

Deep Learning 1

example tasks

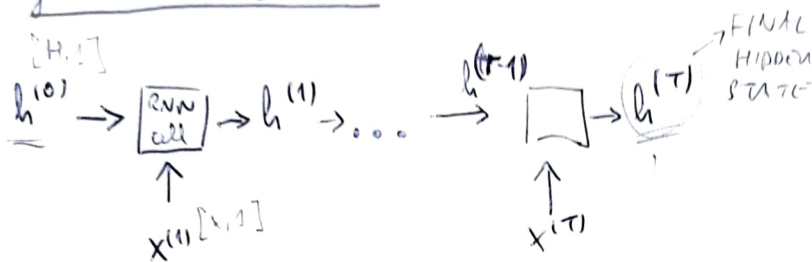
1. We consider a simple recurrent neural network that classifies sequences of two-bit binary numbers into two classes, depending on ratio of odd and even numbers. A sequence belongs to class 0 if the number of even numbers is larger or equal to the number of odd numbers, and to class 1 if there are more odd numbers. For example, a sequence 00,01,10,11,00 contains three even and 2 odd numbers which means that it belongs to class 0.
 - (a) Determine the exact parameters and appropriate activation functions for a basic recurrent cell and a fully connected output layer that would properly solve this classification problem. Model the input $x^{(t)}$ as a binary column vector with two elements. Express the final model with appropriate equations.
 - (b) How many parameters does your model have?
 - (c) Demonstrate the correctness of your model on the following test sequence: 00,10,11.

① $\begin{matrix} 0 & 1 & 2 & 3 & 0 \\ 00, 01, 10, 11, 00 & \rightarrow \# \text{even} > \# \text{odd} & \Rightarrow d=0 \\ 00, 01, 10, 11 & \rightarrow \# \text{even} = \# \text{odd} & \Rightarrow d=0 \\ 00, 01, 10, 11, 01 & \rightarrow \# \text{even} < \# \text{odd} & \Rightarrow d=1 \end{matrix}$

RNN
+
FC on $h^{(T)}$

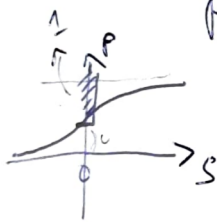
binary cl. \rightarrow sigmoid act. in output layer

general architecture



$$s = W h^{(T)} + b$$

$$p = \sigma(s)$$



$$d: \begin{cases} 0, p < 0.5 \\ 1, p \geq 0.5 \end{cases} \Leftrightarrow \begin{cases} \# \text{even} \geq \# \text{odd} \\ \# \text{even} < \# \text{odd} \end{cases}$$

$$s = (1) \cdot h^{(T)} - 0.5$$

$$s = (1) \cdot h^{(T)} - 0.5$$

$$h^{(T)} = (\# \text{odd} - \# \text{even})$$

$$h^{(t)} = (\# \text{odd} - \# \text{even}) \odot t$$

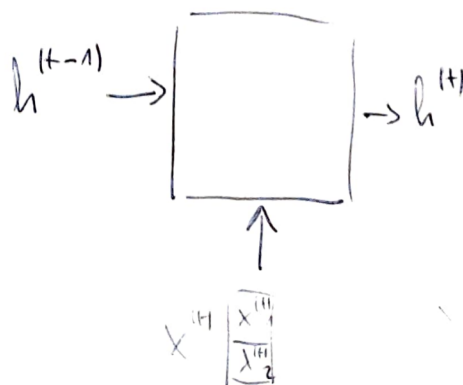
$$h^{(0)} = (\dots) \odot 0 = 0$$

$$h^{(t)} = \begin{cases} h^{(t-1)} + 1 & \text{if } x^{(t)} \text{ odd} \\ h^{(t-1)} - 1 & \text{if } x^{(t)} \text{ even} \end{cases}$$

$$X^{(t)} = \begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \end{bmatrix}$$

$$00 \rightarrow \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$01 \rightarrow \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



$$h^{(t)} = \underset{\substack{[H,H] \\ [1,1]}}{w_{hh}} h^{(t-1)} + \underset{\substack{[H,1] \\ [1,1]}}{w_{hx}} x^{(t)} + \underset{\substack{[H,1] \\ [1,1]}}{b_h}$$

