

# Deep Learning for Natural Language Processing

Juan LONDONO

January 11th 2019

## 1 Monolingual embeddings (/6)

See Code.

## 2 Multilingual word embeddings (/4)

### Question 1

Using the orthogonality and the properties of the trace, prove that, for  $X$  and  $Y$  two matrices:

$$W^* = \underset{W}{\operatorname{argmin}} \|WX - Y\|_F = UV^T \text{ with } U\Sigma V^T = \operatorname{SVD}(YX^T)$$

Let us develop the Frobenius norm  $\|WX - Y\|_F$  :

$$\|WX - Y\|_F = (\operatorname{tr}((WX - Y)^T(WX - Y)))^{1/2}$$

Therefore we only need to minimize  $\operatorname{tr}((WX - Y)^T(WX - Y))$

$$\operatorname{tr}((WX - Y)^T(WX - Y)) = \operatorname{tr}((X^T W^T - Y^T)(WX - Y))$$

$$\iff \operatorname{tr}(X^T X - X^T W^T Y - Y^T W X + Y^T Y)$$

$$\iff \operatorname{tr}(X^T X - X^T W^T Y - (X^T W^T Y)^T + Y^T Y)$$

As  $\operatorname{tr}(X^T X)$  and  $\operatorname{tr}(Y^T Y)$  are constants, we need to minimize:

$$\operatorname{tr}(-X^T W^T Y - (X^T W^T Y)^T) = -\operatorname{tr}(X^T W^T Y) - \operatorname{tr}(X^T W^T Y)^T$$

$$\iff -2 * \operatorname{tr}(X^T W^T Y)$$

We therefore need to maximize  $\operatorname{tr}(Y^T W X)$ .

We know that  $\operatorname{tr}(AB) = \operatorname{tr}(BA)$ , hence:

$$\operatorname{tr}(Y^T W X) = \operatorname{tr}(XY^T W) = \operatorname{tr}((YX^T)^T W) = \operatorname{tr}((U\Sigma V^T)^T W) = \operatorname{tr}(V\Sigma U^T W) = \operatorname{tr}(U^T W V \Sigma)$$

Let us write  $U^T W V = M$ .

We know that  $tr(M\Sigma) = \sum_i (M_{i,i}\Sigma_{i,i}) \leq \sum_i (\Sigma_{i,i})$ , since  $M_{i,i} \leq 1$  for all  $i$ , as  $M$  is orthogonal.

Therefore  $tr(M\Sigma)$  is maximized when each  $M_{i,i} = 1$ .  
 $\iff M = Id \iff W = UV^T$

### 3 Sentence classification with BoV (/4)

#### Question 1

What is your training and dev errors using either the average of word vectors or the weighted-average?

	Training Accuracy	Dev Accuracy
<b>Word vectors average</b>	0.4993561980568887	0.43272727272727274
<b>Weighted average</b>	0.497600374575676	0.4090909090909091

As we may see, our model tends to overfit since the training accuracy is 6 to 9 points higher than the dev accuracy. We also notice that the word vectors model performs 3 points better than the weighted average model in the dev set. This is the one that we chose to predict the test results.

### 4 Deep Learning models for classification (/6)

#### Question 1

Which loss did you use? Write the mathematical expression of the loss you used for the 5-class classification.

For the loss of the classifier, I used the categorical cross entropy. This loss is well adapted to our model since it is useful for multiclass classification problems. It returns the cross-entropy between an approximating distribution and a true distribution.

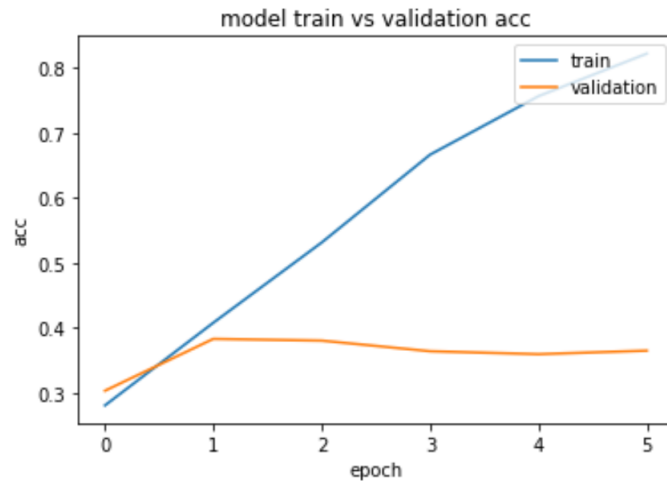
This function computes

$$H(p, q) = -\sum_x p(x) \log(q(x))$$

Where  $p$ =true dist and  $q$ =coding dist.

## Question 2

Plot the evolution of train/dev results w.r.t the number of epochs.



As we may see, our model overfits very quickly. The maximum accuracy for the validation set is obtained in 1 or 2 epochs and then decreases, while the train accuracy continues to increase steadily.

## Question 3

Be creative: use another encoder. Make it work! What are your motivations for using this other model?

I decided to use another encoder. I used the Tokenizer function to preprocess the text.

Each input line would therefore have a size (vocab size) instead of (52). This higher dimension might lead to enhanced prediction.

I used a classical two-layer dense NN, which improved the previous results by 4 points.

Any attempt to perform a Conv1D or an LSTM lead to a very high computational time, and no extra performance.