# HomeWork 2

João André Roque Costa – 99088

Tomás Augusto Vilhena de Oliveira – 90781

## Question 1

### Q1.1

1) $\quad Z = \text{Softmax}\left(QK^{T}\right)V$

$Q, K, V : (L \times D)$

$Q\,K^{T} \longrightarrow$ [matrix $L \times D$] [matrix $D \times L$] $\rightarrow O\left(L \times D \times L\right)$

$\downarrow$

$\text{Softmax}\left(QK^{T}\right) \rightarrow \text{Softmax}\left([\text{matrix } L \times L]\right) \xrightarrow{O(L \times L)} O(L \times D \times L)$

$\downarrow$

$\text{Softmax}\left(QK^{T}\right)V = [\text{matrix } L \times L] \cdot [\text{matrix } L \times D] \xrightarrow{O(L \times L \times D)} O\left(L^{2} \times D\right)$

$Z \rightarrow$ quadratic over sequence length, gets exponentially slower.

Q1.2

2) $\left[1 \quad q^T h \quad \dfrac{q^T h^2}{2}\right]$ $\quad$ $D = 2, \quad q = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad K = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}$

$\text{trace}(q q^T K K^T) = \text{trace}\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\begin{bmatrix} x_1 & x_2 \end{bmatrix}\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}\begin{bmatrix} Y_1 & Y_2 \end{bmatrix}\right)$

$= \text{trace}\left(\begin{bmatrix} x_1^2 & x_1 x_2 \\ x_2 x_1 & x_2^2 \end{bmatrix}\begin{bmatrix} Y_1^2 & Y_1 Y_2 \\ Y_2 Y_1 & Y_1^2 \end{bmatrix}\right)$

$\begin{bmatrix} 1 & , & q^T K \end{bmatrix}$

$= \text{trace}\left(\begin{bmatrix} x_1^2 Y_1^2 + x_1 x_2 Y_1 Y_2 & \\ & x_1 x_2 Y_1 Y_2 + x_2^2 Y_2^2 \end{bmatrix}\right)$

$= x_1^2 Y_1^2 + x_2^2 Y_2^2 + 2 x_1 x_2 Y_1 Y_2$

$\underbrace{\text{Vec}(q q^T)}_{\phi(q)}{}^T \underbrace{\text{Vec}(K K^T)}_{\phi(R)}$

$\phi(t) = \text{Vec}(t t^T)$

$\begin{array}{ccccc} K & 1 & 2 & 3 & n \\ \dim M & 1 & D & D^2 & \underline{\underline{D^{n-1}}} \end{array}$ $\qquad$ $\text{Vec}\left(1\begin{bmatrix} 1 \\ \end{bmatrix}^1\begin{bmatrix} 2 \\ \end{bmatrix}\right)$

$\text{Vec}\left(1\begin{bmatrix} 2 \\ \end{bmatrix}\right)$

$1\begin{bmatrix} 4 \\ \vdots \end{bmatrix}$

$\text{Dim}_M = D^{K-1}$

3) $D^{-1} \phi(Q) \bar{\phi}(K)^T V$

$= \dfrac{\phi(Q) \bar{\phi}(K)^T}{D} V$

$\approx \dfrac{\exp(QK^T)}{D} V$

$\approx \dfrac{\exp(QK^T)}{\text{Diag}(\phi Q \bar{\phi} K^T 1_L)} V$

$\approx \dfrac{\exp(QK^T)}{\text{Diag}(\exp(QK^T) 1_L)} V$

$\approx \dfrac{\exp(QK^T)}{\sum\limits_{i=1}^{L} \exp(q_i K_i^T)} V \qquad \text{Softmax}(t_i) = \dfrac{\exp(t_i)}{\sum\limits_{j=1}^{n} e^{t_i}}$

$\approx \text{Softmax}(QK^T) V \qquad Z = \text{Softmax}(QK^T) V$

$\approx Z$

Q1.4

4) $Z = D^{-1} \underbrace{\phi(Q)}_{L \times M} \underbrace{\phi(K)^T}_{L \times M^T} \underbrace{V}_{L \times D}$

$\underset{\alpha(L \times M)}{\overset{L}{\downarrow}}$

$\underbrace{\qquad}_{O(L \times M \times D)}$

$\underbrace{\qquad\qquad}_{O(L \times M \times D)}$

$\underbrace{\qquad\qquad\qquad}_{O(L \times M \times D)}$

$M \begin{bmatrix} \overbrace{[\cdot \cdots \cdot]}^{L} \end{bmatrix} \quad L\begin{bmatrix} \overbrace{|\,|\,|\,|}^{D} \end{bmatrix}$

