

## Neural networks (NN)

Consider  $f(\vec{x}; \vec{\theta}) = \vec{w} \cdot \vec{s}(\vec{x}) + b$ , where  
 $\vec{\theta} = (\vec{w}, b)$ .

A natural idea is to construct features with their own parameters:

$f(\vec{x}; \vec{\theta}) = \vec{w} \cdot \vec{s}(\vec{x}; \vec{\theta}_2) + b$ , where  
 $\vec{\theta}_1 = (\vec{w}, b)$  &  $\vec{\theta} = (\vec{\theta}_1, \vec{\theta}_2)$ .

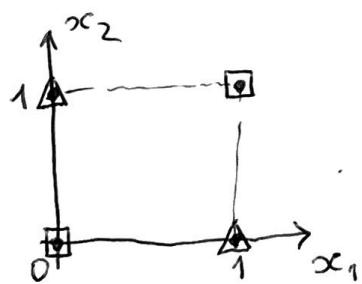
In deep neural networks (DNN), the functions are composed recursively:

$$\begin{array}{ll} f_1(\vec{x}; \vec{\theta}_1) & 1 \\ f_2(f_1(\vec{x}; \vec{\theta}_1); \vec{\theta}_2) & 2 \\ \vdots & \vdots \\ f_e(f_{e-1}(\dots); \vec{\theta}_e) & e \end{array}$$

### The XOR problem

$x_1$	$x_2$	$y$	
0	0	0	□
0	1	1	△
1	0	1	△
1	1	0	□

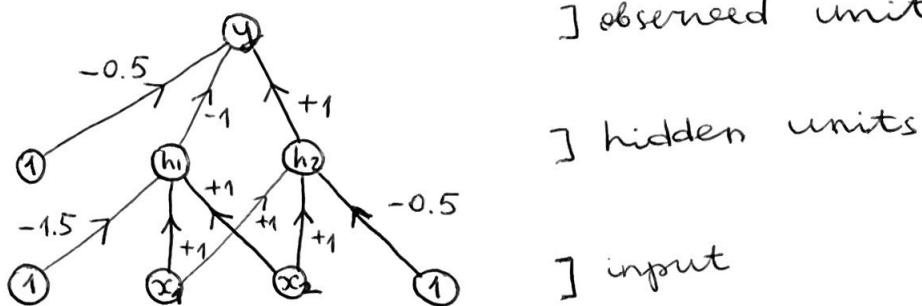
$\Leftarrow$  truth table



The XOR function is not linearly separable: no straight-line DB can separate  $\Delta$  from  $\square$ .

However, we can use linear-DB functions to implement it anyway:

$$f(\vec{x}; \vec{\theta}) = \mathbb{I}(\vec{w} \cdot \vec{x} + b \geq 0) \quad \begin{matrix} \leftarrow \text{perceptron} \\ ] \text{observed unit} \end{matrix}$$



$$h_1 = \mathbb{I}(x_1 + x_2 - 1.5 \geq 0) \quad \text{AND}$$

~~If~~  $x_1 = 1, x_2 = 1 \Rightarrow h_1 = 1,$   
otherwise  $h_1 = 0.$

$$h_2 = \mathbb{I}(x_1 + x_2 - 0.5 \geq 0) \quad \text{OR}$$

~~If~~  $x_1 = 0, x_2 = 0 \Rightarrow h_2 = 0,$   
otherwise  $h_2 = 1.$

$$\text{Finally, } y = \mathbb{I}(h_2 - h_1 - 0.5 \geq 0)$$

$x_1$	$x_2$	$h_1$	$h_2$	$y$
0	0	0	0	0
0	1	0	1	1
1	0	0	1	1
1	1	1	1	0

same  
as before  
(XOR)

# Differentiable multi-layer perceptrons (MLP)

Define a differentiable activation function:

$$\varphi: \mathbb{R} \rightarrow \mathbb{R}$$

Make linear transformations of the inputs,  
or the outputs of the previous layer:

$$\tilde{a}_e = \tilde{b}_e + W_e \underbrace{\tilde{z}_{e-1}}_{\text{or } \tilde{x}} \quad \begin{matrix} \dim = \# \text{ nodes} \\ \text{of layer } e \end{matrix}$$

