

Deep Learning for Natural Language Processing :

Project

Adrien Benamira

Instructor:
TA: Abdulkadir Çelikkanat

January 10, 2019

1 Monolingual embeddings

cf code

2 Multilingual word embeddings

Let's demonstrate the orthogonal procrustes Theorem :

The Frobenius norm can be defined as:

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\text{trace}(A^* A)} = \sqrt{\sum_{i=1}^{\min\{m,n\}} \sigma_i^2(A)} \quad (1)$$

and because we have that minimising a norm and the same norm square is the same problem, we get thank to :

$$\|A + B\|_F^2 = \|A\|_F^2 + \|B\|_F^2 + 2\langle A, B \rangle_F \quad (2)$$

That :

$$\text{argmin}_{W \in \mathbb{R}^{m \times n}} (\|WX - Y\|_F^2) = \text{argmin}_{W \in \mathbb{R}^{m \times n}} (\|WX\|_F^2 + \|Y\|_F^2 - 2\langle WX, Y \rangle_F)$$

We use the orthogonal nature of W (that is, $W^T W = W W^T = \mathbf{I}$) and the cyclic nature of the trace ($\text{trace}(XYZ) = \text{trace}(ZXY)$) to show :

$$\|WX\|_F^2 = \text{trace}(X^T W^T W X) = \text{trace}(X^T X) = \|X\|_F^2, \quad (3)$$

and because $\|X\|_F^2$ and $\|Y\|_F^2$ don't depend on W , we have :

$$\text{argmin}_{W \in \mathbb{R}^{m \times n}} (\|WX - Y\|_F^2) = \text{argmax}_{W \in \mathbb{R}^{m \times n}} (\langle WX, Y \rangle_F) \quad (4)$$

Using the definition of the scalar product and the cyclic propriety of the trace, we have :

$$\langle WX, Y \rangle_F = \text{trace}(Y^T W X) = \text{trace}(W X Y^T) \quad (5)$$

We have : $U\Sigma V^T = \text{SVD}(YX^T)$
and so :

$$\langle WX, Y \rangle_F = \text{trace}(WV\Sigma U^T) = \text{trace}(U^T WV\Sigma) = \text{trace}(Z\Sigma) \quad (6)$$

with $Z = U^T WV$. We have Z orthogonal and because Σ is a diagonal matrix,

$$\text{trace}(Z\Sigma) = \sum \sigma_{i,i} z_{i,i} \quad (7)$$

So minimizing $\|WX - Y\|_F^2$ is equivalent to maximizing Z for $\sum \sigma_{i,i} z_{i,i}$.
And because Z is orthogonal, the norm of every row is 1 and so the maximum for $z_{i,i}$ is 1.

$$Z = U^T WV = I_n \quad (8)$$

and Finally :

$$W^* = UV^T \quad (9)$$

3 Sentence classification with BoV

3.1 Question

Accuracy on training and validation set for $C = 0.1$

	Word vector	Weigth average	Word vector normalized	Weigth average normalized
Trainning	0.418	0.445	0.469	0.458
Validation	0.3906	0.391	0.405	0.395

4 Deep Learning models for classification

4.1 Question

$$\frac{-1}{N} \sum_i^N \sum_c^C \mathcal{I}_{y_i \in C_c} \log(y_{predict}^{(i)}) \quad (10)$$

The double sum is over the observations i , whose number is N , and the categories (classes) c , whose number is C

4.2 Question

cf code

4.3 Question

Référence : Convolutional Neural Networks for Sentence Classification, Yoon Kim