

Deep Learning for Natural Language Processing

M2 MVA Deep Learning

Tong ZHAO (tong.zhao@eleves.enpc.fr)

1 Multilingual Word Embeddings

(Q1) The loss function is defined as following:

$$l(W) = \|WX - Y\|_F \quad (1)$$

$$= \text{tr}\left((WX - Y)^T(WX - Y)\right) \quad (2)$$

$$= \text{tr}\left(X^T W^T W X - Y^T W X - X^T W^T Y + Y^T Y\right) \quad (3)$$

$$= \text{tr}(X^T W^T W X) + \text{tr}(Y^T Y) - \text{tr}(Y^T W X) - \text{tr}((Y^T W X)^T) \quad (4)$$

$$= \text{tr}(X^T X) + \text{tr}(Y^T Y) - 2\text{tr}(X^T W^T Y) \quad (5)$$

Since X and Y are constant matrices, minimizing the objective function $l(W)$ can be solved by maximizing the function :

$$f(W) = \text{tr}(X^T W^T Y) = \text{tr}(W^T Y X^T)$$

We apply an SVD on the $m \times m$ square matrix $Y X^T$. Thus we have:

$$Y X^T = U \Sigma V^T$$

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

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

Let $H = V^T W^T U$, we have:

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

Since H is orthogonal, we have $\|h_{ii}\| \leq 1$, so:

$$\text{tr}(H \Sigma) = \sum_{i=1}^m h_{ii} \sigma_i \leq \sum_{i=1}^m \sigma_i = \text{tr}(I \Sigma)$$

The upper bound is achieved when $H = I$, thus we have:

$$V^T W^T U = I \implies W = (V U^T)^T = U V^T$$

2 Sentence Classification with BoV

(Q2)

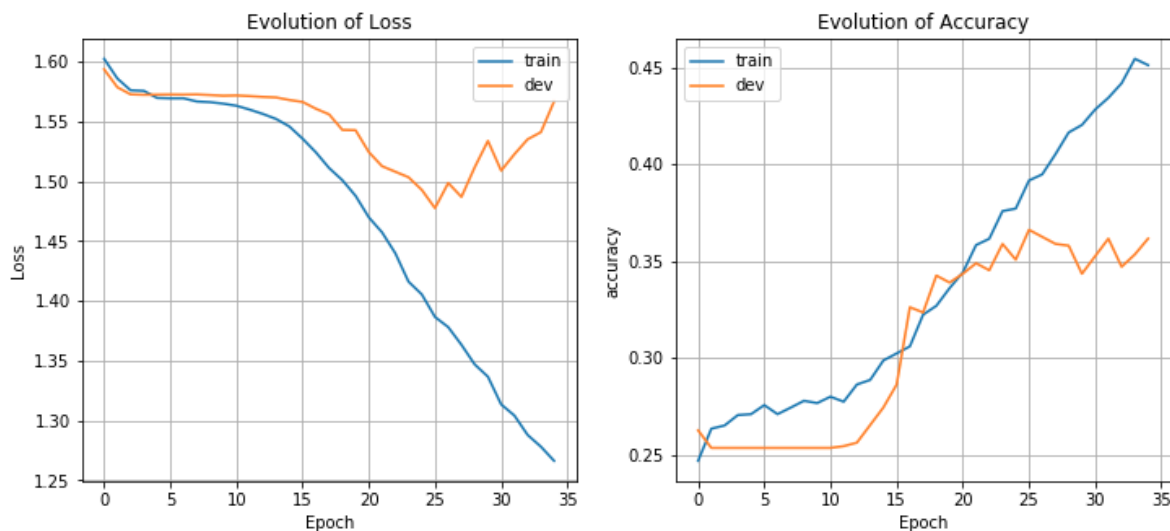
Method	Dev Error	Train Error
Average	55.677%	50.269%
Weighted Average	56.676%	52.797%

3 Deep Learning Models for Classification

(Q3) I use the categorical crossentropy loss in the model. Given the predicted distribution q and the groundtruth distribution p , the formula is as following:

$$l(p, q) = \sum_{i=1}^5 p_i \log q_i$$

(Q4)



(Q5) The model I used includes following features:

- A pretrained embedding matrix gives a good initialization to the first embedding layer, so it makes the convergence faster.
- A bidirectional LSTM lets the network get more information from its context.
- A 2D convolutional layer enlarges its receptive field.