



INSTITUTO SUPERIOR TÉCNICO

DEEP LEARNING

2<sup>nd</sup> QUARTER - 2022/2023

---

## Homework 1

---

*Group 46:*

Santiago QUINTAS - 93179  
Alina KLING - 105208

*Professor:*

Mário FIGUEIREDO

23<sup>st</sup> December, 2022

## Declaration of Contribution

We worked on all of the Questions together. In Question 1 we did together in a Discord Call, the other Questions we worked in separate and compared the results. In Question 2 Santiago was leading the process of development, while Alina did so in Question 3.

### Question 1

1.

a)

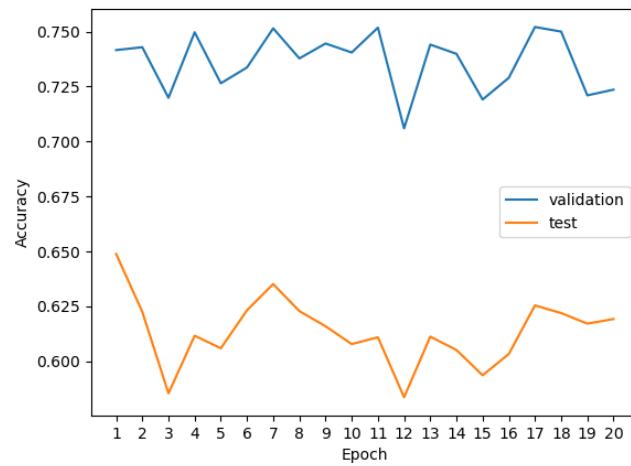


Figure 1: Accuracy of the perceptron as a function of the epoch number

b)

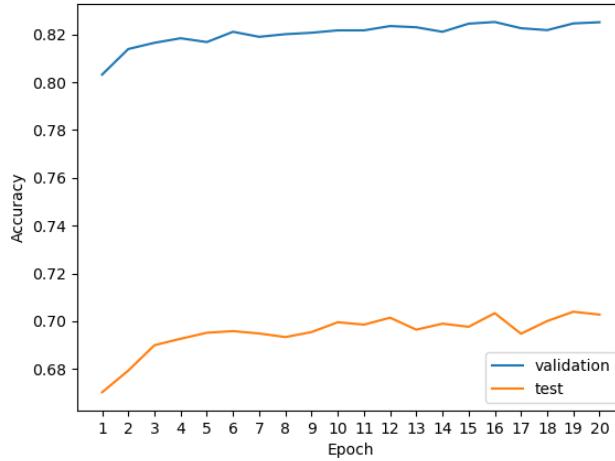


Figure 2: Accuracy of the logistic regression as a function of the epoch number

2.

a)

A multi layered perceptron, (mlp), allows the implementation of non-linear activation functions in its hidden layers, while a simple perceptron is limited to a linear activation function in its output. The main benefit of a mlp is having the capability to learn more complex patterns of data with non linear relationships while a simple perceptron is only useful in data with linear relationships. In this particular task our data is not linearly separable, the mlp will be able to overcome this issue with non linear activation functions.

This advantage goes away if the mlp is implemented with only linear activation functions, the mlp becomes a more complex equivalent of a simple perceptron, the output of a hidden layer is the input of the next one, and if they all possess linear activation functions then the output of the perceptron will be a linear relationship between the weights and biases of the previous layers, which removes the expressiveness necessary for learning more complex patterns.

b)

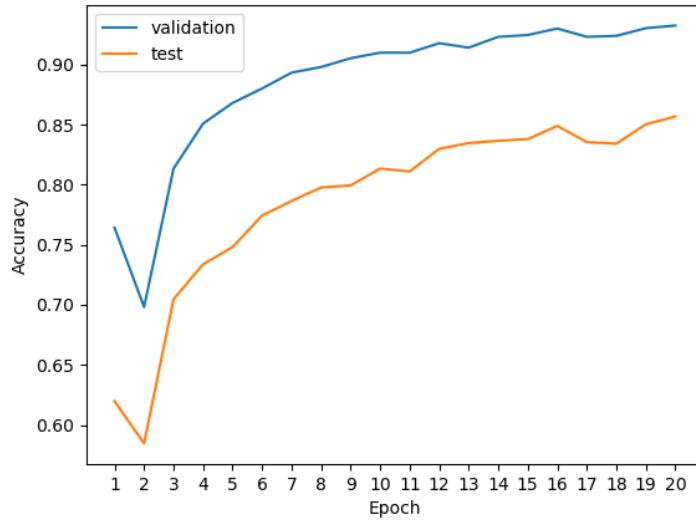


Figure 3: Accuracy of the multi layered perceptron (mlp) as a function of the epoch number

## Question 2

### 1.

The default parameters used in the implementation of the logistic regression linear model were the following:

Number of Epochs	20
Learning Rate	0.01
Batch Size	1

Table 1: Default parameters 1

For different learning rates we achieved different accuracies, shown in the table 2. The best configuration, according to the accuracy test is with a learning rate of 0.001

Learning Rate	Final test accuracy
0.1	0.6133
0.01	0.6806
0.001	0.7019

Table 2: Accuracy test for different learning rates

Below is the plot of both the training loss and accuracy over the number of epochs of the changes for a learning rate of 0.001, which corresponds to our best configuration.