↑  
layer index

Then apply the activation function  
to  $\tilde{a}_e$  element-wise:

$$\varphi_e(\tilde{a}_e) = \varphi_e(\tilde{b}_e + W_e \tilde{z}_{e-1}) \equiv \tilde{z}_e \leftarrow \begin{matrix} \text{same} \\ \uparrow \\ \dim \text{ as } \tilde{a}_e \end{matrix}$$

output of  
layer  $e$

In the final layer, the # nodes is  
dictated by the problem.

[activation functions must be non-linear:]

if  $\varphi_e(a) = C_e a$ ,

$$f(\tilde{x}; \tilde{\theta}) = W_L C_L (W_{L-1} C_{L-1} (\dots (W_1 \tilde{x}) \dots)) \sim$$

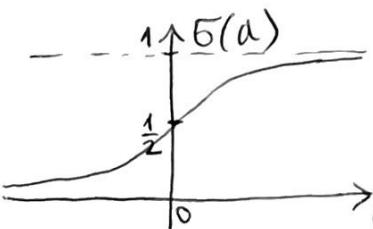
↑  
total # layers

$$\sim W_L W_{L-1} \dots W_1 \tilde{x} = \underbrace{W^T \tilde{x}}_{\text{linear model}}$$

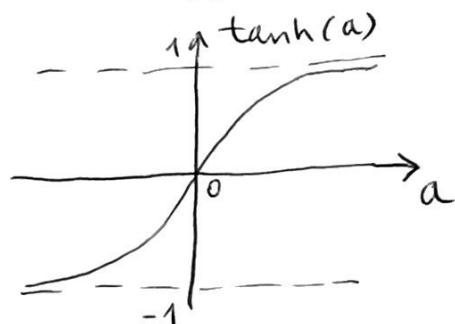
## Activation functions

1. Sigmoid:

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$



2.  $\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$

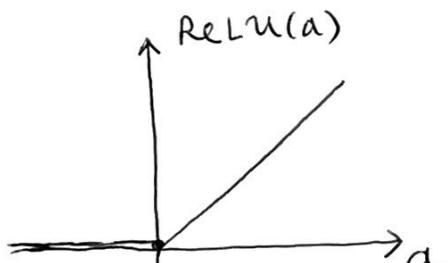


Both  $\sigma(a)$  &  $\tanh(a)$  saturate at  $a \rightarrow +\infty$  and  $a \rightarrow -\infty$   $\Rightarrow$  vanishing gradient problem

