

Deep Learning

NLP Project

MVA

Wassim Bouaziz - wassim.s.bouaziz@gmail.com

1^{er} mars 2020

1 Monolingual embeddings

See `nlp_project.ipynb`.

2 Multilingual word embeddings

Question - We want to solve the problem :

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

And

$$\begin{aligned} \arg \min_{W \in O_d(\mathbb{R})} \|WX - Y\|_F &= \arg \max_{W \in O_d(\mathbb{R})} \text{tr}(X^T W^T Y) \\ &= \arg \max_{W \in O_d(\mathbb{R})} \langle W, Y X^T \rangle_F \\ &= \arg \max_{W \in O_d(\mathbb{R})} \langle W, U \Sigma V^T \rangle_F \\ &= \arg \max_{W \in O_d(\mathbb{R})} \langle U^T W V, \Sigma \rangle_F \end{aligned}$$

The matrix $U^T W V$ is also orthogonal, and the Cauchy-Schwartz inequality gives us :

$$\langle U^T W V, \Sigma \rangle_F \leq \|\Sigma\|_F$$

With equality only if $U^T W V$ is colinear to Σ , hence, if $U^T W V$ is diagonal, so $U^T W V = I_d$, i.e. $W = U V^T$.

	train accuracy	dev accuracy
with idf weight	41.09%	40.05%
without idf weight	48.09%	43.50%

TABLE 1 – Accuracy for train ($C = 10$) and dev ($C = 0.3$) sets.

3 Sentence classification with BoW

Question - Using the BoW and the LogisticRegression model from `scikit-learn`, we get the following results :

4 Deep Learning models for classification

Question - I chose the *categorical crossentropy* which mathematical expression is :

$$L(y) = \sum_{c=1}^C \mathbf{1}_{y=c} \log(\mathbb{P}(y = c))$$

With C being the number of classes.

Question - We can see on the Figure 1 that the model is greatly overfitting on the training set.

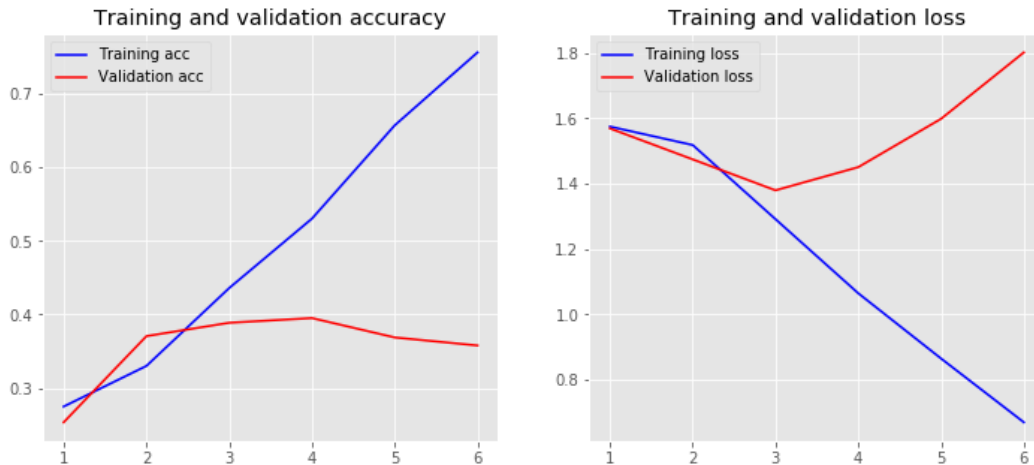


FIGURE 1 – Accuracy and loss on training and development set with the LSTM model.

Question - In the sentence, in order to determine sentiments, we want to be able to extract key words and key expressions which could help us. If using a LSTM allowed the model to observe the whole sentence, looking at distant groups of words can be sufficient. For that purpose, we stack convolutional layers (even tried stacking them with LSTM layers) to be able to extract information in the whole sentence and capture it for dense layers to make the final classification. We notice however that this model does not improve the validation accuracy (see Figure 2). But notice that we used

early stopping and that the training accuracy is higher than it was with the LSTM model, meaning it overfits even more rapidly.

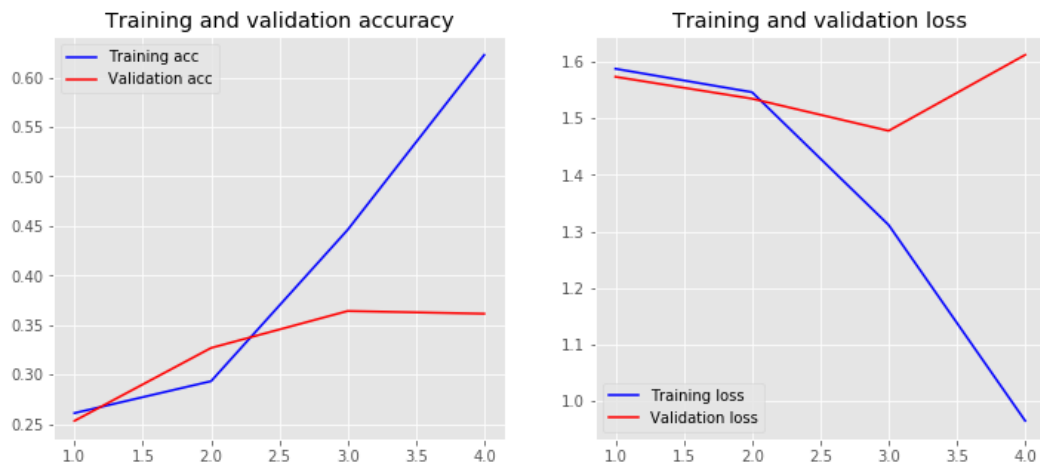


FIGURE 2 – Accuracy and loss on training and development set with the Convolutional model.