$M \times L \times D$

$Z \to (L, M, D)$

Linearly dependent of $L, M$ and $D$

$D = Diag\left( \phi(Q) \underbrace{\phi(B)^T 1_L}_{\alpha(L \times M)} \right)$

$\underbrace{\qquad\qquad}_{O(L \times M)}$

$\underbrace{\qquad\qquad\qquad}_{O(L \times M)}$

# Question 2
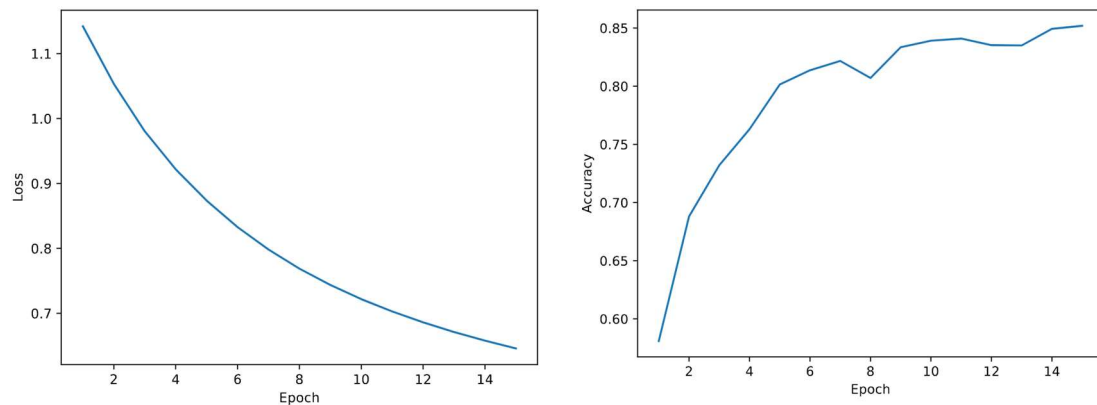
## Q2.1

| 0.1 | 0.01 | 0.001 |
|-----|------|-------|



The best configuration is with a learning rate equal to 0.01 with a final loss of around 0.55 and 0.86 accuracy.