3. Rectified linear unit (ReLU):

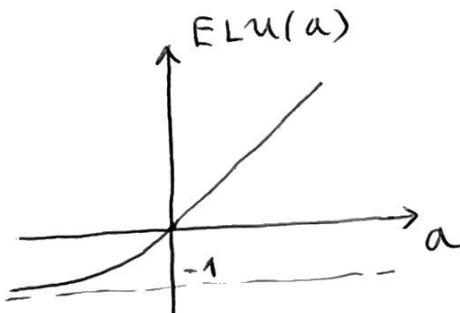
$$\text{ReLU}(a) = \max(a, 0)$$

turns off negative inputs, but passes positive inputs unchanged

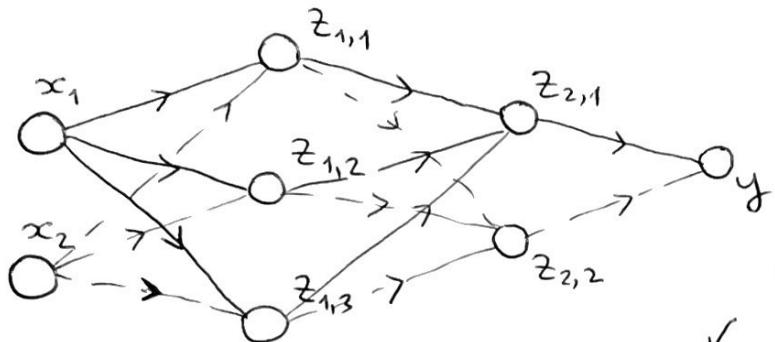


4. Exponential linear unit (ELU):

$$\text{ELU}(a) = \begin{cases} a, & a \geq 0 \\ e^{a-1}, & a < 0 \end{cases}$$



Ex.  $\vec{x} = (x_1, x_2)$  2D input  
output: prob. of class 1 (out of 2 classes)  
possible



[biases omitted for clarity]

interpreted as  
 $p_1 = 1 - p_2$  - prob. of class 1

$y \in [0, 1]$   
with sigmoid activation function

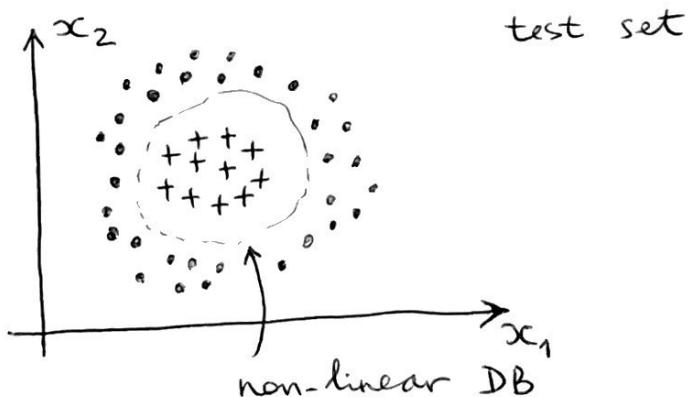
$$\begin{array}{c} \vec{x} \in \mathbb{R}^2 \\ l \\ 1 \end{array} \quad \begin{array}{c} \vec{z}_1 \in \mathbb{R}^3 \\ 2 \end{array} \quad \begin{array}{c} \vec{z}_2 \in \mathbb{R}^2 \\ 3 \end{array} \quad \begin{array}{c} y \in [0, 1] \\ \text{with sigmoid activation} \\ \text{function} \end{array}$$

Mathematically,

$$\begin{cases} \vec{z}_1 = \varphi(\vec{w}_1 \cdot \vec{x} + \vec{b}_1), \\ \vec{z}_2 = \varphi(\vec{w}_2 \cdot \vec{z}_1 + \vec{b}_2), \\ a_3 = \vec{w}_3 \cdot \vec{z}_2 + b_3 \Rightarrow p(y|\vec{x}; \vec{\theta}) = \text{Ber}(y|G(a_3)). \end{cases}$$

$\vec{\theta} = (\vec{w}_1, \vec{b}_1, \vec{w}_2, \vec{b}_2, \vec{w}_3, b_3)$  need to be fitted to maximize  $\log L$ .

Result:



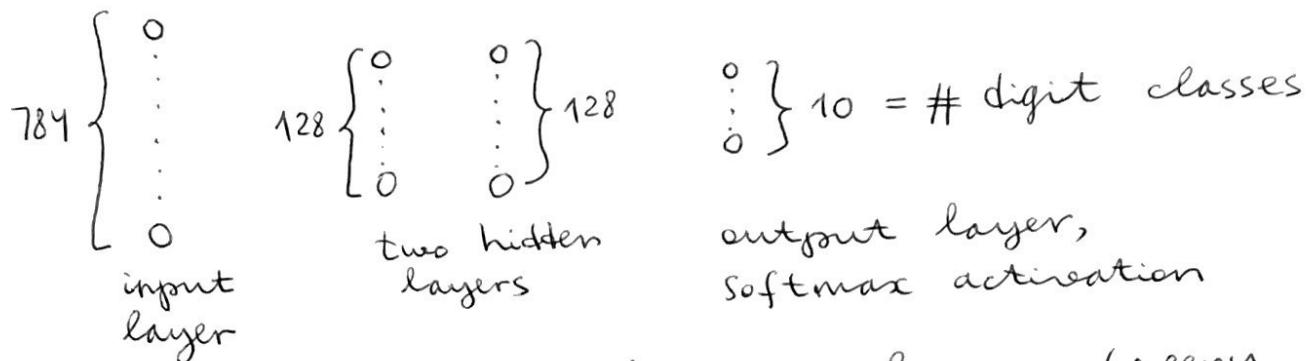
## MLP for image classification

Consider MNIST database of digits:



$\Rightarrow$  flatten into a vector with  $D = 28^2 = 784$ .

