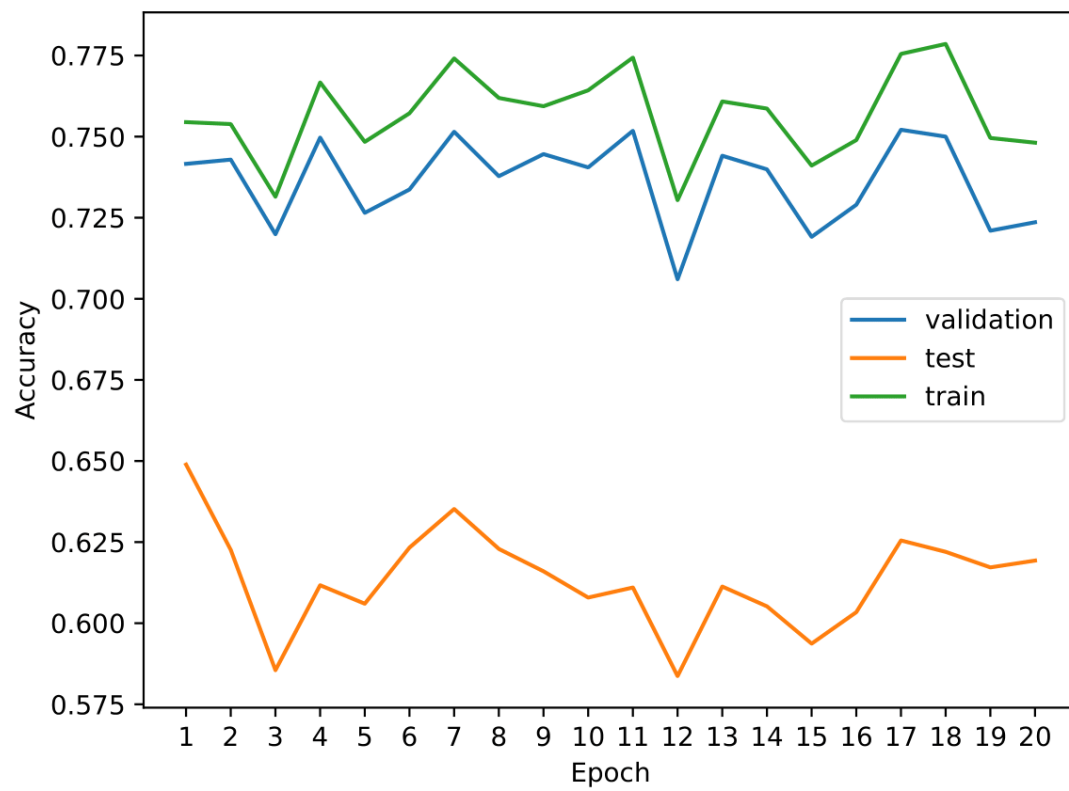
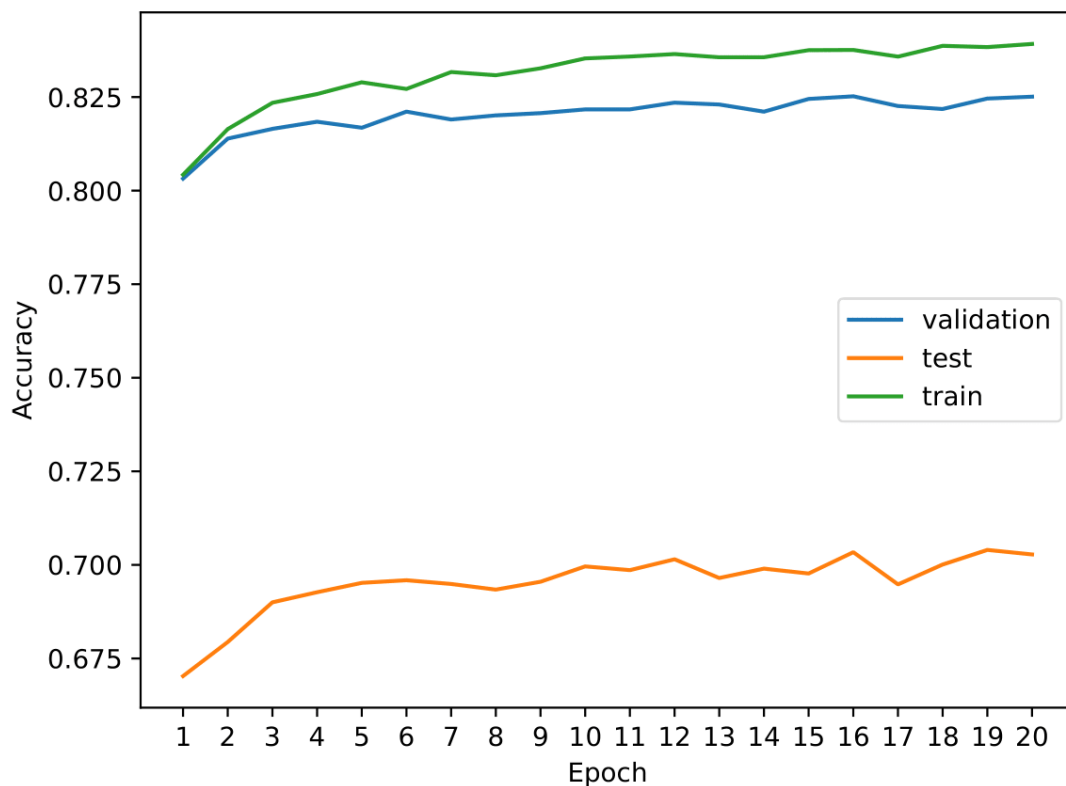


