

DLDE at 2014 Biomedical Summarization Track

Jie Chen

Beijing Institute of Technology
bit_chenjie@bit.edu.cn

Zhendong Niu

Beijing Institute of Technology
zniu@bit.edu.cn

Abstract

In this note paper, we report our experiment method about task1a and task1b at biomedical summarization track 2014.

1 System Overview

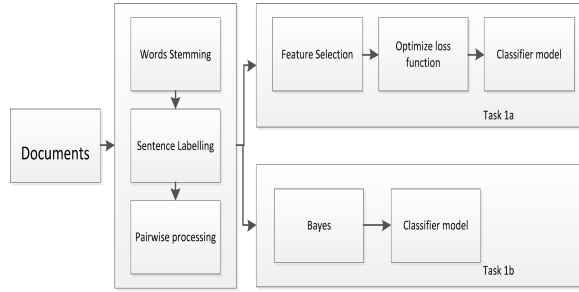


Figure 1: Overview of system framework

Figure 1 outlines the structure of our system framework. It can be divided into two parts. One is prepared processing part which focus on word stemming, sentence labelling and pairwise processing. The other one is model training part which focus on training the classifier model for each task.

2 Training data and processing

Training Data in task1a: Sentences are grouped into different pairs, each pair contains one sentence labeled in the RP and one sentence not labeled from the same RP paper.

Training Data in task1b: Sentences that labeled as the most accurately reflect the citance in RP are selected with their facets to build the bayes classifier training data.

3 Model

Model in task1a: First, We use the LSA to model the latent relationship among sentences in the same paper. Second, the distance between the two sentences is computed. Third, a loss function is introduced to measure the performance.

$$Distance_{s_1, s_2} = \lambda_1 x_{stat} + \lambda_2 cosine(s_1, s_2)$$

where x_{stat} is the statistical feature subset of sentence between s_1, s_2 such as the edit distance, the Jaccard distance.

$$Loss = \sum_i^X \left(\sum_m^M Distance_{c_i, r_m} + \sum_n^N \frac{1}{Distance_{c_i, r_n}} \right)$$

where c_i is the citance sentence in CP, r_m is the labeled sentence in RP paper, r_n is the sentence not labeled in RP paper.

Model in task1b: In our opinion, the latent information about facet hide in some key words so that we choose the native bayes classifier to do this task.

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

where x is the word in the sentences and y is defined by the labeled data.

4 Results

Three runs are submitted for this task, the difference among them is the value selection about λ_1, λ_2 in the distance function.