

Deep Learning for Natural Language Processing

Armand Margerin

February 2020

1 Multilingual Word Embeddings

We have

$$\begin{aligned}\|WX - Y\|_F &= \text{tr}((WX - Y)^T(WX - Y)) \\ &= \text{tr}(X^T W^T W X - X^T W^T Y - Y^T W X + Y^T Y) \\ &= \text{tr}(X^T W^T W X) - 2\text{tr}(X^T W^T Y) + \text{tr}(Y^T Y) \\ &= \text{tr}(X^T X) - 2\text{tr}(X^T W^T Y) + \text{tr}(Y^T Y)\end{aligned}$$

Hence $\text{argmin}_{W \in O_d(R)}(\|WX - Y\|_F) = \text{argmax}_{W \in O_d(R)}(\text{tr}(X^T W^T Y))$

We use the singular value decomposition of YX^T :

$$SVD(YX^T) = U\Sigma V^T$$

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

The matrix $V^T W^T U$ is orthogonal since U and V are unitary matrices:

$$U^T W V V^T W^T U = U^T W W^T U = U^T U = I$$

Hence, $\text{tr}(V^T W^T U \Sigma) \leq \text{tr}(\Sigma)$ and we get the maximum value for $V^T W^T U = I$, i.e $W = UV^T$.

2 Sentence Classification with BoW

Average of word vectors:

Train accuracy: 46,0%

Dev accuracy: 39,4%

Weighted average of word vectors:

Train accuracy: 46,5%

Dev accuracy: 39,2%

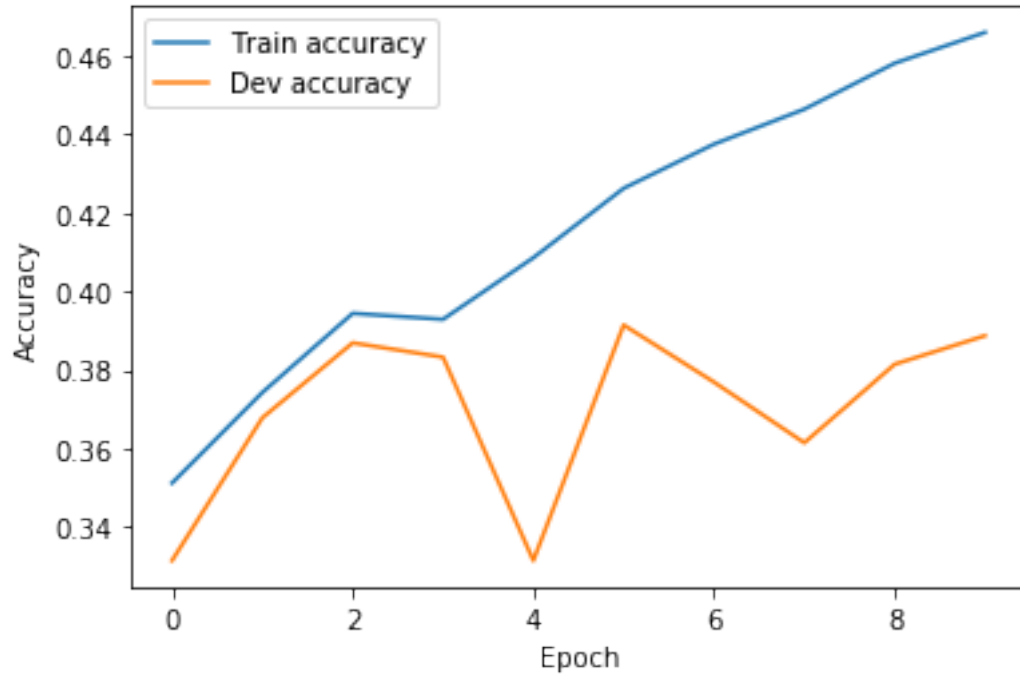
3 Deep Learning models for classification

I chose to use the categorical-crossentropy loss because it is a multi-class classification problem. It is defined by:

$$l(m, n) = \sum_{i=1}^5 m_i \log(n_i)$$

where m is the true class and n the predicted class

Evolution of train/dev results:



Other encoder:

I tried using the same encoder with pre-trained vectors from Word2vec, but it resulted in a poor performance. One possible reason for that is the dimension of the embeddings (300): with this dimension I also had bad results in the original model.