# Question 1

- 1a) Final validation accuracy: 0.7236  
Final test accuracy: 0.6193



- 1b) Final validation accuracy: 0.8251  
Final test accuracy: 0.7028



2a) The use of non-linear activations enables the MLP to capture more complex relationships between the inputs and outputs. If the problem isn't linearly separable, a perceptron is unable to solve it, however an MLP with non-linear activation function is capable of learning a more complex decision boundary, which also depends on the number of hidden units, as stated by the Universal Approximation Theorem.

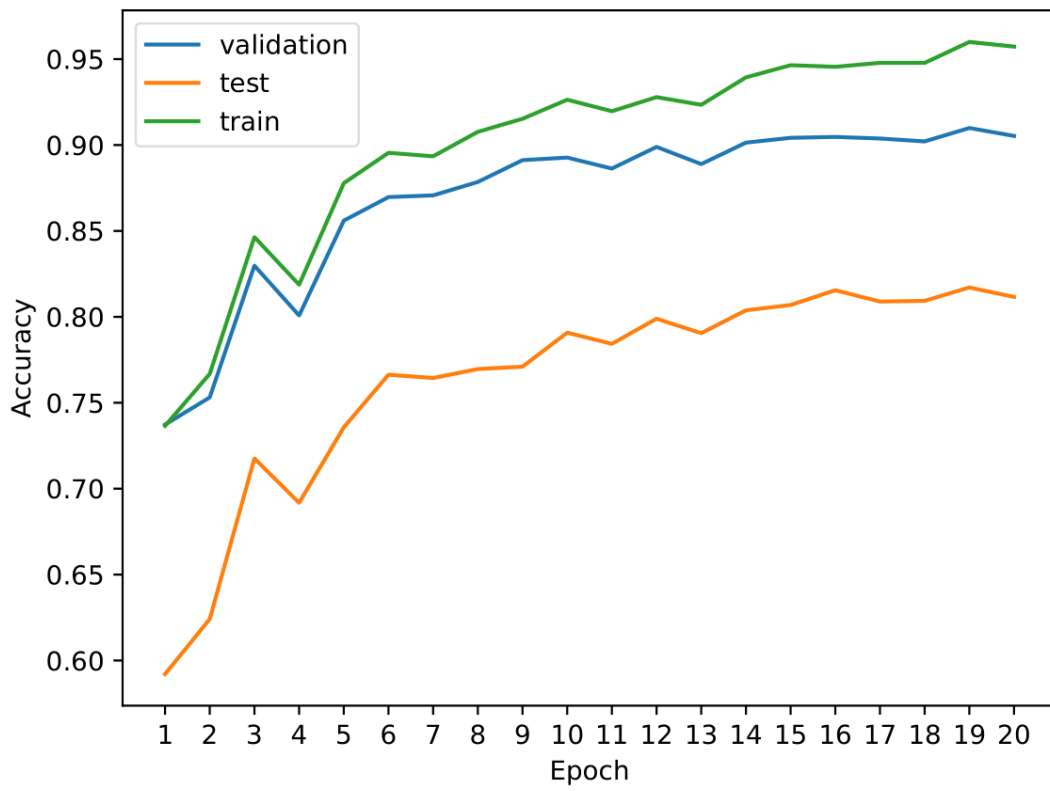
If the activation function is linear, the combination of the various layers can be re-written as being equivalent to a single layer with a linear activation function.

