and } S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T = X^T X$$

We want to find $\max_{U_1} U_1^T S U_1$ s.t. $U_1^T U_1 = 1 \iff \max_{U_1} U_1^T S U_1 + \lambda_1(1 - U_1^T U_1)$

If we set the derivative equal to 0 $\Rightarrow S U_1 = \lambda_1 U_1 \Rightarrow U_1^T S U_1 = \lambda_1$

S is the eigenvalue associated to the eigenvalue λ_1 . S is called first principal component.

2). Yes, because the intrinsic dimensionality of the dataset is 3.

Question 4

$$\text{Layer 1: } |\mathcal{N}_1| = 5 \times 5 \times 3 \times 16 + 16$$

$$\text{Layer 2: } |\mathcal{N}_2| = 3 \times 3 \times 16 \times 32 + 32$$

$$\text{Layer 3: } |\mathcal{N}_3| = 5 \times 3 \times 32 \times 64 + 64$$

$$\text{FC Layer: } |\mathcal{N}_{\text{FC}}| = 200 \times 10$$

A suitable loss function is categorical cross-entropy:

$$J(\mathcal{N}) = -\sum_n \ln p(t|x, \mathcal{N}) \quad \text{If we assume additive gaussian noise:}$$

$$p(t|x, \mathcal{N}) = N(t | f(x, \mathcal{N}), \beta^{-1}) \Rightarrow J(\mathcal{N}) = \sum_n \frac{1}{2} \|t - f(x, \mathcal{N})\|^2$$

Question 5

1). The model can be defined as: $y(x, w) = \sum_{n=1}^N w_n \phi_n(x) = w^T \phi(x)$. We know that the target value is $t = y(x, w) + \varepsilon$ with ε additive gaussian noise. If the noise is gaussian $p(\varepsilon|\beta) = N(\varepsilon|0, \beta^{-1}) \Rightarrow p(t|x_1, \dots, x_N, w, \beta) = N(t|y(x, w), \beta^{-1})$.

$$\text{Now we consider the iid hypothesis: } p(t_1, \dots, t_N|x_1, \dots, x_N, w, \beta) = \prod_{n=1}^N N(t_n|w^T \phi(x_n), \beta^{-1}) =$$

$$= \sum_{n=1}^N \ln N(t_n | w^\top \phi(x_n), \beta^{-1}) = -\beta \cdot \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2 - \frac{N}{2} \ln(2\pi\beta^{-1})$$

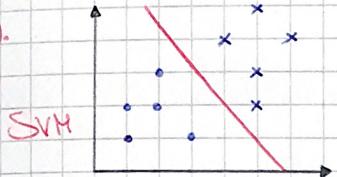
2). We consider the error function $E_\theta(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^\top \phi(x_n))^2$.

We can use a sequential learning algorithm to find $w^* = \arg \min_w E_\theta(w)$.

$$\hat{w} \leftarrow \hat{w} + \eta [t_n - w^\top \phi(x_n)] \phi(x_n)$$

Question 6

1).



PERCEPTRON

2). The main difference is that SVM aims at maximum margin with the better accuracy, while perceptron no. Perceptron stop the iterations when finds a possible separation surface simply summing the vector related to the separation surface with the feature vector of the misclassified point. For this reason very often we find a separation surface that is very close to the points. SVM try to maximize this margin and for this reason we have the plotted separation surface.

3). I prefer SVM because, using the principle of maximum margin, we have a very low probability of misclassification, while in perceptron no. Another advantage of SVM is that we can use it also if the dataset is not perfectly linearly separable.