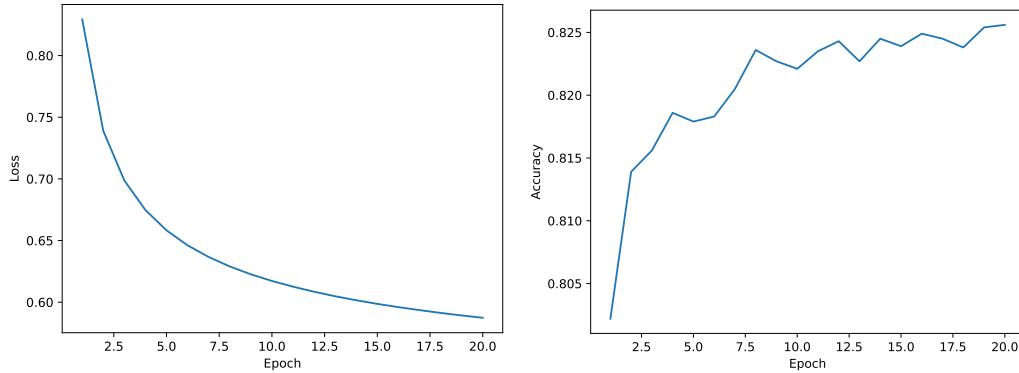


Figure 4: Training loss and accuracy of the logistic regression for a learning rate of 0.001

## 2.

The default parameters used in the implementation of the neural network with a single layer were the following:

Number of Epochs	20
Learning Rate	0.01
Hidden Size	100
Dropout	0.3
Batch Size	16
Activation	ReLU
Optimizer	SGD

Table 3: Default parameters 2

For different parameters the accuracy test ended with these results shown in table 4, keep in mind that only 1 parameter was changed at a time, the rest remained in their default setting. The best configuration was the one with a hidden size of 200.

Parameter	Value	Final test accuracy
Learning Rate	0.1	0.8714
Learning Rate	0.01	0.8593
Learning Rate	0.001	0.7450
Hidden Size	200	0.8811
Dropout	0.5	0.8413
Activation	tanh	0.8257

Table 4: Accuracy test for different parameters

Below is the plot of both the training loss and accuracy over the number of epochs for a hidden size of 200, which corresponds to our best configuration.

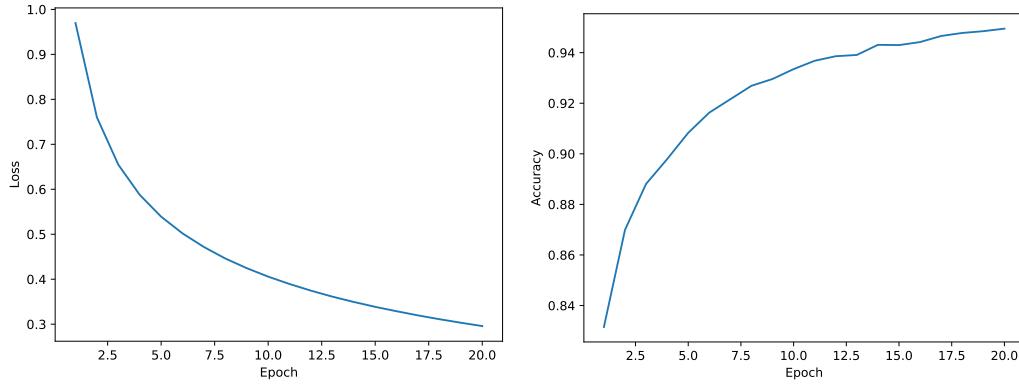


Figure 5: Training loss and accuracy of the mlp for a hidden size of 200

### 3.

Using the same default parameters as before we update the number of layers in our mlp and the results obtained are shown in the table 5, the best configuration was the one with 2 layers.

Number of Layers	Final test accuracy
1 (default)	0.8593
2	0.8698
3	0.8652

Table 5: Accuracy test for different learning rates

Below is the plot of both the training loss and accuracy over the number of epochs for a mlp with 2 layers, which corresponds to our best configuration.