For example, with an activation function  $g(x) = x$ , and with 2 layers, you could rewrite the output  $y = g(z(g(z(x))))$  as

$$g(z(g(W_1x + b_1))) = g(z(W_1x + b_1)) = g(W_2W_1x + W_2b_1 + b_2) = W_2W_1x + W_2b_1 + b_2$$

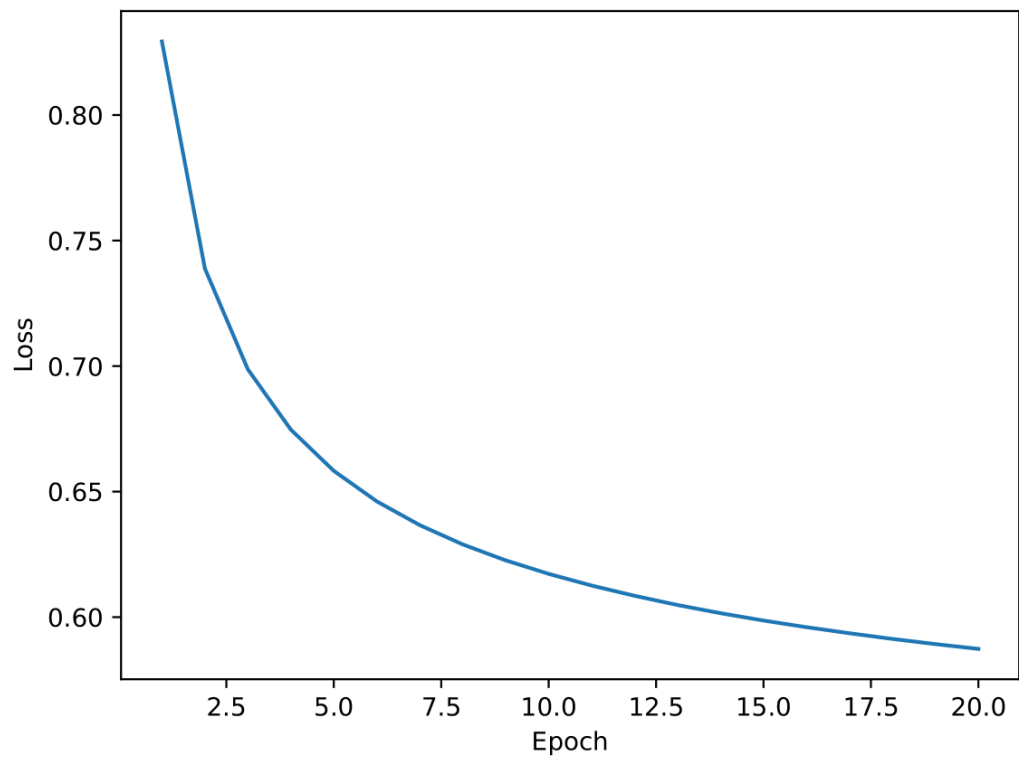
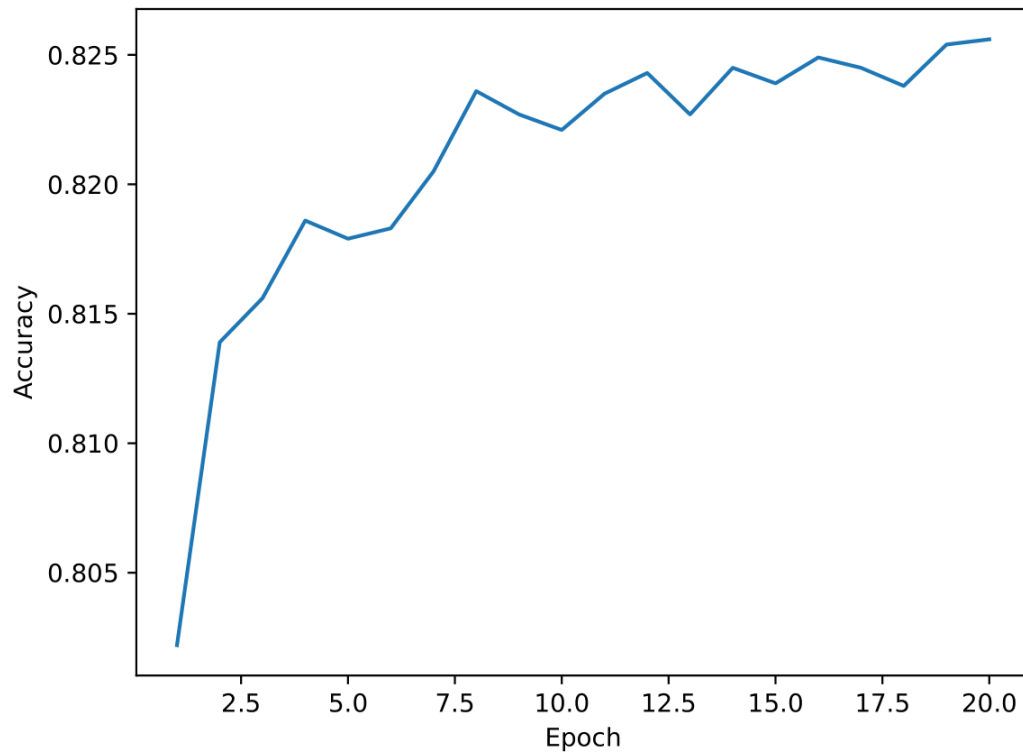
In this particular task the use of MLP instead of the simple perceptron helps us to capture a better representation of the pixels' information.

