

Convolution kernels for natural language

Michael Collins
mcollins@research.att.com
 Nigel Duffy
nigeduff@cse.ucse.edu

AT&T
 Labs-Research
 Department of
 Computer Science
 University of
 California

Kernel methods have been used in many algorithms like perceptron support vector machines and PCA. This paper discusses the use of convolution kernel(tree kernel) in Natural Language Processing Task .

NLP tasks involves complex structures like trees, queues, strings, etc. These structures need to be converted into feature vectors. The authors describes how the tree kernel(convolution kernels) can be used to track such complex and high dimensional features, like parse tree in case of parsing task . They also show how the computations are reduced using tree kernel.

A tree can be represented with the help of n-dimensional vector as

$$h(T) = (h_1(T), h_2(T), h_3(T), \dots, h_n(T))$$

Where $h_1(T), h_2(T), \dots, h_n(T)$ is the number of occurrences of i'th tree fragment in the tree T.

Now the tree kernel is defined as the inner dot product between two trees T1 and T2

$$K(T_1, T_2) = h(T_1) \cdot h(T_2).$$

Let the set of nodes in trees T1 and T2 be N1 and N2. Now an indicator function is defined which gives 1 if the subtree is seen rooted at node n and 0 in other case.

$$h_i(T_1) = \sum_{n_1 \in N_1} l_i(n_1) \text{ and } h_i(T_2) = \sum_{n_2 \in N_2} l_i(n_2)$$

$$h(T_1) \cdot h(T_2) = \sum_{n_1 \in N_1} \sum_{n_2 \in N_2} C(n_1, n_2) \text{ where } C(n_1, n_2) = \sum_i l_i(n_1) \cdot l_i(n_2).$$

If the production at n_1 and n_2 are different

$$C(n_1, n_2) = 0$$

If the production at n_1 and n_2 are same and n_1 and n_2 are pre-terminals, then $C(n_1, n_2) = 1$

Else If the production at n_1 and n_2 are same and n_1 and n_2 are not pre-terminals

$$\text{Then } C(n_1, n_2) = \prod_{j=1}^{nc(n_1)} (1 + C(\text{ch}(n_1, j), \text{ch}(n_2, j)))$$

Where $nc(n_1)$ is the number of children of n_1 in the tree. (this is the recursive definition of our kernel function)

Further there were some issues like the value of kernel was depending greatly on size of the tree which was solved by normalising kernel as

$$K'(T_1, T_2) = K(T_1, T_2) / \sqrt{K(T_1, T_1) K(T_2, T_2)}$$

Secondly the tree fragments were scaled according to their relative importance which depends on the size of the tree. This was done by introducing parameter $0 < \lambda \leq 1$.

This change was modified in the recursive definition of the kernel function.

$$C(n_1, n_2) = \lambda \prod_{j=1}^{nc(n_1)} (1 + C(\text{ch}(n_1, j), \text{ch}(n_2, j)))$$

It was found that the time complexity was also reduced since in this case the time complexity does not depend on the number of nodes in a parse tree (which could be exponentially large) but it depends on the kernel function.

To check the utility of the convolution kernel it was applied to the problem of parsing in ATIS corpus. A PCFG was trained on the data and a voted perceptron with the tree kernel and the results were obtained for both cases. The voted perceptron performed much better than the PCFGs with the tree kernel that was applied to it.