$$h^{(t)} = \boxed{w_{hh}} \cdot h^{(t-1)} + \underset{w_{hx}}{[\quad]} \begin{bmatrix} x^{(t)}_1 \\ x^{(t)}_2 \end{bmatrix} + \boxed{b_h}$$

$$\begin{array}{c|c} 0 & 0 \text{ ev} \\ 0 & 1 \text{ odd} \\ 1 & 0 \text{ ev} \\ 1 & 1 \text{ odd} \end{array}$$

$$\begin{bmatrix} 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} -1 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 2 \end{bmatrix}$$

$$h^{(t)} = (1) \cdot h^{(t-1)} + \begin{bmatrix} 0 & 2 \end{bmatrix} \begin{bmatrix} x^{(t)}_1 \\ x^{(t)}_2 \end{bmatrix} + (-1)$$

b) RNN: $w_{hh} \rightarrow 1$
 $w_{hx} \rightarrow 2$
 $b_h \rightarrow 1$ } 4

FC: $W \rightarrow 1$
 $b \rightarrow 1$ } 2

$\rightarrow \boxed{\underline{6}}$

c) $h^{(0)} = 0$ $h^{(1)} = (1) \cdot 0 + \begin{bmatrix} 0 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + (-1) = (-1)$

0 $x^{(1)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ $h^{(2)} = (1)(-1) + \begin{bmatrix} 0 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + (-1) = (-2)$

2 $x^{(2)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ $h^{(3)} = h^{(3)} = (1)(-2) + \begin{bmatrix} 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + (-1) = (-1)$

3 $x^{(3)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $S = (1) \cdot (-1) - 0.5 = -1.5$
 $P = \sigma(1.5) = \sigma(-1.5) < 0.5 \rightarrow d: 0$

2. We consider siamese metric learning with a quadratic contrastive loss. We perform the embedding with a single convolutional layer without padding and with stride 1, The initial parameters of the convolutional layer are:

$$\mathbf{w} = \begin{bmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{bmatrix}, \quad b = -0.1$$

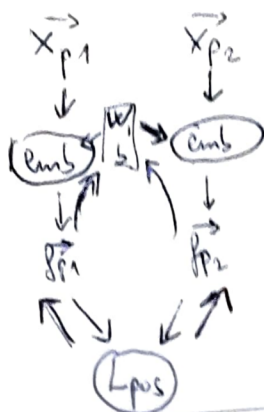
We consider the following positive pair as an input:

$$\mathbf{x}_{p1} = \begin{bmatrix} 1.0 & 1.0 & 1.3 & 1.0 & 1.0 & 1.3 \end{bmatrix}, \quad \mathbf{x}_{p2} = \begin{bmatrix} 3.3 & 3.3 & 2.7 & 3.3 & 3.3 & 2.7 \end{bmatrix}$$

Questions:

- (a) Write the equations that define the model and loss for positive and negative pairs.
- (b) Execute a forward pass for the given input and calculate the loss.
- (c) Determine the gradients with respect to model parameters.

2.



$$\text{emb: } w = \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$$

$$b = -0.1$$

$$x_{p1} = [1, 1, 1.3, 1, 1, 1.3]$$

$$x_{p2} = [3.3, 3.3, 2.7, 3.3, 3.3, 2.7]$$

$$a) L_{pos}(x_{p1}, x_{p2}) = \|f_{p1} - f_{p2}\|^2 = \sum_{i=1}^d (f_{p1,i} - f_{p2,i})^2$$

$$L_{neg}(x_{p1}, x_{p2}) = (\max(0, -\|f_{p1} - f_{p2}\| + m))^2$$

$$f_{p1} = x_{p1} \otimes w^{(1)} + b^{(1)} \quad // \text{conv1d}(x_{p1}, w, b)$$

$$f_{p2} = x_{p2} \otimes w^{(2)} + b^{(2)}$$

$$L_{pos}(f_{p1}, f_{p2})$$

$$b) f_{p1} = [1 \ 1 \ 1.3 \ 1 \ 1 \ 1.3] \otimes \left(\frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}\right) - 0.1$$