2b) Final validation accuracy: 0.9053  
Final test accuracy: 0.8116



## Question 2

- 1) Final validation accuracy: 0.8256  
Final test accuracy: 0.7019  
Best configuration: 0.001 learning rate

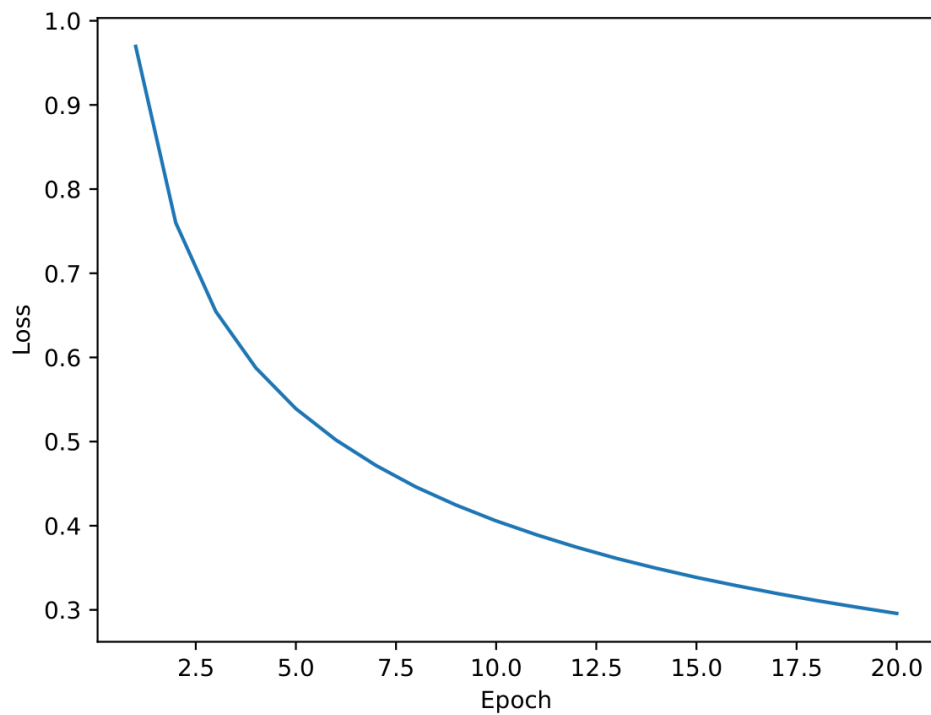
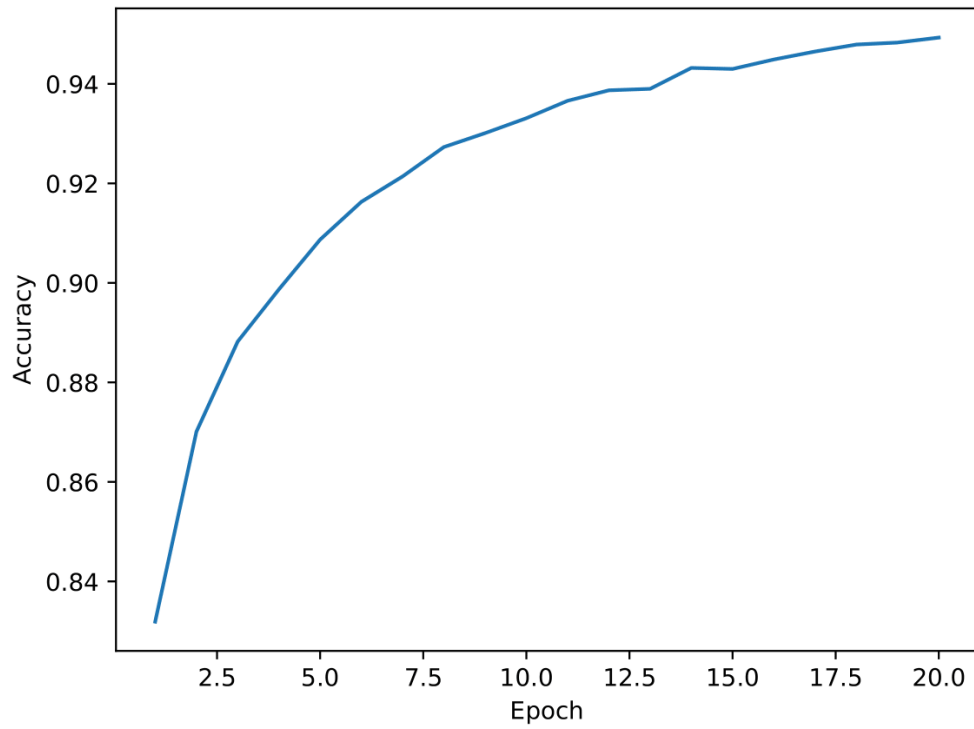