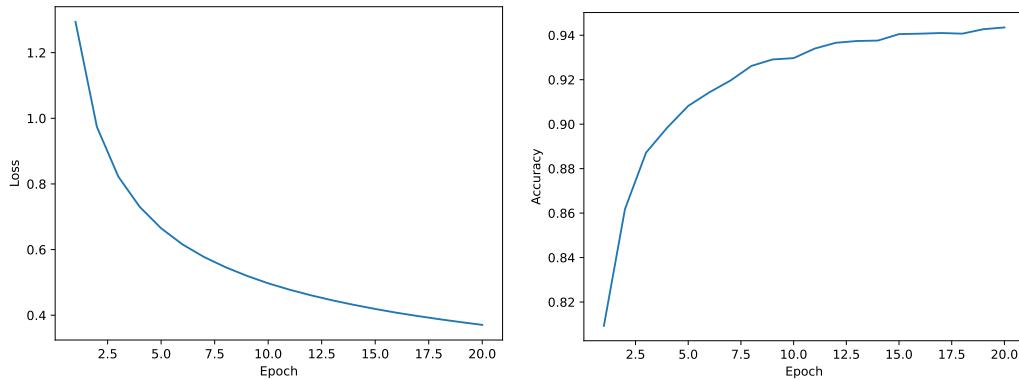


Figure 6: Training loss and accuracy of the mlp using 2 layers

# Homework 1

## Question 3 (30 points)

**Multi-layer perceptron with quadratic activations.** In this exercise, we will consider a feed-forward neural network with a single hidden layer and a quadratic activation function,  $g(z) = z^2$ . We will see under some assumptions, this choice of activation, unlike other popular activation functions such as tanh, sigmoid, or relu, can be tackled as a linear model via a reparametrization.

We assume a univariate regression task, where the predicted output  $\hat{y} \in \mathbb{R}$  is given by  $\hat{y} = \mathbf{v}^\top \mathbf{h}$ , where  $\mathbf{h} \in \mathbb{R}^K$  are internal representations, given by  $\mathbf{h} = g(\mathbf{W}\mathbf{x})$ ,  $\mathbf{x} \in \mathbb{R}^D$  is a vector of input variables, and  $\Theta = (\mathbf{W}, \mathbf{v}) \in \mathbb{R}^{K \times D} \times \mathbb{R}^K$  are the model parameters.

- **feed-forward NN:** 1 single hidden layer  
quadratic activation function  $g(z) = z^2$   
↳ with assumption with reparametrization  $\approx$  lin. model
- **Assumptions :** univariate regression task  
predicted output  $\hat{y} \in \mathbb{R}$  :  $\hat{y} = \mathbf{v}^\top \mathbf{h}$   
internal representations  $\mathbf{h} \in \mathbb{R}^K$  :  $\mathbf{h} = \tilde{g}(\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}})$   
vector of input variables  $\tilde{\mathbf{x}} \in \mathbb{R}^D$   
model parameters  $\Theta = (\tilde{\mathbf{w}}, \tilde{\mathbf{v}}) \in \mathbb{R}^{K \times D} \times \mathbb{R}^K$  :  $\Theta = (\tilde{\mathbf{w}}, \tilde{\mathbf{v}})$

1. (10 points) Show that we can write  $\mathbf{h} = \mathbf{A}_\Theta \phi(\mathbf{x})$  for a certain feature transformation  $\phi: \mathbb{R}^D \rightarrow \mathbb{R}^{\frac{D(D+1)}{2}}$  independent of  $\Theta$  and  $\mathbf{A}_\Theta \in \mathbb{R}^{K \times \frac{D(D+1)}{2}}$ . That is,  $\mathbf{h}$  is a linear transformation of  $\phi(\mathbf{x})$ . Determine the mapping  $\phi$  and the matrix  $\mathbf{A}_\Theta$ .

a)  $\mathbf{h} = \mathbf{A}_\Theta \Phi(\mathbf{x})$  for  $\Phi: \mathbb{R}^D \rightarrow \mathbb{R}^{\frac{D(D+1)}{2}}$  independent of  $\Theta$  and  $\Phi \in \mathbb{R}^{K \times \frac{D(D+1)}{2}}$

Goal:  $\mathbf{h} = g(\mathbf{W}\mathbf{x}) \doteq \mathbf{A}_\Theta \cdot \Phi(\mathbf{x})$

$\downarrow$  Matrix       $\downarrow$  Vector  $\frac{D(D+1)}{2}$   
 (depends only on  $\Theta$ )      (depends only on the input  $\mathbf{x}$ )

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 1 & 3 \\ 3 & 3 & 2 \\ 4 & 1 & 2 \end{bmatrix} \cdot \begin{bmatrix} 2 & 1 & 3 \\ 3 & 3 & 2 \\ 1 & 1 & 2 \end{bmatrix} = \begin{bmatrix} 1 \cdot 2 + 2 \cdot 3 + 3 \cdot 4 \\ 1 \cdot 1 + 2 \cdot 3 + 3 \cdot 1 \\ 1 \cdot 3 + 2 \cdot 2 + 3 \cdot 2 \end{bmatrix} = \begin{bmatrix} 20 \\ 10 \\ 13 \end{bmatrix}$$

$\overbrace{1 \times 3}^{\text{Matrix}} \quad \overbrace{3 \times 3}^{\text{Vector}} \quad \overbrace{1 \times 3}^{\text{Result}}$

