

Improvement of Time Complexity and Space on Optimal Binary Search Trees using post dynamic Programming Methodology and Data Preprocessing

S.Hrushikesava Raju^{1*}, M.Nagabhusana Rao²

¹Research Scholar, Regd.No:PP.CSE.0158, Rayalaseema University, Kurnool, A.P.

²Professor, Department of CSE, K L University, Vijayawada, A.P.

hkesavaraju@gmail.com, mnraosir@gmail.com

Available online at: www.ijcseonline.org

Received: Oct/29/2016

Revised: Nov/10/2016

Accepted: Nov/24/2016

Published: Nov/31/2016

Abstract- There are various methods of handling Optimal Binary search trees in order to improve the performance. One of the methods is Dynamic programming which incurs $O(n^3)$ time complexity to store involved computations in a table. The data mining technique called Data Preprocessing is used in order to remove noise early in the dataset and enhances consistency of the given data. The post dynamic computing is applied using knowledge of dynamic programming principle which starts with only required data and computes only the necessary attributes required to construct Optimal Binary Search Tree with time complexity $O(n)$ if there are n identifiers / integers / any complex objects. This approach avoids computing all necessary table attributes. Hence, the complexity or cost of post dynamic computing using Dynamic Programming is proven to be less than $O(n^3)$ or even less than specified in some cases with experimental results.

Keywords— *Optimal Binary Search Tree (OBST), Data Preprocessing, Post computing, Dynamic Programming, Time Complexity*

I. Introduction

Optimal Binary search Tree is a variety of binary trees in which each node stores maximum of two children and it stores strings as identifiers within it or integers or any complex object as nodes of this binary tree. There were many methods such as greedy, recursion, memorizing are useful to do this work. One of those methods is that randomly drawing all possible binary trees and finds a particular binary tree whose cost is low and considers that result is an OBST. The name optimal binary search tree is titled because of simple reason that is finding a key in a tree incurs least number of comparisons. In this method, the optimal binary search tree is chosen in the possible binary trees which involve minimum cost for searching a key in the tree. There are also Greedy, Memorizing, 0/1 knapsack problem, multiple chain multiplication used but they are not efficient because they are all use strategy recursion and another drawback is solutions are not helpful to compute larger problem solution. Dynamic programming is a popular and efficient method to construct binary search tree by following principles such as saving sub problem solutions, reusing sub problem solutions to construct solution for larger problems, and also it avoid using of recursion strategy. Dynamic programming unnecessarily stores all computations in memory and involves $O(n^3)$ complexity and it is applied on any data that is irrelevant or in improper manner. Hence, A data mining technique named data preprocessing is used to sort or clean the given sequence. Also, a technique is proposed titled Data Post computing using Dynamic Programming concept is

performed that only computes required attributes which are required to construct optimal binary search trees. This second step leads to compute some more low level attributes in order to compute that particular attribute.

II. Related Work

There are many methods involved to construct the optimal Binary search trees (OBST). First approach [4,8,9] called randomization which constructs many binary search trees in order to find out an OBST that has minimum cost. The cost can be calculated by multiplying the key and frequency of that node. This approach has a drawback that wastage of time in computing unnecessary trees in addition to correct tree.

Second approach named Computing OBST using sets[6] arranges elements in a tree using recursive approach in which computing optimal substructures and overlapping sub-problems is became expensive(exponential nature) in terms of computing same sub-problems again and again. This issue is sorted in Using Dynamic Programming by a temporary array in a bottom up manner.

Third approach is Greedy recursive approach [7, 8] which initially finds root through best heuristic using recursion and it applies to remaining elements for computing sub trees. This leads sometimes not a OBST although root of each sub tree is optimal.

To overcome these flaws, Dynamic programming[8,9] is introduced which split the given problem into small instances, proposed algorithm applied to all, later integrate them into a large solution for the given problem. The

advantage of dynamic programming is flexibility in integrating the sub problems and each solution takes a space in memory and solution is obtained using recurrence formulas. But still, some refinement is possible on this work that is refinement of DP algorithm is required which reduce the unnecessary computations. This reduces space in a table in storing the elements as well as computations done in less time compared to other techniques. This work can be done in the proposed work where detailed steps are given. The techniques are summarized in the following table.

| Technique | Advantage | Disadvantage |
|---------------------------------------|----------------------------------------------------|--------------------------------------------------------|
| Randomization | Simple, easy and mandatory task | Takes more time in giving right tree with minimum cost |
| Using sets | Optimal sub-structures | Expensive in avoiding overlapping sub-problems |
| Greedy | Guarantee the optimal in each case | Using Recursion cause |
| Traditional Dynamic Programming | Giving a tree optimally | Leads many Unnecessary computations |
| Proposed Dynamic Programming(Post DP) | Gives a tree optimally with little time complexity | NIL |

