

# Homework 1

## Deep Learning

Duarte Calado de Almeida  
95565

André Lopes Rodrigues  
96576

### Question 1

### Question 2

### Question 3

1. Let  $x_i$  denote the  $i$ -th component of  $\mathbf{x}$  and let  $W_{ij}$  denote the entry in the  $i$ -th row and  $j$ -th column in matrix  $\mathbf{W}$ . We then have that:

$$\begin{aligned}
 h_i(\mathbf{x}) &= g\left(\sum_{j=1}^D W_{ij}x_j\right) = \left(\sum_{j=1}^D W_{ij}x_j\right)^2 = \left(\sum_{j=1}^D W_{ij}x_j\right)\left(\sum_{k=1}^D W_{ik}x_k\right) \\
 &= \sum_{j=1}^D \sum_{k=1}^D W_{ij}x_j W_{ik}x_k = \sum_{j=1}^D W_{ij}^2 x_j^2 + \sum_{j=1}^D \sum_{k \neq j} W_{ij} W_{ik} x_j x_k \\
 &= \sum_{j=1}^D W_{ij}^2 x_j^2 + \sum_{j=1}^D \sum_{k=1}^{j-1} W_{ij} W_{ik} x_j x_k + \sum_{j=1}^D \sum_{k=j+1}^D W_{ij} W_{ik} x_j x_k \\
 &= \sum_{j=1}^D W_{ij}^2 x_j^2 + \sum_{k=1}^D \sum_{j=k+1}^D W_{ij} W_{ik} x_j x_k + \sum_{j=1}^D \sum_{k=j+1}^D W_{ij} W_{ik} x_j x_k \\
 &= \sum_{j=1}^D W_{ij}^2 x_j^2 + \sum_{k=1}^D \sum_{j=k+1}^D W_{ij} W_{ik} (2x_j x_k) \\
 &= \begin{bmatrix} W_{i1}^2 & W_{i1}W_{i2} & \dots & W_{i1}W_{iD} & W_{i2}^2 & W_{i2}W_{i3} & \dots & W_{i(D-1)}^2 & W_{i(D-1)}W_{iD} & W_{iD}^2 \end{bmatrix} \begin{bmatrix} x_1^2 \\ 2x_1x_2 \\ \dots \\ 2x_1x_D \\ x_2^2 \\ 2x_2x_3 \\ \dots \\ x_{D-1}^2 \\ 2x_{D-1}x_D \\ x_D^2 \end{bmatrix}
 \end{aligned}$$

As such,  $\mathbf{h}$  is linear in some feature transformation  $\phi$ , that is,  $\mathbf{h}$  can be written as  $\mathbf{A}_\Theta \phi(\mathbf{x})$ . In particular, we have that such matrix  $\mathbf{A}_\Theta$  can be defined as:

$$\mathbf{A}_\Theta = \begin{bmatrix} -\mathbf{a}_1^T & - \\ -\mathbf{a}_2^T & - \\ \vdots & \\ -\mathbf{a}_D^T & - \end{bmatrix}$$

where

$$\mathbf{a}_i = \begin{bmatrix} W_{i1}^2 & W_{i1}W_{i2} & \dots & W_{i1}W_{iD} & W_{i2}^2 & W_{i2}W_{i3} & \dots & W_{i(D-1)}^2 & W_{i(D-1)}W_{iD} & W_{iD}^2 \end{bmatrix}^T$$

and we can define the feature transformation  $\phi$  as:

$$\phi(\mathbf{x}) = (x_1^2, 2x_1x_2, \dots, 2x_1x_D, x_2^2, 2x_2x_3, \dots, x_{D-1}^2, 2x_{D-1}x_D, x_D^2)$$