(I) Use  $g(z) = z^2$

$$g(\mathbf{W}\mathbf{x}) = g\left(\sum_{i=1}^D w_{ij} x_i\right) = \left(\sum_{i=1}^D w_{ij} x_i\right)^2 = \sum_{l=1}^D \sum_{i=1}^D w_{lj} x_i w_{ij} x_i$$

→ Main idea: Show that the above sum can be rewritten as a multiplication

of a matrix (pure of  $\theta$ ) and a pure vector of  $x$

(II) To understand summation: Example with  $D=3$  and  $K=2$

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{D(D+1)/2}$$

↑  
with same index  
↑ to eliminate  
doubts!  
with another index

(Number of permutations for  $D=3$ :  $n^3 \rightarrow R \frac{3!4!}{2} = n^6$ )  
11, 12, 13  
21, 22, 23

$$W \in \mathbb{R}^{K \times D} \rightarrow W \neq 2 \times 3 \text{ Matrix}$$

$$x \in \mathbb{R}^D \rightarrow x : 3 \times 1 \text{ vector}$$

calculate  $(Wx)^2$

$$Wx = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix}, \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} x_i w_{1i} \\ x_i w_{2i} \end{pmatrix} \text{ with } i \in \{1, 2, 3\}$$

$$(Wx)^2 = \begin{pmatrix} (x_i w_{1i})^2 \\ (x_i w_{2i})^2 \end{pmatrix} = \begin{pmatrix} (x_1 w_{11} + x_2 w_{12} + x_3 w_{13})^2 \\ (x_1 w_{21} + x_2 w_{22} + x_3 w_{23})^2 \end{pmatrix}$$

$$= \begin{pmatrix} (x_1 w_{11})^2 + (x_2 w_{12})^2 + (x_3 w_{13})^2 + 2x_1 x_2 w_{11} w_{12} + 2x_1 x_3 w_{11} w_{13} + 2x_2 x_3 w_{12} w_{13} \\ (x_1 w_{21})^2 + (x_2 w_{22})^2 + (x_3 w_{23})^2 + 2x_1 x_2 w_{21} w_{22} + 2x_1 x_3 w_{21} w_{23} + 2x_2 x_3 w_{22} w_{23} \end{pmatrix}$$

$$\Rightarrow \phi(x) : \mathbb{R}^D \rightarrow \mathbb{R}^{D(D+1)/2}$$

$$(x_1, \dots, x_D) \rightarrow \underbrace{(x_1^2, x_2^2, \dots, x_1 x_2, x_1 x_3, \dots)}_{D \text{-elements}} \underbrace{(x_2^2, \dots, x_1 x_2, x_2 x_3, \dots)}_{D(D+1)/2 \text{-elements}}$$

Try to rewrite this now to  $(Wx)$   $\begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix}$

$$\begin{matrix} j=1 \\ j=2 \end{matrix} \left( \begin{array}{cccccc} w_{11}^2 & w_{12}^2 & w_{13}^2 & 2w_{11}w_{12} & 2w_{13}w_{11} & 2w_{12}w_{13} \\ w_{21}^2 & w_{22}^2 & w_{23}^2 & 2w_{21}w_{22} & 2w_{23}w_{21} & 2w_{22}w_{23} \end{array} \right) \left( \begin{array}{c} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{array} \right)$$

A  $\theta$  : depends only on parameters  $W$

$\phi(x)$ : depends on  $\theta$   
on input  $x$

Can write the  $\phi(x)$  as:

$$\phi(x) = [x_1 x_1, x_1 x_2, \dots, x_2 x_1, x_2 x_3, \dots, x_{D-1} x_D, x_D^2]$$

all possible combinations of  $x_i x_j$

$$\phi(x) = \sum_{j=1}^D \sum_{i=1}^D x_i x_j$$

### III) Write down generalized version of II)

$$h \text{ for each row } j \text{ of the } b \times D \text{ matrix } A_\theta \text{ is given by:}$$

$$h_j = \sum_{i=1}^D \sum_{l=1}^D (w_{jl} w_{li}) (x_i x_l) = \sum_{i=1}^D w_{ji}^2 x_i^2 + 2 \sum_{i=1}^D \sum_{l=i+1}^D x_i x_l w_{jl} w_{li}$$

D-terms (quadratic)      mapping  $\frac{D(D-1)}{2}$

redefine the indices of the sums and re-write them as one sum over m

with  $m = \{(ij)\}_{ij}$  all combinations of ij (without repetitions)

$\begin{matrix} \downarrow D \text{ values} \\ \downarrow D \text{ values} \end{matrix}$

$$\Rightarrow m = \left\{ 1, \dots, \frac{D(D+1)}{2} \right\}$$

for the  $j$ th row of  $A_\theta$ :

$$h_j = \sum_{m=1}^{D(D+1)/2} A_{jm} \Phi_m \Rightarrow A_\theta \Phi(x) \quad \checkmark$$