2) Final validation accuracy: 0.9493

Final test accuracy: 0.8814

Best configuration:

- 0.01 learning rate
- 200 hidden units per layer
- 1 layer
- 0.3 Dropout
- ReLU activation
- SGD optimizer

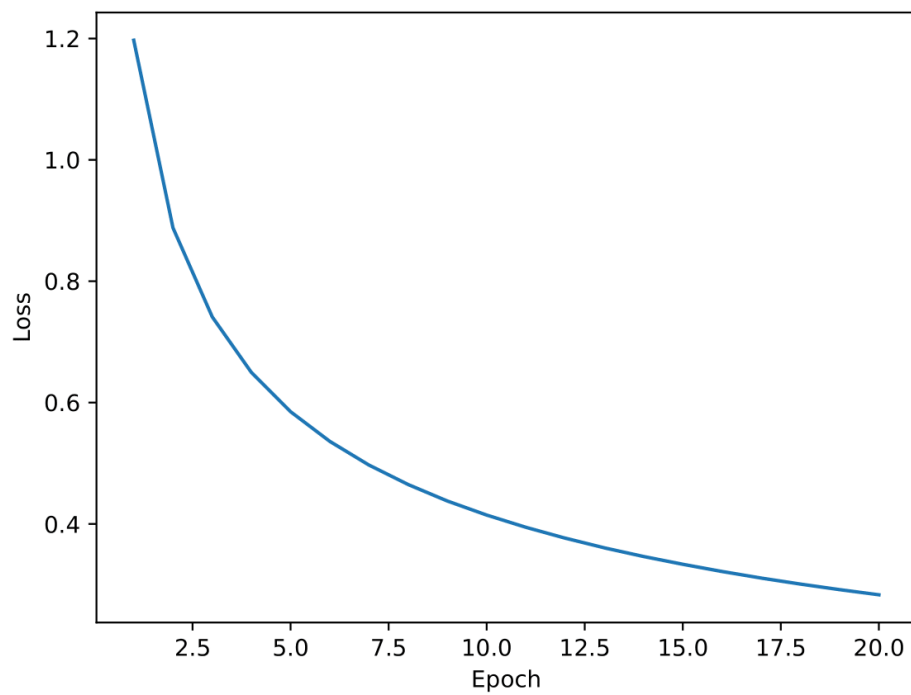
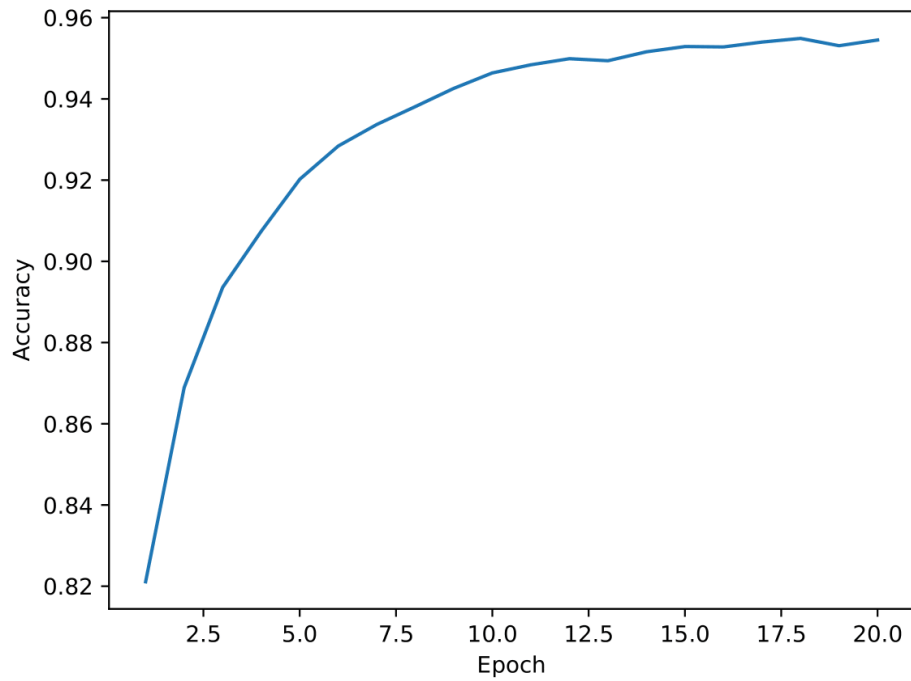


