

ADS Assignment

String & text processing

Pattern Matching Algorithms:-

→ Brute force pattern matching algorithms.

In the classic pattern matching problems we are given a text string T of length n and a pattern string P of length m and want to find whether P is a substring of T .

→ we should find the lowest index j within T at which P begins such that $T[j:j+m]$ equals P , or perhaps to find all indices of T at which pattern P begins.

The brute force algorithmic design pattern is a powerful technique for algorithm design when we have something

we wish to search for or when we wish to optimize some function.

PAGE NO.

DATE

we typically enumerate all possible configurations of the inputs involved & pick the best of all these enumerated configurations (Test all possible placements of P relative to T).

def find_brute(T, P):

"return the lowest index of T at which substring P begins (or else -1)"

$n, m = \text{len}(T), \text{len}(P)$ # introduce convenient notations.

for i in range($n - m + 1$) # try every potential starting index i in T
 $k = 0$ # an index into pattern P .

while $k < m$ and $T[i+k] == P[k]$ # k^{th} character of P matches T
 $k += 1$

if $k == m$: # if we reached the end of pattern.
return i # substring $T[i:i+m]$ matches P .

return -1 # failed to find a match starting with any:

- * Boyer-Moore algorithm: It can sometimes avoid comparisons between P and a sizable fraction of the characters in T .
→ The main goal is to improve running time of the brute force algorithm by adding two potential time-saving heuristics.

1) Looking glass heuristic: When testing a possible placement of P against T , begin the comparisons from the end of P and move backward to front of P .

1) Character Jump Heuristic: During the testing of a possible placement of P within $text$, $P[k]$ is handled as follows.
 If c is not contained anywhere in P , then shift P completely past $T[i]$ (for it cannot anywhere in P then shift P completely past $T[i]$). Otherwise, shift P until an occurrence of character c in P gets aligned with $T[i]$.

→ The efficiency of the Boyer-Moore algorithm relies on creating a lookup table that quickly determines where a mismatched character occurs elsewhere in the pattern.

→ If c is in P , $last(c)$ is the index of the last occurrence of c in P . Otherwise we conventionally define $last(c) = -1$.

def find-boyer-moore(T, P):

"Return the lowest index of T at which substring P begins / or else -1 ".

$n, m = len(T), len(P)$ # introduce convenient notations

if $m == 0$ return 0. # trivial case for empty string

$last = \{-1\}$ # build last dictionary.

for k in range(m):

$last[P[k]] = k$ # later occurrence overwrites

align end of pattern at index $m-1$ of text.

$i = m-1$ # an index into T

$k = m-1$ # an index into P

while $i < n$:

if $T[i] == P[k]$: # a matching character.

if $k == 0$:

return i

pattern begins at index i of text.


```

else:
    i -= 1
    k -= 1
else:
    j =
    i +=

```

$k = m - 1$

return -1

examine previous character of
both T and P

PAGE NO.

DATE / /

last CrCT is -1 if not found.

case analysis for jump step

restart at end of pattern.

The Knuth-Morris-Pratt algorithm: Achieves running time of $O(nm)$ which is asymptotically optimal. The main idea of the KMP algorithm is to precompute self overlaps between portions of the pattern so that when a mismatch occurs at one location, we immediately know the maximum amount to shift the pattern before continuing the search.

To implement the KMP algorithm we will precompute a failure function, f , that indicates the proper shift of P upon a failed comparison. The failure function $f(k)$ is defined as the length of the longest prefix of P .

that is a prefix of $P[0:k+1]$.

def find-kmp(T, P):

"Return the lowest index of T at which substring P begins (or else -1)"

$n, m = \text{len}(T), \text{len}(P)$ # introduce convenient notation

if $m = 0$ return 0 # trivial search for empty string

The Huffman Coding Algorithm begins with each of the d distinct characters of the string x to encode being the root node of a single node binary tree. The algorithm proceeds in a series of rounds the algorithm takes the two binary trees.

Algorithm Huffman(x):

Input: String x of length n with d distinct characters

Output: coding tree of x .

compute the frequency $f(c)$ of each character of x . Initialize a priority queue Q .

for each character c in x do.

create a single node binary tree T storing c

insert T into Q with key $f(c)$

while $(n(Q)) > 1$ do.

$(S_1, T_1) = Q.\text{remove-min}()$

$(S_2, T_2) = Q.\text{remove-min}()$

create a new binary tree T with left subtree T_1 and right subtree T_2

insert T into Q with key $f(T_1) + f(T_2)$

$(S, T) = Q.\text{remove-min}()$

return tree T .

greedy method: Huffman's algorithm for building an optimal encoding is example of greedy method. This design pattern is applied to optimization problems,

while we are trying to construct some structure while minimizing or maximizing some property of that structure. The general formula for the greedy method pattern is almost as simple as the brute force method. This approach does not always lead to an optimal solution.

