

Std-ID: Sol

One of the main limitations of static word embedding's or word vector space models is that they failed to represents out-of-vocabulary term for example a term that is not available on training time when these representations are learned. Moreover, words with multiple meanings are conflated into a single representation (a single vector in the semantic space). In other words, polysemy and homonymy are not handled properly.

Question No.2

How Neural Information Retrieval different from traditional Information Retrieval? Outline some of the benefits of Neural IR? [5]

Traditional IR model are based on assumptions about common representation scheme for both document and query. The retrieval process is very much dependent on these assumptions and does not support contextual retrieval based on semantics of the natural language text generally used in documents.

Neural Information Retrieval (NeurIR) is based on Neural Network. The representation is autonomously learned as per task setting. NeurIR model consists of at least one input and one output layers and possibly zero or more hidden layers with different numbers of neurons in each layers. The connection and weight of these neurons automatically set to optimized through supervised or unsupervised learning for the retrieval task. Hence it is more generalized model.

Benefits:

- Context based word representation can be autonomously learning for any NLP specific task like Ad hoc Information Retrieval.
- More generalized model and It can be quite robust to support long/complex queries

Question No.3

Consider the following examples for the task of text classification [5+5]

	docID	words in document	in $c = \text{China}$?
training set	1	Taipei Taiwan	yes
	2	Macao Taiwan Shanghai	yes
	3	Japan Sapporo	no
	4	Sapporo Osaka Taiwan	no
test set	5	Taiwan Taiwan Sapporo	?

- a. Using feature vectors as term frequency and the k-Nearest Neighbors (KNN) with $k=3$ identify the class of test instance docID=5?

Consider the features

$V = \langle \text{japan; macao; osaka; sapporo; shanghai; taipei; taiwan} \rangle$

$D1 = \langle 0, 0, 0, 0, 1, 1 \rangle$ magnitude $|D1| = \sqrt{2}$

$D2 = \langle 0, 1, 0, 0, 1, 0, 1 \rangle$ magnitude $|D1| = \sqrt{3}$

$$\begin{aligned} D3 &= \langle 1, 0, 0, 1, 0, 0, 0 \rangle & \text{magnitude } |D1| &= \sqrt{2} \\ D4 &= \langle 0, 0, 1, 1, 0, 0, 1 \rangle & \text{magnitude } |D1| &= \sqrt{3} \\ D5 &= \langle 0, 0, 0, 1, 0, 0, 2 \rangle & \text{magnitude } |D1| &= \sqrt{5} \end{aligned}$$

Now distance as (dot-product) with D5 for each document can be given as

$$\begin{aligned} \text{Dot-Product } (D1, D5) &= 2 \\ \text{Dot-Product } (D2, D5) &= 2 \\ \text{Dot-Product } (D3, D5) &= 1 \\ \text{Dot-Product } (D4, D5) &= 3 \end{aligned}$$

For K=3 the three neighbors of D5 will be D1, D2 and D3 from these D1 and D2 majority belong to class c=china is YES

- b. Using the same feature vectors and the Rocchio's algorithm, classify the test instance docID=5?

Consider the two mass vectors of the classes

$$\begin{aligned} \mu(\text{class}=\text{china}=\text{Yes}) &= \frac{1}{2} (D1+D2) = \frac{1}{2} \{ \langle 0, 0, 0, 0, 0, 1, 1 \rangle + \langle 0, 1, 0, 0, 1, 0, 1 \rangle \} \\ &= \langle 0, 1/2, 0, 0, 1/2, 1/2, 1 \rangle \end{aligned}$$

$$\begin{aligned} \mu(\text{class}=\text{china}=\text{No}) &= \frac{1}{2} (D3+D4) = \frac{1}{2} \{ \langle 1, 0, 0, 1, 0, 0, 0 \rangle + \langle 0, 0, 1, 1, 0, 0, 1 \rangle \} \\ &= \langle 1/2, 0, 1/2, 1, 0, 0, 1/2 \rangle \end{aligned}$$

Now angle between vectors from $\mu(\text{class}=\text{china}=\text{Yes})$ to D5 can be computed as using:

$$\text{angle } (\mu(\text{class}=\text{china}=\text{Yes}), D5) = 0.8283$$

$$\text{angle } (\mu(\text{class}=\text{china}=\text{No}), D5) = 0.8283$$

As the D5 has the same angle from the two center mass Rocchio's cannot be able to decide about this example.