## Q2.2



## Q2.3

no_maxpool: False

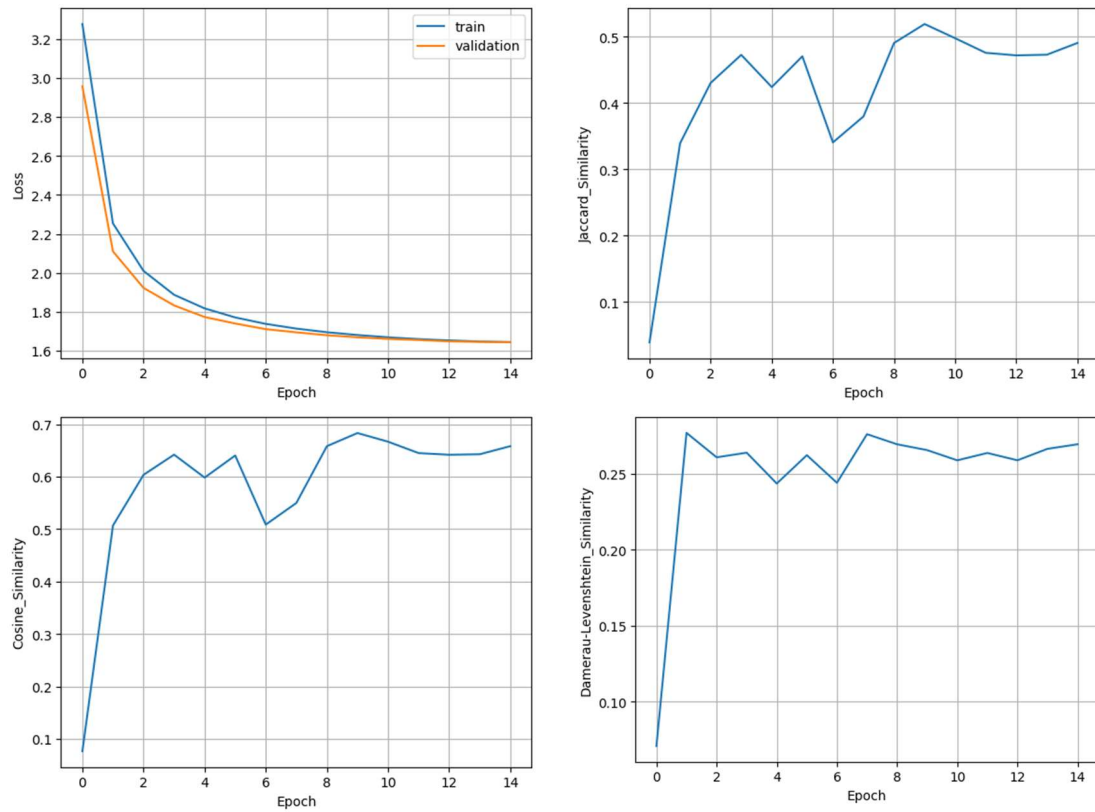Number of trainable parameters: 168572

no_maxpool: True

Number of trainable parameters: 224892

Having more parameters can enhance the model's capacity, enabling it to capture intricate patterns in the data. Nevertheless, an elevated parameter count also escalates the likelihood of overfitting, particularly when the dataset lacks sufficient size.
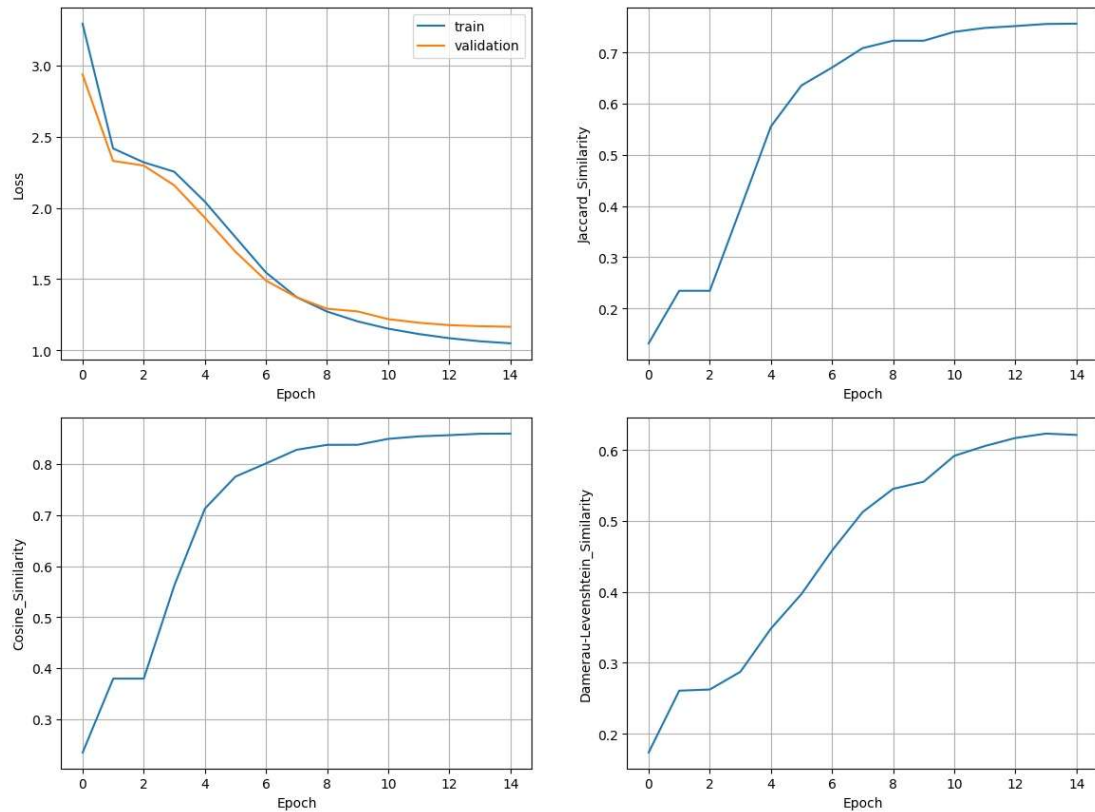
# Question 3

## Q3.1



{'jaccard_similarity': 0.503583607220501, 'cosine_similarity': 0.669388120203678, 'damerau-levenshtein_similarity': 0.2736332793584701, 'loss': 1.6690384796479854}

## Q3.2



{'jaccard_similarity': 0.7636617752753647, 'cosine_similarity': 0.8647155264952522, 'damerau-levenshtein_similarity': 0.6314035820584282, 'loss': 1.160844937330339}

### Q3.3

Unlike the sequential input processing of the LSTM, the attention mechanism allows the model to focus on different parts of the input sequence simultaneously. It assigns different weights to different positions in the input sequence, emphasizing the relevant parts for predicting the current token.

### Q3.4

There is not a very big difference between the three measures in use. When maximizing one, the other will come to a near maximum state as well. In particular, Jaccard Similarity, when calculating the percentage of common elements in a sequence, can be considered alike Demerau-Levenstein which counts the number of edits for two sequences to reach equality.

Finally, the cosine turns the text sequences into vectors of words and its frequencies and calculates the cosine between two given vectors.

While Tomás Oliveira came up with the coding resolutions to the second and third questions, João Costa solved the hole question 1, ran the necessary code and created the report.