Table 1: details of each technique used for OBST

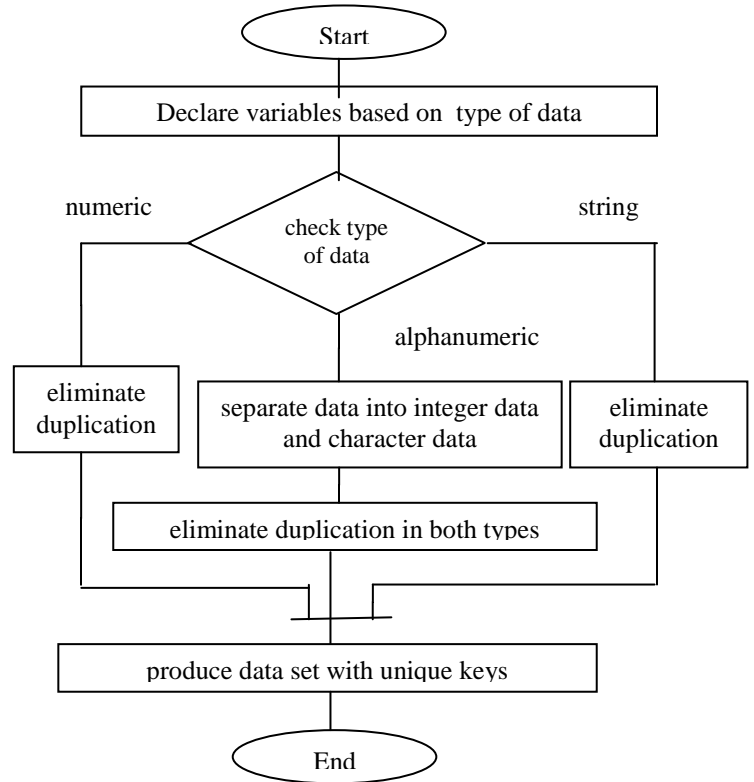
III. Proposed Work

Consider the given data elements are as $a_1, a_2, a_3, \dots, a_{n-1}$, and a_n . These labels denote the keywords in the given data. Suppose the raw data is given, which can be classified into the text keywords, individual characters, and numeric data. For numeric data, it is easy to construct binary search tree by using rules such as left node element should be less than root node and right sub tree element should be greater than root element and this procedure is repeated until last element in the given data is processed.

For alphabets, it is also easy to construct binary search trees by using same logic but left sub tree element is alphabetically comes before the root element, the right sub tree element is alphabetically comes after the root element, and this procedure is applied till last character is processed.

For text keywords to make as an OBST, there are predefined procedures also called algorithms which pick the root node first by a notation t_{ij} which is the last level calculation that helps to pick the root node which can be indexed by the rank value calculated using DP_OBST algorithm. From this notation t_{ij} , root can be picked using the rank r which takes an element in the given data. The sub trees in the next low level are taken based on this rank r . The first level sub trees are $t_{i,r-1}$

(left sub tree) and t_{ij} (right sub tree). The second level nodes can be estimated from $t_{i,r-1}$ and t_{ij} . The rank of entry $(i, r-1=k)$ is $r1$ assume. The sub trees of this are $t_{i,r1-1}$ (left sub tree) and right sub tree t_{r1k} (right sub tree). Assume the rank of first level right sub tree is $r2$. The sub trees of t_{r-ij} are $t_{i,r2-1}$ and t_{r2j} . This is repeated until all elements in the given data are taken a place in the tree.



Flow Graph 1: Data Preprocessing Steps

Based on above data preprocessing flowchart, binary search tree is constructed for numeric, character(alphabet) data separately. Optimal binary search trees is constructed for textual keywords for searching quickly compared to traditional optimal binary search trees which consume more space and take more time for building optimal binary search trees.

The following are the procedures to process the data and construct the OBST from the category of data like numeric, individual characters (alphabets), and text words.

The following is the flowchart which demonstrates the OBST in bottom up fashion from last level to initial level which reduce unnecessary computations.