2. Given that the predicted output  $\hat{y}$  is defined as:

$$\hat{y} = \mathbf{v}^T \mathbf{h}$$

the linearity of  $\mathbf{h}$  in the feature transformation  $\phi(\mathbf{x})$  proven above leads to following equality:

$$\hat{y} = \mathbf{v}^T \mathbf{A}_\Theta \phi(\mathbf{x}) = (\mathbf{A}_\Theta^T \mathbf{v})^T \phi(\mathbf{x}) = \mathbf{c}_\Theta^T \phi(\mathbf{x})$$

where we take  $\mathbf{c}_\Theta$  to be equal to  $\mathbf{A}_\Theta^T \mathbf{v}$ , thereby proving that  $\hat{y}$  is also a linear transformation of  $\phi(\mathbf{x})$ . Nevertheless,  $\hat{y}$  is **not** linear in terms of the original parameters  $\Theta$ . To see this, note that the model is now a linear combination of **products** of entries of  $\mathbf{W}$  and  $\mathbf{v}$  rather than the entries by themselves:

$$\hat{y} = \mathbf{v}^T \mathbf{A}_\Theta \phi(\mathbf{x}) = \sum_{i=1}^D v_i (w_i^T \mathbf{x})^2 = \sum_{i=1}^K \sum_{j=1}^D \sum_{k=1}^D v_i W_{ij} W_{ik} x_i x_k$$

where we define  $w_i$  to be the  $i$ -th row vector of matrix  $\mathbf{W}$ .

3. To prove the desired result, for  $\mathbf{c}_\Theta$  defined in the previous subquestion and for any  $\mathbf{c} \in \mathbb{R}^{\frac{D(D+1)}{2}}$ , we make the observation that the inner products  $\mathbf{c}_\Theta^T \phi(\mathbf{x})$  and  $\mathbf{c}^T \phi(\mathbf{x})$  actually correspond to quadratic forms in  $\mathbf{x}$ :

$$\begin{aligned} \mathbf{c}_\Theta^T \phi(\mathbf{x}) &= \sum_{i=1}^K v_i (\mathbf{A}_\Theta \phi(\mathbf{x}))_i^2 = \sum_{i=1}^K v_i (w_i^T \mathbf{x})^2 \\ &= \begin{bmatrix} w_1^T \mathbf{x} & w_2^T \mathbf{x} & \dots & w_K^T \mathbf{x} \end{bmatrix} \text{diag}(\mathbf{v}) \begin{bmatrix} w_1^T \mathbf{x} \\ w_2^T \mathbf{x} \\ \vdots \\ w_K^T \mathbf{x} \end{bmatrix} \\ &= (\mathbf{W}\mathbf{x})^T \text{diag}(\mathbf{v}) \mathbf{W}\mathbf{x} = \mathbf{x}^T \mathbf{W}^T \text{diag}(\mathbf{v}) \mathbf{W}\mathbf{x} \end{aligned}$$

and

$$\begin{aligned} \mathbf{c}^T \phi(\mathbf{x}) &= \sum_{i=1}^{\frac{D(D+1)}{2}} c_i \phi_i(\mathbf{x}) = \sum_{i=1}^D c_{(i-1)D+i} x_i^2 + \sum_{i=1}^D \sum_{j=i+1}^D c_{(i-1)D+j} (2x_i x_j) \\ &= \mathbf{x}^T \mathcal{M}(\mathbf{c}) \mathbf{x} \end{aligned}$$

where  $\mathcal{M}(\mathbf{c}) \in \mathbb{R}^{D \times D}$  is a symmetric matrix obtained through  $\mathbf{c}$  such that:

- the diagonal and the part above the diagonal of the matrix  $\mathcal{M}(\mathbf{c})$  is filled row-wise with the elements of vector  $\mathbf{c}$ , i.e.:

$$\mathcal{M}(\mathbf{c}) = \begin{bmatrix} c_1 & c_2 & c_3 & \dots & c_D \\ c_2 & c_{D+1} & c_{D+2} & \dots & c_{2D-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ c_D & c_{D-1} & c_{D-2} & \dots & c_{\frac{D(D+1)}{2}} \end{bmatrix}$$

- for  $1 \leq i \leq D$  and  $j < i$ ,  $(\mathcal{M}(\mathbf{c}))_{ij} = (\mathcal{M}(\mathbf{c}))_{ji}$