# prms    100480    16512    1290



Dense connections between layers (every node in layer  $l$  connected to every node in layer  $l+1$ ).

This model can reach test accuracy of ~97% after training for 2 epochs.

## MLP for text classification

Idea: convert a variable-length sequence of words  $\tilde{v}_1, \dots, \tilde{v}_T$  into a fixed-size vector.  $T$   $\leftarrow$  length of the document

$\tilde{v}_j$  = 'one-hot' vector of length  $V$ ,  
 $V$  = vocabulary size.

The simplest way to interpret a text is to treat it as a 'bag of words', i.e. ignore word order.

- ① Maps each word to a token via preprocessing steps :

- drop punctuation
- convert all words to lowercase
- ~~drop~~ drops "a", "the", "and", etc. (stop word removal)
- replace conjugations with base forms,  
e.g. "running"  $\rightarrow$  "run" (word stemming)

Tokens come from a vocabulary.

Now, consider  $x_t = \text{token at location } t \text{ in the document}$

Compute  $\tilde{x}_v = \sum_{t=1}^T \mathbb{I}(x_t = v)$   
↑ taken  $v$

Now we have  $\tilde{x} = D\text{-dim vector of token counts}$

Note that  $D \leq V$  vocab. size

$\tilde{x}$  is the vector space model of the text.

With  $N$  texts, we can construct

a  $D \times N$  (or  $V \times N$ ) term frequency matrix  
 $D_{\max} \times N$ , where  $D_{\max}$  is the union of all  $D$ 's  
 TF, where  $TF_{ij} = \text{freq. of term } i \text{ in document } j$

$$\sum_{i=1}^D TF_{ij} = 1, \forall j \quad \text{each column is normalized}$$

what to do with common (i.e., high-freq.) words? Define  $\log(TF_{ij} + 1)$  - log frequencies.

Also, define inverse document frequency:

$IDF_i = \log \frac{N}{1 + DF_i}$ , where  $1 \leftarrow$  in the  $D_{max}$  scheme  
 $0 \leq DF_i \leq N$   
 $N$  is the number of documents with term  $i$ .

$$\begin{cases} DF_i = 1 \Rightarrow IDF_i = \log\left(\frac{N}{2}\right) \\ DF_i = N \Rightarrow IDF_i = \log\left(\frac{N}{1+N}\right) \rightarrow 0 \text{ as } N \rightarrow \infty \end{cases}$$

Thus,  $IDF_i$  can be used to penalize common words:

compute  $\underbrace{TFIDF_{ij}}_{\text{TF-IDF matrix}} = \log(TF_{ij} + 1) \times IDF_i$

The TF-IDF transform does not place semantically similar words together in vector space.

Idea: use word embeddings

$\vec{v}_t$  = one-hot vector of <sup>the</sup> token at position  $t$  in the document  
 $\dim V$

is mapped into

$\vec{e}_t \in \mathbb{R}^K$  = linear-dimensional dense vector ( $K \ll V$ )  
 $\dim K$

$\vec{e}_t = E_{K \times V} \vec{v}_t$   
 $E_{K \times V}$  embedding matrix, places semantically similar words together

With  $N$  documents, we have

can compute  $\vec{e}_n = \sum_{t=1}^T \vec{e}_{nt}$

$\vec{e}_{nt} \in \mathbb{R}^K$ ,  
↑ token index  
document index

bag of word embeddings, can be used as input to a NN classifier document

Word embeddings should be similar for semantically similar words.