$\Rightarrow$  so it was shown, that h is a linear transformation  $\Phi(x)$  for the feature transformation with mapping

$$\Phi[x_1, \dots, x_D] \rightarrow \Phi[x_1 x_1, x_1 x_2, \dots, x_1 x_D, x_2 x_2, \dots, x_{D-1} x_D, x_D^2]$$

and that the matrix  $A_\theta$  has the j rows

$$\sum_{j=1}^D \sum_{l=1}^D w_{jl} w_{li}$$

### III Alternative zu III above in Matrix-Form

$$h = g(WX) = (WX)^2 = \left( \sum_{j=1}^D w_{ij} x_j \right)^2$$

$$= \begin{pmatrix} (\sum_{j=1}^D w_{1j} x_j)^2 \\ \vdots \\ (\sum_{j=1}^D w_{Dj} x_j)^2 \end{pmatrix} = \begin{pmatrix} \sum_{j=1}^D w_{1j}^2 x_1^2 & + & 2 \sum_{i=1}^D \sum_{l=i+1}^D x_i x_l w_{1i} w_{li} \\ \vdots & & \vdots \\ \sum_{j=1}^D w_{Dj}^2 x_D^2 & + & 2 \sum_{i=1}^D \sum_{l=i+1}^D x_i x_l w_{Di} w_{Dl} \end{pmatrix}$$

D-Terms       $\frac{D(D-1)}{2}$  Terms

$$= \begin{pmatrix} w_{11}^2 & \dots & w_{1D}^2 & w_{11} w_{12} & \dots & w_{11} w_{D-1} w_{1D} \\ \vdots & & \vdots & \vdots & & \vdots \\ w_{D1}^2 & \dots & w_{DD}^2 & w_{D1} w_{D2} & \dots & w_{D1} w_{D-1} w_{D2} \end{pmatrix} \cdot \begin{pmatrix} x_1^2 \\ \vdots \\ x_D^2 \\ 2x_1 x_2 \\ \vdots \\ 2x_{D-1} x_D \end{pmatrix}$$

$=: A_\theta \in \mathbb{R}^{N \times D(D+1)/2}$

$=: \Phi(x) \in \mathbb{R}^{\frac{D(D+1)}{2}}$

$$\Rightarrow h = A_\theta \cdot \Phi(x) \quad (\text{so it's linear + we can separate})$$

2. (5 points) a) Based on the previous claim, show that  $\hat{y}$  is also a linear transformation of  $\phi(x)$ , i.e., we can write  $\hat{y}(x; \theta) = c^T \phi(x)$  for some  $c \in \mathbb{R}^{\frac{D(D+1)}{2}}$ . Does this mean this is a linear model in terms of the original parameters  $\theta$ ?

Idea:

$$\text{predicted output } \hat{y} \in \mathbb{R} : \hat{y} = v^T h$$

$$\text{internal representations } h \in \mathbb{R}^k : h = g(\tilde{w} \tilde{x})$$

$$\text{vector of input variables } \tilde{x} \in \mathbb{R}^D$$

$$\text{Model parameters } \theta = (\tilde{w}, \tilde{v}) \in \mathbb{R}^{k \times D} \times \mathbb{R}^k : \theta = (\tilde{w}, \tilde{v})$$

$\Rightarrow$  so should be straight forward: if  $h$  from above linear can just plug it in?

(I) Use definition of  $\hat{y}$  and plug  $h$  from 3.1. in

( $h$  from 3.1.  $h = A\phi \cdot \psi(x)$  with  $A\phi = \text{matrix with elements } A_{ij} = w_i \phi_j(x)$  and  $\psi(x) = \text{vector of monomials of } x \text{ up to 2nd order}$ )

$$\Rightarrow \hat{y} = v^T h = v^T A\phi \cdot \psi(x)$$

$$\Rightarrow c^T = v^T A\phi$$

$$\begin{array}{c} v^T \\ \hline 1 & \cdots & m \end{array} \cdot \begin{array}{c} A \\ \hline n \\ \hline 1 & \cdots & m \end{array} = \begin{array}{c} c^T \\ \hline 1 & \cdots & n \end{array}$$

(II) Therefore we can write  $\hat{y}$  as a linear transformation of  $\phi(x)$

b) Model = linear in terms of the original parameters  $\theta = (\tilde{w}, \tilde{v}) \in \mathbb{R}^{k \times D} \times \mathbb{R}^k$ ,

NOT necessarily, cause we can see in 3.1 (II) that  $A\phi$  and therefore also  $c^T$  which depends on  $A\phi$  is not linear in the parameters in the original parameter  $v$ .

(because the weights are multiplied with each other (e.g. squared) and multiplied with  $v$ )

3. (10 points) Assume  $K \geq D$ . Show that for any real vector  $\mathbf{c} \in \mathbb{R}^{\frac{D(D+1)}{2}}$  there is a choice of the original parameters  $\Theta = (\mathbf{W}, \mathbf{v})$  such that  $\mathbf{c}_\Theta = \mathbf{c}$ . That is, we can equivalently parametrize the model with  $\mathbf{c}_\Theta$  instead of  $\Theta$ . Does this mean this is a linear model in terms of  $\mathbf{c}_\Theta$ ? Show that an equivalent parametrization might not exist if  $K < D$ .

single answer

<https://cs229.stanford.edu/section/cs229-linalg.pdf>

In Q2, you show that for any  $\Theta$  there is a corresponding  $c$  that makes the network equivalent to a linear model.