3) Final validation accuracy: 0.9545

Final test accuracy: 0.8949

Best configuration:

- 0.01 learning rate
- 200 hidden units per layer
- 2 layers
- 0.3 Dropout
- ReLU activation
- SGD optimizer



### Question 3

1)

20:04 Fri 23 Dec

< DeepLearning



$$h = \begin{bmatrix} h_1 \\ \vdots \\ h_i \\ \vdots \\ h_K \end{bmatrix} = \begin{bmatrix} g(w_1^T x) \\ \vdots \\ g(w_i^T x) \\ \vdots \\ g(w_K^T x) \end{bmatrix}, \quad h_i = g(w_i^T x)$$

$$\begin{aligned} g(w_i^T x) &= (w_i^T x)^2 = (w_i^T x)^T (w_i^T x) = \text{Tr}((x^T w_i)(w_i^T x)) = \\ &= \text{Tr}(x^T w_i w_i^T x) = \text{Tr}(w_i w_i^T x x^T) = \langle w_i w_i^T, x x^T \rangle_F \\ &\quad (\text{Tr}(ABC) = \text{Tr}(CAB)) \end{aligned}$$

Symmetric matrix:  $A^T = A$

$$(AB)^T = B^T A^T$$

$$A \in \mathbb{R}^{m \times n}$$

$$(AA^T)^T = (A^T)^T A^T = AA^T \Rightarrow \underline{AA^T \text{ is symmetric.}}$$

$w_i w_i^T$  is symmetric  $\mathbb{R}^{D \times D}$ , also  $x x^T$

Both have  $\frac{D(D+1)}{2}$  independent terms.  
(upper triangle of matrix)

$$\begin{aligned} \langle w_i w_i^T, x x^T \rangle_F &= \underbrace{\text{vec}(w_i w_i^T)^T}_{\frac{D(D+1)}{2} \text{ indep. terms}} \underbrace{\text{vec}(x x^T)}_{\frac{D(D+1)}{2} \text{ indep. terms.}} \\ &\downarrow \\ &\text{Linear combination} \\ &\text{of } \frac{D(D+1)}{2} \text{ independent terms.} \end{aligned}$$

terms.  $\frac{D(D+1)}{2}$  independent

$$h = \begin{bmatrix} h_1 \\ \vdots \\ h_i \\ \vdots \\ h_K \end{bmatrix} = \begin{bmatrix} \text{vec}(w_1 w_1^T) \text{vec}(x x^T) \\ \vdots \\ \text{vec}(w_i w_i^T) \text{vec}(x x^T) \\ \vdots \\ \text{vec}(w_K w_K^T) \text{vec}(x x^T) \end{bmatrix} = \begin{bmatrix} a_{1\theta} \\ \vdots \\ a_{i\theta} \\ \vdots \\ a_{K\theta} \end{bmatrix} \phi(x)$$

