

# Master M2 MVA 2018/2019

## Deep Learning - MP2

Souhaib ATTAIKI

December 29, 2018

### 2 - Multilingual word embeddings

We want to prove that :

$$W^* = \operatorname{argmin}_{W \in \mathcal{O}_d(\mathbb{R})} \|WX - Y\|_F = UV^T$$

with  $U\Sigma V^T = \operatorname{SVD}(YX^T)$ .

This is equivalent to :

$$\begin{aligned} W^* &= \operatorname{argmin}_{W \in \mathcal{O}_d(\mathbb{R})} \|WX - Y\|_F^2 \\ &= \operatorname{argmin}_{W \in \mathcal{O}_d(\mathbb{R})} (\|WX\|_F^2 + \|Y\|_F^2 - 2\langle WX, Y \rangle) \\ &= \operatorname{argmin}_{W \in \mathcal{O}_d(\mathbb{R})} (\|X\|_F^2 + \|Y\|_F^2 - 2\operatorname{Tr}(X^T W^T Y)) \\ &= \operatorname{argmax}_{W \in \mathcal{O}_d(\mathbb{R})} \operatorname{Tr}(X^T W^T Y) \\ &= \operatorname{argmax}_{W \in \mathcal{O}_d(\mathbb{R})} \operatorname{Tr}(Y X^T W^T) \end{aligned}$$

Where we used the fact that  $\|W\|_F^2 = 1$  and that  $\|X\|_F^2 + \|Y\|_F^2$  does not depend on  $W$ .

However, we have  $YX^T = U\Sigma V$  and  $W = U_w \Sigma_w V_w^T$  ( $W \in \mathcal{O}_d(\mathbb{R})$ ), so :

$$\begin{aligned} \operatorname{Tr}(Y X^T W^T) &= \operatorname{Tr}(U\Sigma V^T V_w \Sigma_w U_w^T) \\ &= \operatorname{Tr}(\tilde{U} \Sigma \tilde{V} \Sigma_w) \quad \text{with } \tilde{U} = U_w^T U \text{ and } \tilde{V} = V^T V_w \\ &\leq \operatorname{Tr}(\Sigma \Sigma_w) \quad \text{by applying Von Neumann's trace inequality} \end{aligned}$$

So the optimal value is such that  $\tilde{U} = \mathbb{1}$  and  $\tilde{V} = \mathbb{1}$ . Thereby,  $U = U_w$  and  $V = V_w$ , so since  $W = U\Sigma_w V^T \in \mathcal{O}_d(\mathbb{R})$ , we have  $U\Sigma_w^2 U^T = \mathbb{1}$  so  $\Sigma_w^2 = \Sigma_w = \mathbb{1}$ , which leads us to  $W = UV^T$ .

### 3 - Sentence classification with BoV

See *notebook*.

### 4 - Deep Learning models for classification

**Which loss did you use? Write the mathematical expression of the loss you used for the 5-class classification** I've used the '*categorical\_crossentropy*' loss. The mathematical expression of the loss is :

$$L = \frac{1}{n_{obs}} \sum_{i=1}^{n_{obs}} \sum_{c=1}^5 \mathbb{1}_{y_i \in C_c} \log(p_{model}(y_i \in C_c))$$

where  $y_i$  is the output of the model,  $C_c$  is the category  $c$ , and  $p_{model}(y_i \in C_c)$  is the probability predicted by the model for the 'i'th observation to belong to the  $c$  th category.

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