

Multilingual word embeddings

Question

Let's show that

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

Let $W \in O_d(\mathbb{R})$. We first notice that :

$$\begin{aligned} \|WX - Y\|_F^2 &= \|WX\|_F^2 + \|Y\|_F^2 - 2\langle WX, Y \rangle \\ &= \|X\|_F^2 + \|Y\|_F^2 - 2\langle WX, Y \rangle \end{aligned}$$

Where $\langle WX, Y \rangle = \operatorname{Tr}((WX)^T Y) = \operatorname{Tr}(W^T Y X^T)$. Hence,

$$\begin{aligned} \operatorname{argmin}_{W \in O_d(\mathbb{R})} \|WX - Y\|_F &= \operatorname{argmin}_{W \in O_d(\mathbb{R})} \|WX - Y\|_F^2 \\ &= \operatorname{argmax}_{W \in O_d(\mathbb{R})} \operatorname{Tr}(W^T Y X^T) \end{aligned}$$

Let $U\Sigma V^T$ be the SVD decomposition of YX^T .

$$\begin{aligned} \operatorname{Tr}(W^T Y X^T) &= \operatorname{Tr}(W^T U \Sigma V^T) \\ &= \operatorname{Tr}(V^T W^T U \Sigma) \end{aligned}$$

Let $Z = W^T U \Sigma$. Z is orthogonal as the product of orthogonal matrices. $\operatorname{Tr}(Z\Sigma) = \sum_{i=1}^n Z_{i,i} \Sigma_{i,i}$. $(\Sigma_{i,i})_i$ are non-negative and Z is orthogonal hence the previous quantity is maximized if $Z_{i,i} = 1$ which means that $Z = I_n$. We then conclude that $W^* = UV^T$.

Sentence classification with BoW

Question

We test the sentence classification task using both the average of word vectors and the idf-weighted-average .

We get the following results :

Dataset	Average	idf
Train	0.4641	0.4904
Dev	0.4078	0.4078

Table 1: BoW Accuracy

We get better results using the idf-weighted average of word vectors on Train set, on dev set we get the same performance, we'll keep the idf-weighted results.

Deep Learning models for classification

0.1 Question

We use the cross-entropy loss function (per Pytorch's documentation) :

$$\text{loss}(x, \text{class}) = -\log \left(\frac{\exp(x[\text{class}])}{\sum_j \exp(x[j])} \right) = -x[\text{class}] + \log \left(\sum_j \exp(x[j]) \right)$$

It is good to notice here the we don't need any one-hot coding or a softmax layer on our model as all of this is taken care of by the CrossEntropyLoss in pytorch.

0.2 Question

The model seems to start overfitting quickly as we can see that the Validation error starts going up while Train error keeps going down.



Figure 1: Train and Validation Losses .

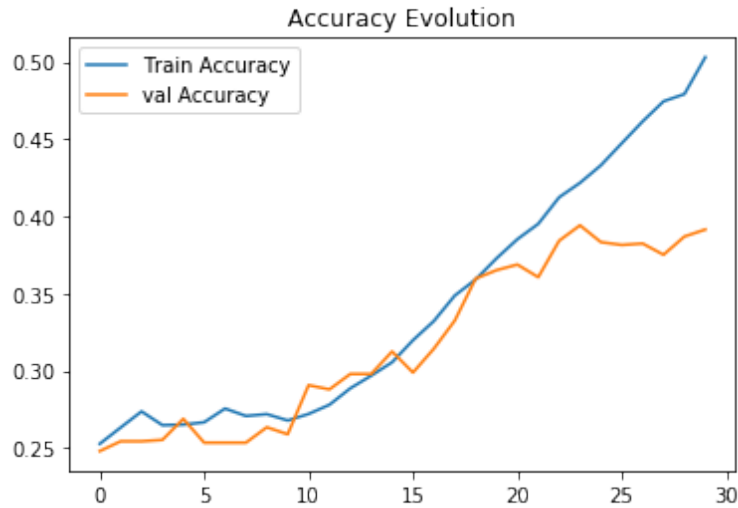


Figure 2: Train and Validation accuracy .

0.3 Innovate

model

we'll Use more or less the same model as before the only change we'll introduce is by using pre-trained word embeddings (we'll use Glove, which is publically available at [Glove Project](#))

One thing worth mentioning is that since we use pretrained embedding we need to recreate our datasets such that the words and indexes match the ones used in the pretrained embedding. We get a better performance on dev set by using this as the accuracy goes from **40.3%** to **45.04%**