Furthermore, we also recurr to the following lemma:

**Lemma 1.** *Two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{\frac{D(D+1)}{2}}$  are equal if and only if  $\mathbf{a}^T \phi(\mathbf{x}) = \mathbf{b}^T \phi(\mathbf{x})$ , for all  $\mathbf{x} \in \mathbb{R}^D$*

*Proof.* If  $\mathbf{a} = \mathbf{b}$ , then  $\mathbf{a}^T \phi(\mathbf{x}) = \mathbf{b}^T \phi(\mathbf{x})$  is trivially verified. For the reverse implication, note that, according to the previously made observation, for any  $\mathbf{x} \in \mathbb{R}^D$ :

$$\mathbf{a}^T \phi(\mathbf{x}) = \mathbf{b}^T \phi(\mathbf{x}) \Rightarrow \mathbf{x}^T \mathcal{M}(\mathbf{a}) \mathbf{x} = \mathbf{x}^T \mathcal{M}(\mathbf{b}) \mathbf{x}$$

Since both  $\mathcal{M}(\mathbf{a})$  and  $\mathcal{M}(\mathbf{b})$  are symmetric and the associated quadratic forms are twice differentiable continuous, taking the hessian on both sides of the equation yields:

$$\mathcal{M}(\mathbf{a}) = \mathcal{M}(\mathbf{b})$$

that is,  $\mathcal{M}(\mathbf{a})$  and  $\mathcal{M}(\mathbf{b})$  are entrywise equal. In particular, we have that  $a_i = b_i$ , for  $i = 1, \dots, \frac{D(D+1)}{2}$ .  $\square$

We are now equipped with the tools needed for the proof. Since  $\mathcal{M}(\mathbf{c})$  is symmetric, the **Spectral Decomposition Theorem** tells us that there is a matrix  $\mathbf{Q}$  orthonormal and  $\mathbf{\Lambda}$  diagonal such that  $\mathcal{M}(\mathbf{c}) = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T$ . Let  $\mathbf{q}_i$  denote the eigenvector of  $\mathcal{M}(\mathbf{c})$  that is present in  $i$ -th column of  $\mathbf{Q}$  and let  $\lambda_i$  be the corresponding eigenvalue (note that  $\{\mathbf{q}_i\}_{i=1}^D$  forms an orthonormal basis of  $\mathbb{R}^D$ ). Then, we can write  $\mathbf{c}^T \phi(\mathbf{x})$  as:

$$\mathbf{c}^T \phi(\mathbf{x}) = \mathbf{x}^T \mathcal{M}(\mathbf{c}) \mathbf{x} = (\mathbf{Q}^T \mathbf{x})^T \mathbf{\Lambda} (\mathbf{Q}^T \mathbf{x}) = \sum_{i=1}^D \lambda_i (\mathbf{q}_i^T \mathbf{x})^2$$

Now, if we assume that  $K \geq D$ , then we can construct the matrix  $\mathbf{W}$  and vector  $\mathbf{v}$  that make  $\mathbf{c}_{\Theta}^T \phi(\mathbf{x})$  equal to  $\mathbf{c}^T \phi(\mathbf{x})$  in the following way:

- we make  $\mathbf{v}$  to be equal to  $(\lambda_1, \lambda_2, \dots, \lambda_D, \underbrace{0, \dots, 0}_{K-D \text{ times}})$ ;
- we make  $\mathbf{W}$  to be equal to the vertical concatenation of  $\mathbf{Q}^T$  with a  $(K-D) \times D$  matrix of zeros, i.e.:

$$\mathbf{W} = \begin{bmatrix} \mathbf{Q}^T \\ \mathbf{0}_{(K-D) \times D} \end{bmatrix}$$

