

28.10 Avg Word2Vec, TF-IDF Weighted Word2Vec

Average Word2Vec

Let us assume we have a document/review 'r₁'. Let its corresponding average vector form be denoted as 'v₁'. Let the total number of words in 'r₁' be 'n₁'.

r₁: w₁ w₂ w₁ w₃ w₄ w₅

So now the average word2vec for 'r₁' is

$$\mathbf{v}_1 = (1/n_1)[w_2\mathbf{v}(w_1) + w_2\mathbf{v}(w_2) + w_2\mathbf{v}(w_1) + w_2\mathbf{v}(w_3) + w_2\mathbf{v}(w_4) + w_2\mathbf{v}(w_5)]$$

w₂v(w_i) represents the word vector of the word 'w_i'. This is called the Average Word2Vec representation of 'r₁'. Average Word2Vec fairly works well in practice, but is not perfect all the time. It still works well enough.

Average Word2Vec is a simple way to leverage Word2Vector to build sentence vectors.

TF-IDF Weighted Word2Vec

Let us assume we have 7 words (say 'w₁', 'w₂', 'w₃', 'w₄', 'w₅', 'w₆', 'w₇') and the review/document 'r₁'.

r₁: w₁ w₂ w₁ w₃ w₄ w₅

The TF-IDF representation of the 'r₁' vector is given as below

w1	w2	w3	w4	w5	w6	w7
t1	t2	t3	t4	t5	t6=0	t7=0

Here t_i → TF-IDF(w_i, r₁)

Now, let 'v₁' be the TF-IDF Weighted Word2Vec representation of the vector 'r₁', then

$$\mathbf{v}_1 = (t_1 * w_2\mathbf{v}(w_1) + t_2 * w_2\mathbf{v}(w_2) + t_3 * w_2\mathbf{v}(w_3) + t_4 * w_2\mathbf{v}(w_4) + t_5 * w_2\mathbf{v}(w_5)) / (t_1 + t_2 + t_3 + t_4 + t_5)$$

It can simply be written as

$$\text{TFIDF Weighted Word2Vec (r}_i) = \sum_{i=1}^n (t_i * w_2\mathbf{v}(w_i)) / (\sum_{i=1}^n t_i)$$

Note - Special Case:

If t_i=1 (ie., t₁ = t₂ = t₃ = = 1), then the TF-IDF Weighted Word2Vec is the Average Word2Vec.

Average Word2Vec and TF-IDF Weighted Word2Vec are two simple weighting strategies to convert sentences into vectors. They both serve the same purpose. TF-IDF Weighted Word2Vec weights each word differently as compared to Average Word2Vector. In practice, we try both the options and choose the one that performs better at our task.