

Question 1

1). The gram matrix is $U = XX^T$ that represents the inner product between the input vector and the transpose of the input vector. If we use a kernel function.

$$k(x, x') = x^T x' \rightarrow y(n) = \sum_{n=1}^N d_n x_n^T x \text{ and } K = \begin{pmatrix} x_1^T x_1 & \dots & x_1^T x_N \\ \vdots & \ddots & \vdots \\ x_N^T x_1 & \dots & x_N^T x_N \end{pmatrix}$$

$$\text{If we have a generic kernel function } k(x, x') \rightarrow y(x) = \sum_{n=1}^N d_n k(x_n, x)$$

$$\text{and } K = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_N) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \dots & k(x_N, x_N) \end{pmatrix}$$

2). It's possible to obtain a kernelized version of regression. The error function

$$\text{becomes } J(w) = \sum_{n=1}^N E(y_n, t_n) + \lambda \|w\|^2 \text{ with } E(y_n, t_n) = (y_n - t_n)^2$$

We can use the equations defined previously and we obtain: $y(x, w^*) = \sum_{n=1}^N d_n k(x_n, x)$

$$\text{We apply the kernel trick} \rightarrow y(n, w^*) = \sum_{n=1}^N d_n k(x_n, x)$$

This is really computationally expensive, in fact requires $|D|^2$ matrix multiplications.

Question 2

$$1). \dim(W_1) = 128 \times 10 = 1280 \quad \dim(W_2) = 10 \times 50 = 500$$

$$2). h = g(W^T x + c) \text{ with } g(z) = \max(0, z) \text{ ReLU function}$$

$$y(n, W) = w^T \max(0, W^T m + c) + b$$

Question 3

1). The goal of an autoencoder is dimensionality reduction. An autoencoder is composed by 2 neural networks, an encoder and a decoder. The structure is the following:



In autoencoder usually we have hidden layers with reduced size that are also called bottlenecks. The training is based on reconstruction loss.

Given a dataset $\{x_n\}$, autoencoders are trained with the sample x_n in input and in output.

1). We compute the K nearest neighbors of a new sample $x \notin D$ and then we assign to x the most common label among the majority of neighbors.

$$P(c|x_n, D, K) = \frac{1}{K} \sum_{x_n \in N_K(x_n, D)} I(t_n = c) \quad \text{with } I(c) = \begin{cases} 1 & \text{if } c \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

$N_K(x_n, D)$ K nearest neighbors of x_n .

- 2).
-
- $k=3$ The sample will be classified as +.

Question 6

- 1). We can use boosting, a simple ensemble method. We split our dataset D into M bootstrapped datasets D_1, \dots, D_M and we train each model $y_i(x)$ using the dataset $D_i \quad \forall i = 1 \dots M$

$$y_{\text{ensemble}}(x) = \frac{1}{M} \sum_{k=1}^M y_k(x)$$

- 2). No, I think that there are no properties that have to be satisfied.