[https://www.deeplearningbook.org/contents/linear\\_algebra.html](https://www.deeplearningbook.org/contents/linear_algebra.html)

In Q3, the goal is to show the reverse direction -- that for any given  $c$  there is a corresponding  $\Theta$  in the original model so that the two models are equivalent. This is the most difficult of that group of questions. The requirement of  $K \geq D$  will arise if you are on the right path.

What should be the bridge between that conclusion for  $D(D+1)/2$  and for  $D$ ? I've worked out a potential proof but I ended up reaching the same conclusion as my colleague. Thanks in advance.

helpful | 0



Mário Figueiredo 2 days ago

Please keep in mind that the relationship between the parameters  $\Theta$  and  $c$  is not linear. Thinking in terms of the number of equations versus the number of unknowns is only valid for linear systems of equations.

good comment | 0

- The relationship between  $y$  and  $x$  is certainly non-linear.
- In Q3.1 and Q3.2 you prove that there is a transformation  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{\frac{D(D+1)}{2}}$  such that  $y$  is a linear function of  $\phi(x)$ , that is, it can be written as  $y = c_\Theta^T \phi(x)$ .
- The parameter  $c_\Theta$  that appears in the previous paragraph is a non-linear function of the original network parameters  $\Theta = (\mathbf{W}, \mathbf{v})$ , that is, we can obtain it via some non-linear function  $\psi : \mathbb{R}^{K \times D} \times \mathbb{R}^K \rightarrow \mathbb{R}^{\frac{D(D+1)}{2}}$ , such that  $c_\Theta = \psi(\mathbf{W}, \mathbf{v})$ .
- Question Q3.3 can be stated as follows: (1) show that if  $K \geq D$ , then, for any  $c \in \mathbb{R}^{\frac{D(D+1)}{2}}$ , there exist a pair  $(\mathbf{W}, \mathbf{v})$  such that  $c = \psi(\mathbf{W}, \mathbf{v})$ .  
(2) Show that if  $K < D$ , there may exist some  $c \in \mathbb{R}^{\frac{D(D+1)}{2}}$ , such that there is no  $(\mathbf{W}, \mathbf{v})$  that satisfies  $c = \psi(\mathbf{W}, \mathbf{v})$ .

I hope this helps to clarify the question.

$$\mathbf{c}_\Theta^T = \mathbf{v}^T \mathbf{A}_\Theta \quad (1)$$

$$(\mathbf{c}_\Theta^T)^T = (\mathbf{v}^T \mathbf{A}_\Theta)^T$$

$$\mathbf{c}_\Theta = \mathbf{A}_\Theta \mathbf{v} = \mathbf{C} \quad \text{and} \quad \mathbf{C} = \psi(\mathbf{W}, \mathbf{v})$$

Model parameters  $\Theta = (\mathbf{W}, \mathbf{v}) \in \mathbb{R}^{K \times D} \times \mathbb{R}^K : \Theta = (\mathbf{W}, \mathbf{v})$

$$\hat{y}(x; \mathbf{c}_\Theta) = \mathbf{c}_\Theta^T \phi(x) \text{ for some } \mathbf{c}_\Theta \in \mathbb{R}^{\frac{D(D+1)}{2}}$$

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

$$\mathbf{c}_\Theta^T = \mathbf{v}^T \mathbf{A}_\Theta$$

can I divide the vector

$\mathbf{c}$  in Matrix-times Vector? → with right dimensions?

hmhhh???

Look at dimensions:

internal representations  $\mathbf{h} \in \mathbb{R}^K : \mathbf{h} = \tilde{g}(\tilde{\mathbf{w}} \cdot \tilde{\mathbf{x}})$

Vector of input variables  $\mathbf{x} \in \mathbb{R}^D$

Model parameters  $\Theta = (\mathbf{W}, \mathbf{v}) \in \mathbb{R}^{K \times D} \times \mathbb{R}^K : \Theta = (\mathbf{W}, \mathbf{v})$

if  $K \geq D$ : more internal representations than input variables

then? :

## • Look at definition of linearity:

### Deep-learning-book:

In order for the system  $\mathbf{Ax} = \mathbf{b}$  to have a solution for all values of  $\mathbf{b} \in \mathbb{R}^m$ , we therefore require that the column space of  $\mathbf{A}$  be all of  $\mathbb{R}^m$ . If any point in  $\mathbb{R}^m$  is excluded from the column space, that point is a potential value of  $\mathbf{b}$  that has no solution. The requirement that the column space of  $\mathbf{A}$  be all of  $\mathbb{R}^m$  implies immediately that  $\mathbf{A}$  must have at least  $m$  columns, that is,  $n \geq m$ . Otherwise, the dimensionality of the column space would be less than  $m$ . For example, consider a  $3 \times 2$  matrix. The target  $\mathbf{b}$  is 3-D, but  $\mathbf{x}$  is only 2-D, so modifying the value of  $\mathbf{x}$  at best enables us to trace out a 2-D plane within  $\mathbb{R}^3$ . The equation has a solution if and only if  $\mathbf{b}$  lies on that plane.