$(K \times 1)$   $(K \times \frac{D(D+1)}{2})$

$(\frac{D(D+1)}{2} \times 1)$

Each row  $h_i$  in vector  $h$  is a sum of the  $\frac{D(D+1)}{2}$  independent terms resulting from  $\text{vec}(w_i w_i^T) \text{vec}(x x^T)$ .

Therefore we can write  $h$  as  $A_\theta \phi(x)$  where  $A_\theta = \begin{bmatrix} a_{1\theta} \\ \vdots \\ a_{i\theta} \\ \vdots \\ a_{K\theta} \end{bmatrix}$  (symmetric) and each  $a_{i\theta}$  is a vector  $(\frac{D(D+1)}{2} \times 1)$  containing all the independent terms of  $\text{vec}(w_i w_i^T)$  multiplied by their frequency on the  $\text{vec}(w_i w_i^T)$ ,  $a_{i\theta} = [w_1^2 \ w_2^2 \ \dots \ w_3^2 \ \dots \ 2w_{K-1}w_K]$ .

That means:  $a_{i\theta} \in \mathbb{R}^{\frac{D(D+1)}{2}}$ ,  $i \in \{1, \dots, K\}$   $w_{K-1}w_K = w_K w_{K-1}$

$\hookrightarrow \underline{A_\theta \in \mathbb{R}^{K \times \frac{D(D+1)}{2}}}$

By the same logic  $\phi(x) \in \mathbb{R}^{\frac{D(D+1)}{2} \times 1}$ ,  $\phi(x) = \begin{bmatrix} \phi_1(x) \\ \vdots \\ \phi_i(x) \\ \vdots \\ \phi_{\frac{D(D+1)}{2}}(x) \end{bmatrix}$

$\phi_i(x) \in \mathbb{R}$  and corresponds to the  $i^{\text{th}}$  term of the  $\frac{D(D+1)}{2}$  independent terms in  $\text{vec}(x x^T)$ , where each repeated term is multiplied by 2, as in the case of  $a_{i\theta}$ .

$\underline{h = A_\theta \phi(x)}$ ,  $A_\theta \in \mathbb{R}^{K \times \frac{D(D+1)}{2}}$ ,  $\phi(x) \in \mathbb{R}^{\frac{D(D+1)}{2}}$



2)

20:11 Fri 23 Dec

...

81%



$$\hat{y} = v^T h = v^T A_{\theta} \phi(x)$$

$$\underline{v^T A_{\theta}} = \begin{bmatrix} v_1 & \dots & v_i & \dots & v_K \end{bmatrix}_{(1 \times K)} \begin{bmatrix} a_{1\theta} \\ \vdots \\ a_{i\theta} \\ \vdots \\ a_{K\theta} \end{bmatrix}_{(K \times \frac{D(D+1)}{2})} =, \quad a_{i\theta} = \begin{bmatrix} w_{i1}^2 \\ \vdots \\ w_{iD}^2 \\ \vdots \\ 2w_{iK}w_{K-1} \end{bmatrix}$$

$$= \begin{bmatrix} \sum_{i=1}^K v_i w_{i1}^2 \\ \vdots \\ \sum_{i=1}^K v_i w_{iD}^2 \\ \vdots \\ 2 \sum_{i=1}^K v_i w_{iK} w_{K-1} \end{bmatrix}^T = \underline{C_{\theta}^T} \Rightarrow \underline{\hat{y}} = \underline{C_{\theta}^T} \phi(x), \quad C_{\theta} \in \mathbb{R}^{\frac{D(D+1)}{2}}_{\theta \in (w, v)}$$

$C_{\theta}$  is not linear in terms of  $\theta$  because its rows do not consist of linear combinations of the parameters of  $\theta \in (w, v)$ .

3)

Our neural network structure consists of a single hidden layer which receives  $x \in \mathbb{R}^D$  as input and has a quadratic activation function  $g(z) = z^2$ ; and an output layer which receives the internal representation of the network ( $h$ ).