$$= \frac{1}{3} [3.3 \ 3.3 \ 3.3 \ 3.3] - 0.1$$

$$= [1 \ 1 \ 1 \ 1]$$

$$f_{p2} = [3.3 \ 3.3 \ 2.7 \ 3.3 \ 3.3 \ 2.7] \otimes \left(\frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}\right) - 0.1$$

$$= \frac{1}{3} [9.3 \ 9.3 \ 9.3 \ 9.3] - 0.1$$

$$= [3 \ 3 \ 3 \ 3]$$

$$L_{pos}(f_{p1}, f_{p2}) = \sum_{i=1}^4 (1-3)^2 = 16$$

$$\nabla_{f_{p1}} L = 2 \begin{bmatrix} -2 & -2 & -2 & -2 \end{bmatrix}$$

$$= (-4) \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\nabla_{f_{p2}} L = (4) \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}$$

$$c) \nabla_{f_{p1}} L = ? , \nabla_{f_{p2}} L = ?$$

$$(\nabla_{f_{p1}} L)_k = \frac{\partial L}{\partial f_{p1,k}} = \frac{\partial}{\partial f_{p1,k}} \left(\sum_{i=1}^4 (f_{p1,i} - f_{p2,i})^2 \right) = 2(f_{p1,k} - f_{p2,k}) \cdot (1)$$

$$\nabla_{f_{p1}} L = 2(\vec{f}_{p1} - \vec{f}_{p2})$$

$$(\nabla_{f_{p2}} L)_k = \frac{\partial L}{\partial f_{p2,k}} = \frac{\partial}{\partial f_{p2,k}} \left(\sum_{i=1}^4 (f_{p1,i} - f_{p2,i})^2 \right) = 2(f_{p1,k} - f_{p2,k}) \cdot (-1)$$

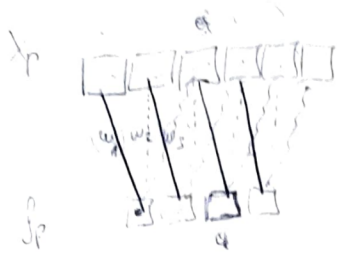
$$\nabla_{f_{p2}} L = 2(\vec{f}_{p2} - \vec{f}_{p1})$$

C) [cont'd]

$$\nabla_{w^{(1)}} L = ?$$

$$\nabla_{w^{(2)}} L = ?$$

$$\frac{\partial L}{\partial w_a} = \sum_{a=1}^4 \frac{\partial L}{\partial p_a} \cdot \frac{\partial p_a}{\partial w_a} = \sum_{a=1}^4 \left(\frac{\partial L}{\partial p_a} \right) \cdot X_{p,a+2}$$



$$\frac{\partial L}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \sum_{a=1}^4 X_{p,a+2} \cdot \frac{\partial L}{\partial p_a}$$

$$- (X_{jL})^T X_{p,a+2}$$

$$\frac{\partial L}{\partial w} = X_p \odot \frac{\partial L}{\partial p}$$

$$\nabla_w L = X_p \odot \nabla_p L$$

$$\nabla_{w^{(1)}} L = [1 \ 1 \ 1.3 \ 1 \ 1 \ 1.3] \odot (-4 [1 \ 1 \ 1 \ 1])$$

$$= (-4) [4.3 \ 4.3 \ 4.6]$$

$$\nabla_{w^{(2)}} L = [3.3 \ 3.3 \ 2.7 \ 3.3 \ 3.3 \ 2.7] \odot (4 [1 \ 1 \ 1 \ 1])$$

$$= (4) [12.6 \ 12.6 \ 12]$$

$$\nabla_w L = \nabla_{w^{(1)}} L + \nabla_{w^{(2)}} L =$$

$$= (4) [8.3 \ 8.3 \ 7.4]$$

$$= [33.2 \ 33.2 \ 29.6]$$

$$\nabla_{b^{(1)}} L = ?, \nabla_{b^{(2)}} L = ?$$

$$\nabla_{b^{(1)}} L = \sum_{i=1}^4 \left(\frac{\partial L}{\partial p_i} \right) \cdot 1$$

