

---

**TP 5:** Deep Learning for Natural Language Processing

---

victor.busa@ens-paris-saclay.fr

## 1 Monolingual embeddings

## 2 Multilingual word embedddings

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

$$W^* = \arg \min_{W \in O_d(\mathbb{R})} \|WX - Y\|_F = UV^T, \text{ with } U\Sigma V^T = \text{SVD}(YX^T) \quad (1)$$

**Answer:** We want to minimize  $\|WX - Y\|_F$ . Without loss of generality we can focus on minimizing  $\|WX - Y\|_F^2$ . Let  $X, Y \in \mathbb{R}^{(d \times n)}$ . We then have:

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

Where we have used the fact that  $W \in O_d(\mathbb{R})$ , i.e  $W^T W = I$ . So finally we have that  $(\text{tr}(X^T X) + \text{Tr}(Y^T Y))$  being constant w.r.t  $W$ :

$$\begin{aligned} \min_{W \in O_d(\mathbb{R})} \|WX - Y\|_F &= \max_{W \in O_d(\mathbb{R})} \text{Tr}(YX^T W^T) \\ &= \max_{W \in O_d(\mathbb{R})} \text{Tr}(U\Sigma V^T W^T U) \end{aligned}$$

Where we have used the singular value decomposition of  $YX^T$ :  $U\Sigma V^T = \text{SVD}(YX^T)$

Now, as  $U, V, W$  are in  $O_d(\mathbb{R})$ , so their products  $\widehat{W} = V^T W^T U$  is in  $O_d(\mathbb{R})$  because  $(O_d(\mathbb{R}), \times)$  is a group.

So finally, we have:

$$\begin{aligned} \min_{W \in O_d(\mathbb{R})} \|WX - Y\|_F &= \max_{W \in O_d(\mathbb{R})} \text{Tr}(YX^T W^T) \\ &= \max_{W \in O_d(\mathbb{R})} \text{Tr}(\Sigma \widehat{W}) \\ &= \max_{W \in O_d(\mathbb{R})} \sum_{i=0}^d \sigma_{ii} \widehat{w}_{ii} \end{aligned}$$

Where we have used the fact that  $\Sigma$  is a diagonal matrix. So at the end to maximize the sum  $\sum_{i=0}^d \sigma_{ii} \hat{w}_{ii}$ , We should choose  $\hat{W} = I_{d \times d}$  because  $\hat{W} \in O_d(\mathbb{R})$  and  $\forall i \in [1, \dots, d], \sigma_{ii} \geq 0$ .

Hence:  $\hat{W} = V^T W^T U = I \Rightarrow W^T = V U^T \Rightarrow W = U V^T$

So, at the end we have that the maximum is attained for  $W = U V^T$  where  $U \Sigma V^T = \text{SVD}(Y X^T)$

### 3 Sentence classification with BoV

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

**Answer:** Quite weirdly we can see on table 1 that the logistic regression using mean average achieves a better accuracy on the testing set then if we use the *idf* average. Maybe it would have more sense to use a *idf-tf* weighted average and not just an *idf* weighted average.

accuracy	mean average	idf average
<b>Train</b>	0.497	0.493
<b>Test</b>	0.447	0.415

Table 1: Accuracy on the Stanford Sentiment Treebank using logistic regression and a BoV with either mean average or idf weighted average

**Bonus question:** I have used the xgboost classifier as it usually leads to good results (cf Kaggle competition). Yet I wasn't able to achieve better results with it.

### 4 Deep Learning model for classification

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

**Answer:**

- This is a multi-class classification problem so I use a **categorical cross-entropy** loss function.
- The mathematical expression is just (for 5 classes):

$$H(p, q) = - \sum_{i=1}^5 y_i \log(y'_i)$$

where  $y_i$  is the true probability distribution and  $y'_i$  is the predicted probability