We then have that:

$$\mathbf{c}_{\Theta}^T \phi(\mathbf{x}) = \sum_{i=1}^K \lambda_i (\mathbf{w}_i^T \mathbf{x})^2 = \sum_{i=1}^D \lambda_i (\mathbf{w}_i^T \mathbf{x})^2 = (\mathbf{Q}^T \mathbf{x})^T \mathbf{\Lambda} (\mathbf{Q}^T \mathbf{x}) = \mathbf{x}^T \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T \mathbf{x} = \mathbf{c}^T \phi(\mathbf{x})$$

and, by Lemma 1, we prove that the previous choice of  $\mathbf{W}$  and  $\mathbf{v}$  originate a vector  $\mathbf{c}_{\Theta}$  such that  $\mathbf{c}_{\Theta} = \mathbf{c}$   $\square$

In fact, we can relax the requirement  $K \geq D$  to be  $K \geq D - N$ , where  $N$  is the dimension of the nullspace of  $\mathcal{M}(\mathbf{c})$ . Since  $\mathcal{M}(\mathbf{c})$  is symmetric and thus diagonalizable, the set of indices  $\mathcal{I}$  such that the eigenvectors  $\{\mathbf{q}_i | i \in \mathcal{I}\}$  are not associated with eigenvalue zero has  $D - N$  elements. As such:

$$\mathbf{c}^T \phi(\mathbf{x}) = \sum_{i=1}^D \lambda_i (\mathbf{q}_i^T \mathbf{x})^2 = \sum_{i \in \mathcal{I}} \lambda_i (\mathbf{q}_i^T \mathbf{x})^2$$

and thus we can take  $\mathbf{W}$  to be equal to the concatenation of the matrix  $\tilde{\mathbf{Q}}^T$  with a matrix of  $N \times K$  zeros (where  $\tilde{\mathbf{Q}} \in \mathbb{R}^{D \times (D-N)}$  and  $\tilde{\mathbf{q}}_k = \mathbf{q}_{i_k}$ ,  $1 \leq k \leq D - N$ ). Following the previously made argument, we would obtain again a vector  $\mathbf{c}_\Theta$  equal to  $\mathbf{c}$ .

Now, if  $K < D - N$ , then the nullspace of  $\mathbf{W}$  has at least dimension  $D - (D - N - 1) = N + 1$ . Since the rank of  $\mathcal{M}(\mathbf{c})$  is  $D - N$ , the dimensions of the row space of  $\mathcal{M}(\mathbf{c})$  and the nullspace of  $\mathbf{W}$  sum up to  $D + 1$  and thus there exists some non-null vector  $\mathbf{x}^*$  that is in the rowspace of  $\mathcal{M}(\mathbf{c})$  and in the nullspace of  $\mathbf{W}$ . Let  $\mathbf{y}^* = \mathbf{Q}^T \mathbf{x}^*$  and choose  $\mathbf{c}$  to be a vector such that  $\mathcal{M}(\mathbf{c})$  is positive semidefinite, namely:

$$\mathcal{M}(\mathbf{c}) = \text{diag}(\underbrace{1, \dots, 1}_{D-N \text{ times}}, \underbrace{0, \dots, 0}_{N \text{ times}})$$

we have:

$$\mathbf{c}_\Theta^T \phi(\mathbf{x}^*) = \mathbf{x}^{*T} \mathbf{W}^T \text{diag}(\mathbf{v}) \mathbf{W} \mathbf{x}^* = \mathbf{0}$$

and

$$\mathbf{c}^T \phi(\mathbf{x}^*) = \sum_{i \in \mathcal{I}} \lambda_i (\mathbf{q}_i^T \mathbf{x}^*)^2 > 0$$

since  $\lambda_i > 0$  ( $i \in \mathcal{I}$ ), and  $(\mathbf{q}_i^T \mathbf{x}^*)^2 > 0$  for at least one  $i$ , as  $\mathbf{x}^*$  belongs to the rowspace of  $\mathcal{M}(\mathbf{c})$ , and, consequently, it does not belong to its nullspace.

We have thus constructed an instance where  $K < D$  and there is no choice of parameters  $\mathbf{W}$  and  $\mathbf{v}$  that make  $\mathbf{c}_\Theta^T$  equal to  $\mathbf{c}$ .

4.