$$\nabla_{b^{(1)}} L = (-4) \cdot 4 = -16$$

$$\nabla_{b^{(2)}} L = (4) \cdot 4 = 16$$

$$\nabla_b L = \nabla_{b^{(1)}} L + \nabla_{b^{(2)}} L$$

$$= 0$$

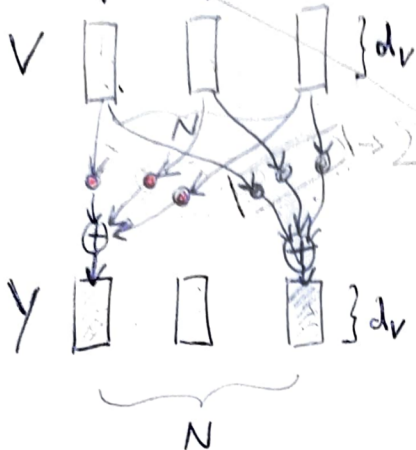
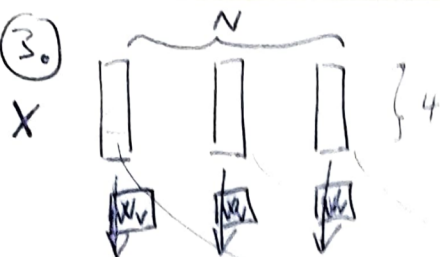
3. Consider a complex datum X that denotes a collection of N tokens of dimension 4 that should be processed with single-head self-attention. The keys and queries both correspond to the first two dimensions of the input tokens while the values correspond to the last two dimensions of the input tokens.

Determine the number of parameters of the layer.

Propose an implementation in numpy to recover the result X_{att} consisting of:

- (a) implementation of the function `self_attn(Q, K, V)`
- (b) definition of the projection matrices W_q, W_k, W_v
- (c) correct invocation of the function `self_attn` that solves the task

3.



$$W_v \downarrow [d_{in}, d_v]$$

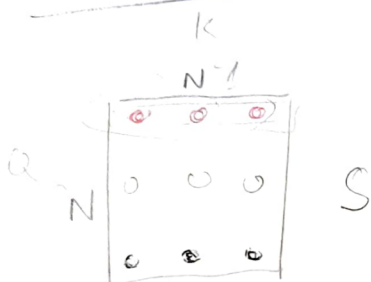
$$X \rightarrow [N, d_{in}]$$

$$V \rightarrow [N, d_v]$$

$$V = X W_v$$

$$Q \rightarrow \{ \text{columns} \} d_k$$

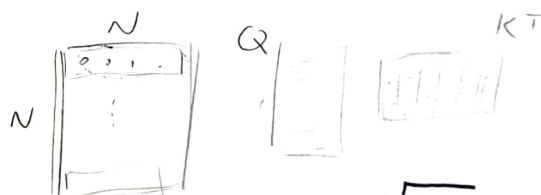
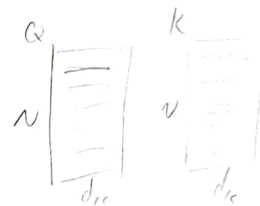
$$K \rightarrow \{ \text{columns} \} d_k$$



$$A = \text{Softmax}(S, \text{axis}=1)$$

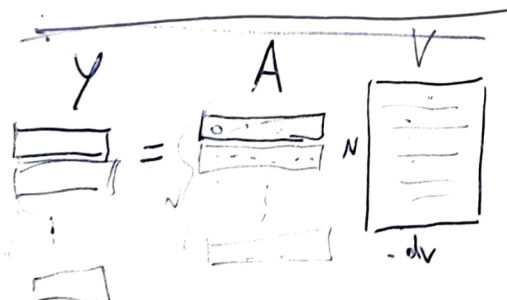
$$[N, d_k] \quad Q = X W_q$$

$$[N, d_k] \quad K = X W_k$$

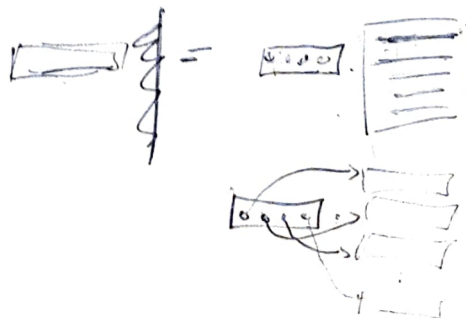


$$S = Q K^T / (\dots) \rightarrow \sqrt{d_k}$$

$$A = \text{softmax}(S, \text{axis}=1)$$



$$Y = AV$$



a) def self_attn(Q, K, V):