Having  $n \geq m$  is only a necessary condition for every point to have a solution. It is not a sufficient condition, because it is possible for some of the columns to be redundant. Consider a  $2 \times 2$  matrix where both of the columns are identical. This has the same column space as a  $2 \times 1$  matrix containing only one copy of the replicated column. In other words, the column space is still just a line and fails to encompass all of  $\mathbb{R}^2$ , even though there are two columns.

Formally, this kind of redundancy is known as **linear dependence**. A set of vectors is **linearly independent** if no vector in the set is a linear combination of the other vectors. If we add a vector to a set that is a linear combination of the other vectors in the set, the new vector does not add any points to the set's span. This means that for the column space of the matrix to encompass all of  $\mathbb{R}^m$ , the matrix must contain at least one set of  $m$  linearly independent columns. This condition is both necessary and sufficient for equation 2.11 to have a solution for every value of  $\mathbf{b}$ . Note that the requirement is for a set to have exactly  $m$  linearly independent columns, not at least  $m$ . No set of  $m$ -dimensional vectors can have more than  $m$  mutually linearly independent columns, but a matrix with more than  $m$  columns may have more than one such set.

I think we need this to  
ensure (column) exist?

↓  
no but it's not linear?  
 $\mathbf{w}^T \mathbf{A} \mathbf{x} = \sum_i w_i A_{i,j} x_j$



4. (5 points) Suppose we are given training data  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$  with  $N > \frac{D(D+1)}{2}$  and that we want to minimize the squared loss

$$L(\mathbf{c}_\Theta; \mathcal{D}) = \frac{1}{2} \sum_{n=1}^N (\hat{y}_n(\mathbf{x}_n; \mathbf{c}_\Theta) - y_n)^2.$$

Let the matrix  $\mathbf{X} \in \mathbb{R}^{N \times \frac{D(D+1)}{2}}$  have  $\phi(\mathbf{x}_n)$  as rows and assume that  $\mathbf{X}$  has full column-rank. Can we find a closed form solution  $\hat{\mathbf{c}}_\Theta$ ? Comment on the fact that global minimization is usually intractable for feedforward neural networks – what makes our problem special?

### (I) What is a closed solution: lecture

As described in lecture 2 in slide 16 the

closed form of a model, that is linear w.r.t the model parameters  $w$ :  $\hat{y} = \mathbf{w}^\top \Phi(\mathbf{x})$

with the feature vector  $\Phi(\mathbf{x})$  (e.g.  $\Phi(\mathbf{x}) = [1, x_1, x_1^2, \dots, x^D]^T \in \mathbb{R}^{D+1}$  bias term)

is given as

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \text{ with } \mathbf{X} = \begin{bmatrix} \Phi(\mathbf{x}_1)^T \\ \vdots \\ \Phi(\mathbf{x}_n)^T \\ \vdots \\ \Phi(\mathbf{x}_N)^T \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \\ \vdots \\ y_N \end{bmatrix}.$$

This solution can be found by minimizing the squared loss with

$$L(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}}) \text{ with } \hat{\mathbf{y}} = \mathbf{w}^\top \Phi(\mathbf{x})$$

with respect to  $w$ .

In order to do so write  $L$  in matrix-vector notation:

$$L(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{n=1}^N L(y_n, \hat{y}_n) = \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \Phi(\mathbf{x}_n))^2$$

$$= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

$$\text{Find minimum } \min_{\mathbf{w} \in \mathbb{R}^{D+1}} \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \Phi(\mathbf{x}_n))^2 :$$

$$0 = \nabla_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

$$= \nabla_{\mathbf{w}} \|\mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \|\mathbf{y}\|^2\|$$

$$0 = 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y}$$

$$\Rightarrow \text{closed form solution is } \hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

### (II) Apply this to our example

From (I) we can conclude, that it should be also possible to find a closed form solution  $\hat{c}_\theta$  should be found, since

$\hat{g}$  is a linear model in terms of the parameters  $c_\theta$

as shown in Question 3.2:  $\hat{g}(x; c_\theta) = c_\theta^\top \Phi(x)$

with  $X = \begin{bmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_N) \end{bmatrix}$  and  $y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$  and  $\hat{y}(x_n; c_\theta) = c_\theta^\top \Phi(x_n)$

and the Loss function

$$L(c_\theta; D) = \frac{1}{2} \sum_{n=1}^N (\hat{y}_n(x_n; c_\theta) - y_n)^2 = \frac{1}{2} \|X \cdot c_\theta - y\|^2$$

the closed form solution should be

proof in appendix :)

$$\hat{c}_\theta = (X^\top X)^{-1} X^\top y$$

according to (I)

we can get this from  $\nabla_{c_\theta} L(c_\theta; D) = 0$  and its okay to change the order to free  $c_\theta$  since the inverse of  $X^\top X$  exist because  $X$  has full rank!

