

# Deep Learning for Natural Language Processing

Louis TREZZINI

23 décembre 2018

## 2 Multilingual word embeddings

### Question

We want to minimize  $\|WX - Y\|$  under the constraint that  $W \in O_d(\mathbb{R})$ .

We notice :

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

So minimizing  $\|WX - Y\|$  is equivalent to maximizing  $\langle WX, Y \rangle$ .

Let's write  $U\Sigma V^T = SVD(YX^T)$ . We have :

$$\begin{aligned}\langle WX, Y \rangle &= \langle W, YX^T \rangle \\ &= \langle W, U\Sigma V^T \rangle \\ &= \langle U^T W V, \Sigma \rangle\end{aligned}$$

Let's denote  $Q = U^T W V$ . We have :

$$\begin{aligned}\langle Q, \Sigma \rangle &= \text{Trace}(Q^T \Sigma) \\ &= \sum_{k=1}^d Q_{kk} \Sigma_{kk}\end{aligned}$$

Because  $\Sigma$  is diagonal.

Notice that  $Q = U^T W V \in O_d(\mathbb{R})$  because  $O_d(\mathbb{R})$  is a group and  $U, V, W$  belong to it. It implies that  $Q_{kk} \in [-1, 1]$ .

Thus, this expression is maximized when  $Q = I$ , since we need to choose  $Q_{kk} = 1$  for all  $k$  because the entries of  $\Sigma$  are non-negative.

So :

$$Q = I \Leftrightarrow U^T W V = I \Leftrightarrow W = UV^T$$

### 3 Sentence classification with BoV

#### Question

### 4 Deep Learning models for classification

#### Question 1

I used the categorical cross-entropy loss.

Its formula is

$$-\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C \mathbf{1}_{\{y_i=c\}} \log(p(\hat{y}_i = c))$$

where  $N$  is the number of examples,  $C$  is the number of classes,  $y_i$  is the true label of example  $i$ , and  $\hat{y}_i$  is the predicted label of example  $i$ .

#### Question 2

FIGURE 1 – Evolution of train/dev results w.r.t the number of epochs.

#### Question 3