S = (Q @ K.T) / np.sqrt(Q.shape[1])
A = sp.special.softmax(S, axis=1)
return A @ V

b)



$N, 4, 4, 2$
~~X~~ $X W_Q$

$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \leftarrow W_Q = W_K$

$W_Q = W_K = \text{np.array}([[1, 0], [0, 1], [0, 0], [0, 0]])$
 $W_V = \dots \dots \dots \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \leftarrow W_V$

c) $X_{ATT} = \text{self_attn}(X @ W_Q, X @ W_K, X @ W_V)$

of param.s :

$W_Q \rightarrow 4 \cdot 2$
 $W_K \rightarrow 4 \cdot 2$
 $W_V \rightarrow 4 \cdot 2$
 $\left. \begin{matrix} \rightarrow 8 \\ \rightarrow 8 \\ \rightarrow 8 \end{matrix} \right\} 24$

4. We consider siamese metric learning with triplet loss. There are 7 examples in the train set. Matrix E contains their embeddings and semantic classes. Each row of the matrix corresponds to one example, with the embedding given by the first two elements, and the third element designating the class.

Tasks:

- (a) Define the L2 distance between two vectors $d_{L2}(x, z)$. Fill in the missing values in the following definition of triplet loss if you are given the embeddings of the anchor a , the positive example p , and the negative example n .

$$L(a, p, n) = \max(d_{L2}(_, _)^2 - d_{L2}(_, _)^2 + m, _)$$

$$E = \begin{bmatrix} -1 & 4 & 0 \\ -1 & 2 & 0 \\ 3 & -3 & 1 \\ 1 & -1 & 1 \\ 1 & 4 & 2 \\ 4 & 1 & 2 \\ 1 & 1 & 2 \end{bmatrix}$$

- (b) Classify the embedding $z = [0, 2]$ using the training set embeddings by comparing with class centroids as in Lab 4. Show the calculation procedure.
- (c) We consider using the learned metric embedding model as part of an image retrieval system. For any given query, the system returns the **five** most similar examples, according to the inverse L1 distance. Let the input query to the system be the embedding $q = [1, 0]$, which belong to class $y = 2$. Which examples would be retrieved by the system? What would be the precision (P) and recall (R) for such a response.

5. Your task is to use a vanilla recurrent network with a ReLU nonlinearity for the task of tracking a users' bank account balance. You will use a 3-dimensional hidden state representation, where the first element will track the total amount of money available for withdrawal (account balance + the allowed overdraft), the second element will store the allowed overdraft, while the third element will serve as a flag checking whether the account is in negative balance (0 if the account **is** in negative balance, any positive value otherwise). The case where a user exceeds the allowed overdraft is currently not of concern.

The initial hidden state of a test users' account is:

$$h^{(0)} = \begin{bmatrix} 10000 \\ 5000 \\ 765 \end{bmatrix}$$

The 2-dimensional input vectors will contain the amount of money added to the account on the first element, and the amount of money withdrawn on the second element.

- (a) Determine the values of the remaining parameters of the recurrent neural network (W_{hh} , W_{xh} , b_h) so that the network would function as intended
- (b) Run the forward pass of the network with the determined parameters for the following inputs:

$$x^{(1)} = \begin{bmatrix} 1000 \\ 0 \end{bmatrix} \quad x^{(2)} = \begin{bmatrix} 1000 \\ 3500 \end{bmatrix} \quad x^{(3)} = \begin{bmatrix} 0 \\ 4000 \end{bmatrix}$$

- (c) Consider using an LSTM network for the same problem, which has the same input representation dimension and hidden state dimension. Determine the total number of parameters of such a network.