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

**Answer:** the evolution of the train/dev results w.r.t the number of epochs is shown on Figure 4.1

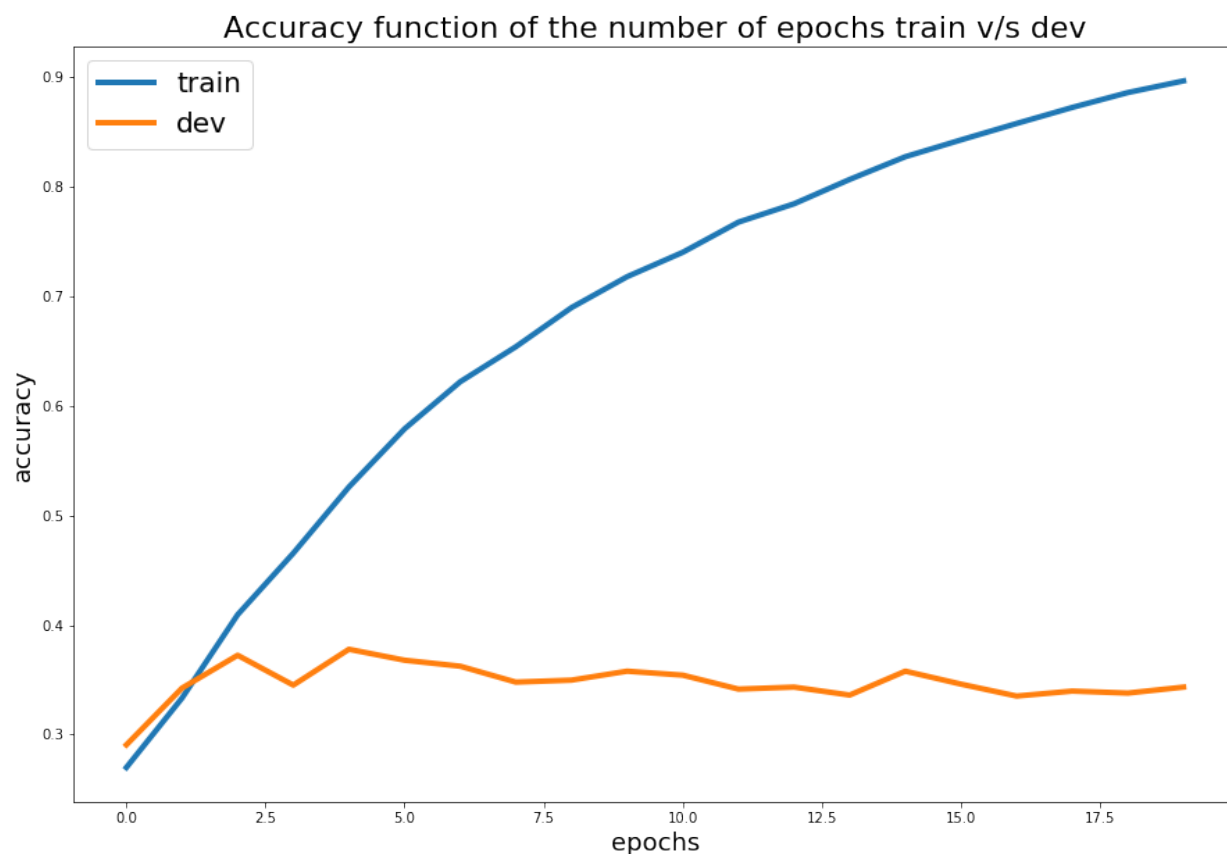


Figure 4.1: Evolution of train and dev accuracy w.r.t the number of epochs

**Question** Be creative: use another encoder. What are your motivations for using this other model?

We can use pretrained word-embeddings. The neural network will converge faster as it doesn't need to learn the word-embeddings from scratch. At the end I have a better model, yet it is unable to achieve outstanding performance overall. Figure 4.2 shows the evolution of the validation and training accuracy w.r.t to the number of epochs.

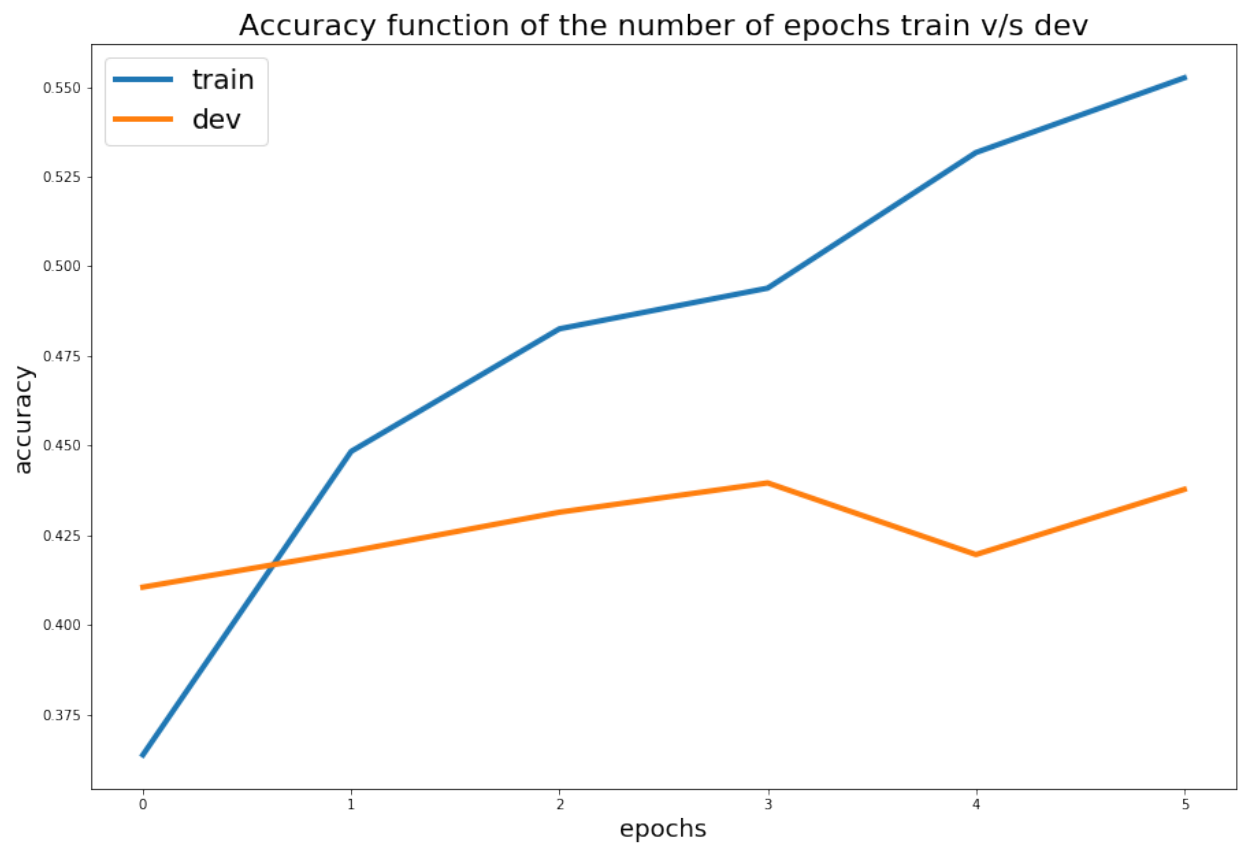


Figure 4.2: Evolution of train and dev accuracy w.r.t the number of epochs on the final model