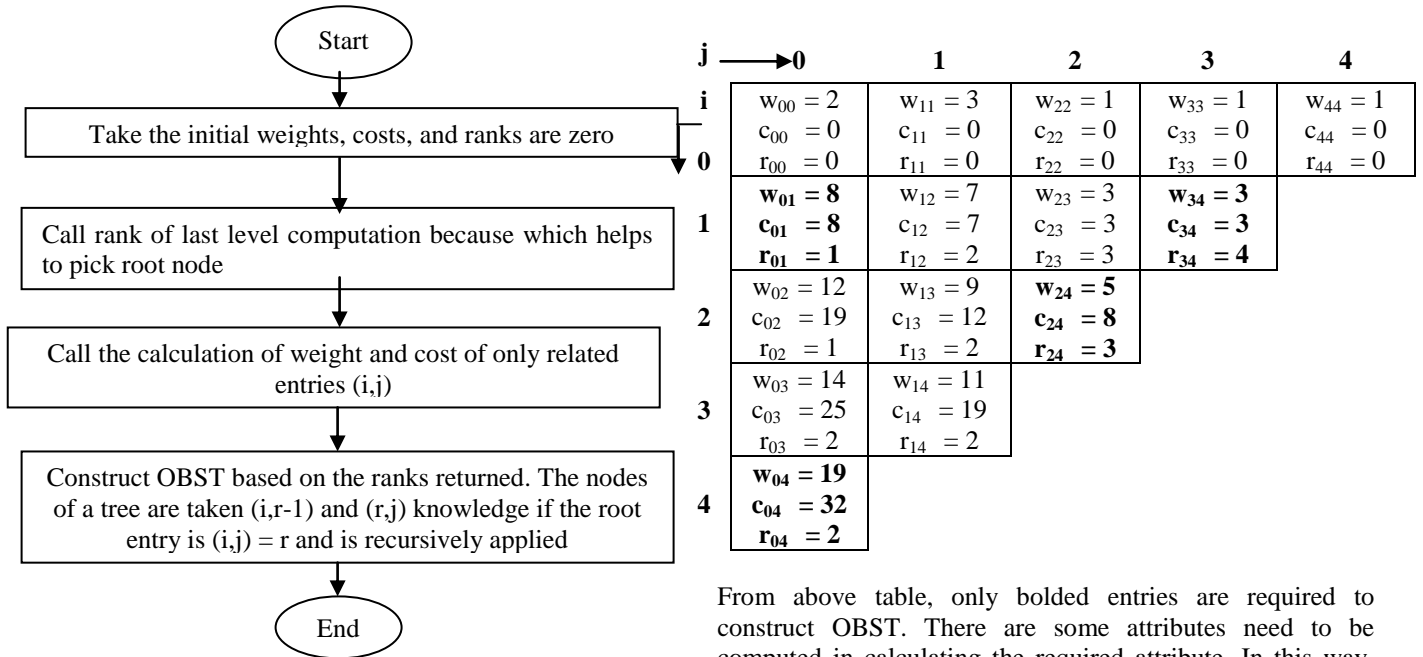


Fig. 1: Procedure for constructing OBST using Dynamic Programming knowledge

IV. Experimental Results

The results are taken by considering some real sample example scenarios.

Example: Take a set (do, if, int, while),

Let p(1: 4) = (3, 3, 1, 1) and q(0: 4) = (2, 3, 1, 1, 1).

The following table shows the comparison between traditional Dynamic Programming and Proposed Dynamic Programming using Data Preprocessing and Post Computing from bottom up manner.

| Technique | Space | Time |
|----------------------------------------------------------------------------------------------|-------------------|---------------------------------------------------------------------------------------------|
| Traditional Dynamic Programming | 45 * sizeof(data) | O(15*3)=45 |
| Proposed Dynamic Programming using Data preprocessing and post computing in bottom up manner | 4 * sizeof(data) | O(4 * 3 + 5 for initial variables + some intermediate variables but not all variables) < 45 |

Table 2: Performance details comparison using traditional and proposed techniques

The following table shows the calculations using traditional Dynamic Programming

From above table, only bolded entries are required to construct OBST. There are some attributes need to be computed in calculating the required attribute. In this way, computing of some unnecessary attributes are going be avoided using Proposed Dynamic Programming using Data Preprocessing and bottom up post computing method.

By observing a simple example, the complexity is more. If you the data set is large, huge number of entries to be computed in a dynamic programming table which not only leads to occupying space and also huge time. To avoid this overhead, proposed technique named Modified Dynamic programming using Data Preprocessing and attribute post computing in bottom up fashion from last level to first level depending on required attributes.

The following is also a graph that shows efficiency of the traditional and proposed Techniques:

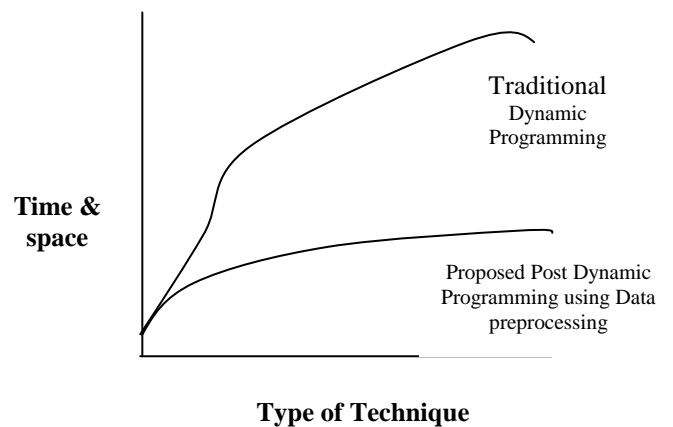


Fig.2: Efficiency of Traditional and Proposed Techniques

V. Conclusion

This approach is very useful early in removing the redundancy in the dataset and later as saving space and time in terms of not to compute many unnecessary attributes using Post Dynamic Programming Approach. The results are described through examples specified in which unnecessary space is removed and also minimized the time to construct optimal binary search tree. This approach may be extended in future using either backtracking or any other advanced designed algorithm to get better time efficiency.

References

- [1]. Micheline Kamber, Hei, "Data Mining: Concepts and Techniques", Second Edition, The Morgan Kaufmann Series, ISBN 13: 978-1-55860-901-3, ISBN 10: 1-55860-901-6.
- [2]. Data Preprocessing Techniques for Data Mining, iasri.res.in/ebook/win_school_aa/notes/Data_Preprocessing.pdf
- [3]. Nguyen Hung Son, Data cleaning and data preprocessing, www.mimuw.edu.pl/~son/datamining/DM/4-preprocess.pdf
- [4]. Andy Mirzaian , DYNAMIC PROGRAMMING: OPTIMAL STATIC BINARY SEARCH TREES, <http://www.cse.yorku.ca/~andy/courses/3101/lecture-notes/OptBST.pdf>
- [5]. ROLF KARLSSON, ALGORITHM THEORY, [HTTP://FILEADMIN.CS.LTH.SE/CS/PERSONAL/ROLF_KARLSSON/LECT5.PDF](http://fileadmin.cs.lth.se/cs/personal/rolf_karlsson/lect5.pdf)
- [6]. DYNAMIC PROGRAMMING | SET 24 (OPTIMAL BINARY SEARCH TREE), [HTTP://WWW.GEEKSFORGEEKS.ORG/DYNAMIC-programming-set-24-optimal-binary-search-tree/](http://www.gEEKSFORGEEKS.ORG/DYNAMIC-programming-set-24-optimal-binary-search-tree/)
- [7]. <http://www.sciencedirect.com/science/article/pii/S0020019081901435>
- [8]. E.Horowitz, S.Sahni, Dinesh Mehta, "Fundamentals of data structures in C++", Second Edition.
- [9]. Ellis Horowitz, Sartaj Sahni, Sanguthevar Rajasekharan, http://vitconference.com/vit_mca/images/resources/DAOA/Fundamentals-of-Computer-Algorithms-By-Ellis-Horowitz-1984.pdf
- [10]. Dynamic Programming- Optimal Binary Search Trees, <http://www.radford.edu/~nokie/classes/360/dp-opt-bst.html>
- [11]. Dynamic Programming, <http://www.cs.utsa.edu/~bylander/cs3343/chapter8handout.pdf>
- [12]. optimal binary search trees, https://en.wikipedia.org/wiki/Optimal_binary_search_tree
- [13]. Data Preprocessing, http://www.cs.ccsu.edu/~markov/ccsu_courses/datamining-3.html
- [14]. Data processing techniques for data mining, http://iasri.res.in/ebook/win_school_aa/notes/Data_Preprocessing.pdf
- [15]. Data Mining: Data and Preprocessing, <http://Staffwww.itn.liu.se/~aidvi/courses/06/dm/lectures/lec2.pdf>
- [16]. "Data Preprocessing in Data Mining" by, ISBN : 9783319102474 & 9783319102467.
- [17]. Salvador Garcia, Julian Luengo, Francisco Hurrera , <http://www.enggjournals.com/ijcse/doc/IJCSE16-08-01-009.pdf>
- [18]. Efficient construction of Optimal Binary Search Trees using Data Preprocessing to improve Quality and Attribute Post

computing to save Space and time through modified Dynamic Programming Technique, <http://www.ijcse.net/issue.php?file=vol05issue1,40-46>.

- [19]. Application of data preprocessing on given data and efficient construction of OBSTs using post dynamic Programming, http://journals.indexcopernicus.com/issue.php?id=737&id_issue=882666

About Authors:



Mr. S. Hrushikeshava Raju, working as a Professor in the Dept. of CSE, SIETK, Narayanavanam Road, Puttur. He is pursuing Ph.D from Rayalaseema University. His areas of interest are Data Mining, Data Structures, and Networks.



Dr. M. Nagabhushana Rao, working as Professor in the Dept. of CSE, K L University, Vijayawada, A.P. He had completed Ph.D from S.V. University in the area of Data mining. He is presently guiding many scholars in various disciplines.