Assumption: two words are similar if they occur in similar contexts (the distributional hypothesis).

Thus, it makes sense to learn the word's embedding vector from its context.

## → Latent semantic indexing (LSI)

Compute  $C_{ij} = \# \text{ times token } i \text{ occurs in document } j \quad (j=1, \dots, N)$

$C_{M \times N}$  = count matrix  
↑  
vocab. size

approximate  $C$  with rank  $K$  (truncated)

SVD approximation:

$$C_{ij} \approx \sum_{k=1}^K u_{ik} s_k v_{jk}$$

Then  $\tilde{u}_i$  = embedding for word  $i = 1, \dots, M$   
length  $K$

$\tilde{w}_j = \sum_{t=1}^{T_j} \tilde{u}_t$  = embedding for document  $j = 1, \dots, N$   
 $\uparrow$   
element-wise

One can compute a 'bag of words' vector:  
for each document

$$\tilde{q}_{fj} = \frac{1}{T_j} \sum_{t=1}^{T_j} \tilde{u}_t$$

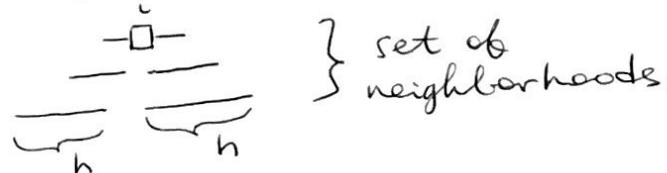
$\uparrow j$   
# terms in document  $j$

and use cosine similarity to compare

$$\tilde{q}_{fj} \text{ with } \tilde{w}_j : \frac{\tilde{q}_{fj} \cdot \tilde{w}_j}{\|\tilde{q}_{fj}\| \|\tilde{w}_j\|}$$

## → Latent semantic analysis (LSA)

Do the same SVD analysis as in LSI,  
but for each word  $i$ , define local neighborhoods:



Thus,  $c_{ij} = \# \text{ times token } i \text{ occurs in neighborhood } j$

$h$  is a hyperparameter.

## → positive pointwise MI (PPMI)

use  $\text{PPMI}(i,j) = \max(\text{PMI}(i,j), 0)$ ,  
 $\wedge$  remove negative correlations

where

$$\text{PMI}(i,j) = \log \frac{p(i,j)}{p(i)p(j)}$$

Do k-rank SVD decomposition on  
 $\text{PPMI}(i,j)$  instead of  $C_{ij}$ .

→ word2vec

~~Question~~

Finally, use word embedding as part  
of the NN architecture itself:

input = unordered bag of words  
 $\vec{v}_1, \dots, \vec{v}_T$  'one-hot'  
 $E \times V$  document length

(1) Embedding:  $\tilde{\ell}^t = \underbrace{W_1}_{\text{'do all words simultaneously'}} \vec{v}_t \quad t=1, \dots, T$

$$\text{pooling: } \langle \tilde{\ell} \rangle = \frac{1}{T} \sum_{t=1}^T \tilde{\ell}^t$$

(2) global average layer:  $\tilde{h} = \underbrace{W_2}_{H \times E} \langle \tilde{\ell} \rangle + \tilde{b}_2$

(3) Hidden layer:  $p(y| \vec{v}_1, \dots, \vec{v}_T) =$

(4) Binary classification:  $= \text{Ber}(y| G(\tilde{w}_3 \cdot \tilde{h} + b_3))$

This can be used to e.g. classify movie reviews into pos. & neg.

Ex.  $V=10^4$  words,  $E=16$ ,  $H=16$   
embedding size hidden layer size

Layer	# prms	
1	$16 \times 10^4$	
2	0	
3	$16 \times 16 + 16$ <sup>bias</sup> = 272	
4	$16 + 1$ <sup>bias</sup> = 17	
		160,289 in total

However,  $W_i$  can be obtained ('pre-trained') separately using the techniques described above, to reduce the # prms & prevent overfitting.