### (III) Global minimization in feed forward NNs

In general: Loss functions of feed forward NNs = non-convex

$\Rightarrow$  local minima exist

$\Rightarrow$  global minima are intractable

our case: Is special, because we apply the quadratic activation  $g(z) = z^2$  and as shown in questions 3.1 and 3.2 with certain assumptions a feature transformation can be used to write the model as a linear model  $\hat{g}(x; c_\theta)$  in terms of the parameters  $c \in \mathbb{R}^{D(D+1)/2}$  (as said in 3.3.2, n)

$\Rightarrow$  squared loss  $L_g(c_\theta; D)$  is convex

$\Rightarrow$  closed form solution exists

Proof for closed form in 3.4

$$\omega = \begin{bmatrix} \text{vectors} \\ \vdots \end{bmatrix}, \mathbf{y} = \begin{bmatrix} \text{vectors} \\ \vdots \end{bmatrix}$$

$$\begin{aligned}
 & \|X\omega - \mathbf{y}\| \\
 &= (\mathbf{y}^T - \mathbf{y})^T (X\omega - \mathbf{y}) \\
 &= (\omega^T X^T - \mathbf{y}^T)(X\omega - \mathbf{y}) \\
 &= (\omega^T X^T X \omega - \omega^T X^T \mathbf{y} - \mathbf{y}^T X \omega + \mathbf{y}^T \mathbf{y}) \\
 &\quad \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \\
 &\quad = \text{scalar} \quad = \text{scalar} \\
 &\quad \Rightarrow \text{am allowed to write } 2\omega^T X \omega \\
 &= \underline{\omega^T X^T X \omega} - 2\omega^T X \mathbf{y} + \mathbf{y}^T \mathbf{y} \\
 &\quad \text{---} \quad \text{---} \quad \text{Definition norm} \\
 &\quad \text{quadratic}
 \end{aligned}$$

$$D_{\omega} \|X\omega - \mathbf{y}\| = 2X^T X \omega - 2\mathbf{y}^T X + 0 = 0 \quad \leftarrow \frac{1}{2}$$

$$\omega = (X^T X)^{-1} X^T \mathbf{y}$$

or

$$\omega = (X^T X)^{-1} X^T \mathbf{y}$$

Slides for 3.4

### Linear Regression

Often, linear dependency of  $\hat{y}$  on  $x$  is a bad/simplistic assumption

More general model:  $\hat{y} = w^T \phi(x)$ , where  $\phi(x)$  is a feature vector

- e.g.  $\phi(x) = [1, x, x^2, \dots, x^D]^T \in \mathbb{R}^{D+1}$  (monomials degree  $\leq D$ )
- the bias term  $b$  is captured by the constant feature  $\phi_0(x) = 1$

Fit/learn  $w$  by solving  $\min_{w \in \mathbb{R}^{D+1}} \sum_{n=1}^N (y_n - w^T \phi(x_n))^2$

- Closed-form solution:

$$\hat{w} = (X^T X)^{-1} X^T \mathbf{y}, \text{ with } X = \begin{bmatrix} \phi(x_1)^T \\ \vdots \\ \phi(x_N)^T \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Still called **linear regression** – linear w.r.t. the model parameters  $w$ .

J. Martino, F. Melo, M. Figueiredo (IST) Lecture 2 DL IST Fall 2022 16 / 50

### One-Slide Proof

Write the objective function in matrix-vector notation:

$$\sum_{n=1}^N (y_n - w^T \phi(x_n))^2 = \|\mathbf{y} - Xw\|^2$$

Equate the gradient to zero and solve the resulting equation:

$$\begin{aligned}
 0 &= \nabla_w \|\mathbf{y} - Xw\|^2 \\
 &= \nabla_w (w^T X^T X w - 2w^T X^T \mathbf{y} + \|\mathbf{y}\|^2) \\
 &= 2X^T X w - 2X^T \mathbf{y}
 \end{aligned}$$

Therefore

$$\hat{w} = (X^T X)^{-1} X^T \mathbf{y}$$

### Squared Loss Function

Linear regression with least squares criterion corresponds to a loss function

$$L(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2, \text{ where } \hat{y} = w^T \phi(x)$$

This is called the **squared loss**.

The model is fit to the training data by **minimizing the loss function**:

$$\hat{w} = \arg \min_w \sum_{n=1}^N L(y_n - \hat{y}_n) = \arg \min_w \frac{1}{2} \sum_{n=1}^N (y_n - w^T \phi(x))^2$$

(the factor 1/2 is irrelevant but convenient)

More later.

J. Martino, F. Melo, M. Figueiredo (IST) Lecture 2 DL IST Fall 2022 19 / 50