Given model parameters  $\Theta = (W, v) \in \mathbb{R}^{K \times D} \times \mathbb{R}^K$ , the matrix  $W \in \mathbb{R}^{K \times D}$  represents the weights between the input vectors (of Size  $D$ ) and the hidden layer (of size  $K$ ), and the vector  $v \in \mathbb{R}^K$  represents the weights between the hidden layer and the output layer.

Our model receives the internal representation of the network as input and outputs a prediction, given weights  $v$ :  $\hat{y} \in \mathbb{R}$  is given by  $\hat{y} = v^T h$ , where  $h \in \mathbb{R}^K$ ,  $h = g(Wx)$ . Since  $g(z) = z^2$ , this means that  $h$  is a vector of quadratic functions of the input variables  $x \in \mathbb{R}^D$ .

The feature transformation  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}^{D(D+1)/2}$  takes  $x \in \mathbb{R}^D$  as input and maps it to a higher-dimensional space, represented by the vector  $\phi(x) \in \mathbb{R}^{D(D+1)/2}$ .

We also proved we can write the predicted output  $\hat{y}$  as  $\hat{y}(x; c\Theta) = c^T \Theta \phi(x)$ , where  $c\Theta \in \mathbb{R}^{D(D+1)/2}$ .

Considering  $K \geq D$ , the number of hidden units  $K$  is greater than or equal to the number of input variables  $D$ , so there must exist a  $c \in \mathbb{R}^{D(D+1)/2}$  such that  $c = c\Theta$  because there are enough hidden units to represent the features of each input variable and their relations in the higher-dimensional space represented by the feature transformation  $\phi(x)$ .

This does not mean that it is a linear model in terms of  $c\Theta$  given.

On the other side if  $K < D$ , the number of hidden units  $K$  is smaller than the number of input variables, and as such there might not be a choice of the original parameters such that  $c = c\Theta$  for  $c \in \mathbb{R}^{D(D+1)/2}$ .

4)

20:32 Fri 23 Dec

&lt; DeepLearning



80%

$$L(\theta; \mathcal{D}) = \frac{1}{2} \sum_{n=1}^N (\hat{y}_n(x_n, \theta) - y_n)^2 \quad \mathcal{D} = \{(x_n, y_n)_{n=1}^N\}$$

$$\nabla_{\theta} L(\theta; \mathcal{D}) = 0 ? \quad (\text{To find closed form solution})$$

$$\Phi = \begin{bmatrix} \phi_1(x) \\ \vdots \\ \phi_N(x) \end{bmatrix}, \quad \phi_i \in \mathbb{R}^{\frac{\mathcal{D}(\mathcal{D}+1)}{2}}$$

$$\hat{y} = C_{\theta}^T \phi(x)$$

$$\nabla_{\theta} L(\theta; \mathcal{D}) = \nabla_{\theta} \sum_{n=1}^N \frac{1}{2} (C_{\theta}^T \phi(x_n) - y_n)^2 = \nabla_{\theta} \frac{1}{2} \|y - \Phi C_{\theta}^T\|^2 =$$

$$= \frac{1}{2} \nabla_{\theta} (C_{\theta} \Phi^T \Phi C_{\theta}^T - 2 C_{\theta} \Phi^T y + \|y\|^2) = \frac{1}{2} (2 \Phi^T \Phi C_{\theta}^T - 2 \Phi^T y) \Rightarrow$$

$$\Rightarrow \nabla_{\theta} L(\theta; \mathcal{D}) = \Phi^T \Phi C_{\theta}^T - \Phi^T y$$

$$\nabla_{\theta} L(\theta; \mathcal{D}) = 0 \Leftrightarrow \hat{C}_{\theta}^T = (\Phi^T \Phi)^{-1} \Phi^T y \quad (\text{CLOSED FORM SOLUTION})$$

$$\downarrow$$

$$\Phi^T \Phi \text{ is invertible if and only if } N \geq \frac{\mathcal{D}(\mathcal{D}+1)}{2} \quad (\text{Given full column-rank})$$

Considering an  $N > \frac{\mathcal{D}(\mathcal{D}+1)}{2}$ , then we can find the closed form solution for  $\hat{C}_{\theta}$ .

In addition to that, all features must NOT be linear combinations of others, and that is guaranteed, since  $X$  has full column-rank.