For some problems it works for which passes the greedy-choice property which means that a global optimal condition which can be achieved by a series of locally optimal choices starting from a well defined starting condition.

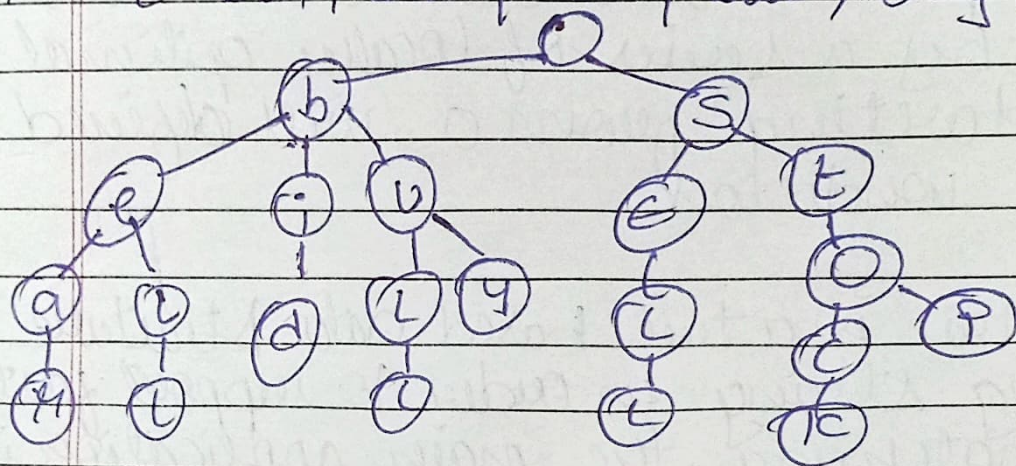
Trees: A tree is a tree based data structure for storing string in order to support fast pattern matching. The main application is information retrieval. The primary query operations that trees support are pattern matching & prefix matching.

Standard trees: Let S be a set of S strings from alphabet Σ such that no string in S is a prefix of another string. A standard tree for S is S an order tree T with the following properties.

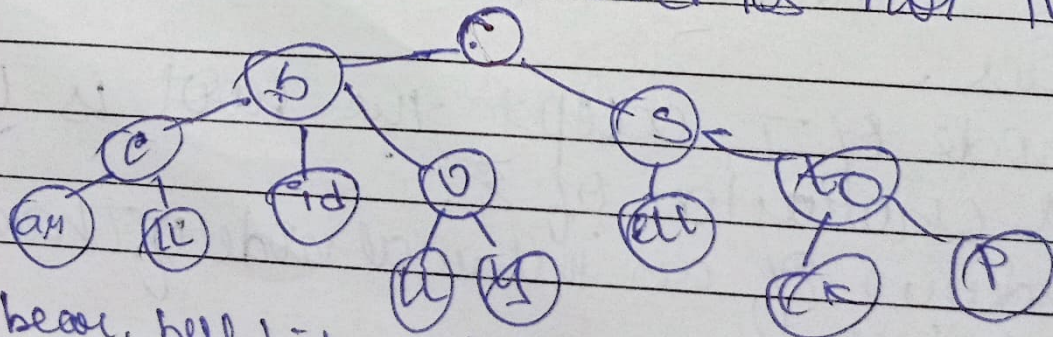
- Each node of T , except the root is labelled with a character of Σ .
- The children of an internal node of T have distinct labels.

\Rightarrow T has s values, each associated with a string S_i such that the concatenation of the labels of the nodes on the path from root to a leaf v of T yields the string of S associated with v . Thus a tree T represents the strings of S with paths from the root to the leaves of T .

Ex: $S = \{ \text{bear, bell, bid, bull, buy, sell, stock, stop} \}$

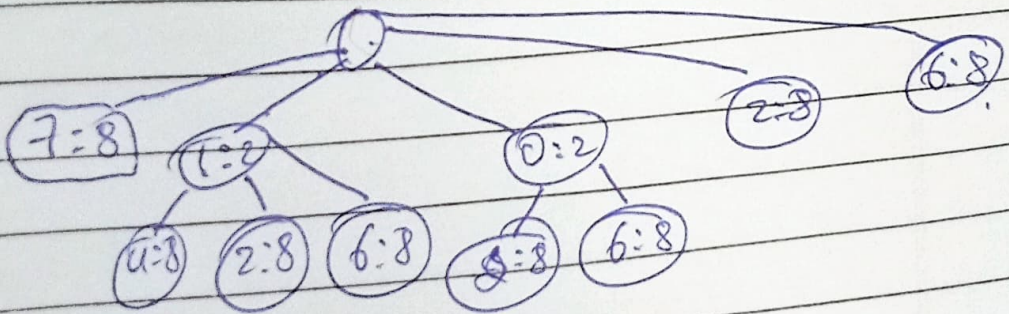
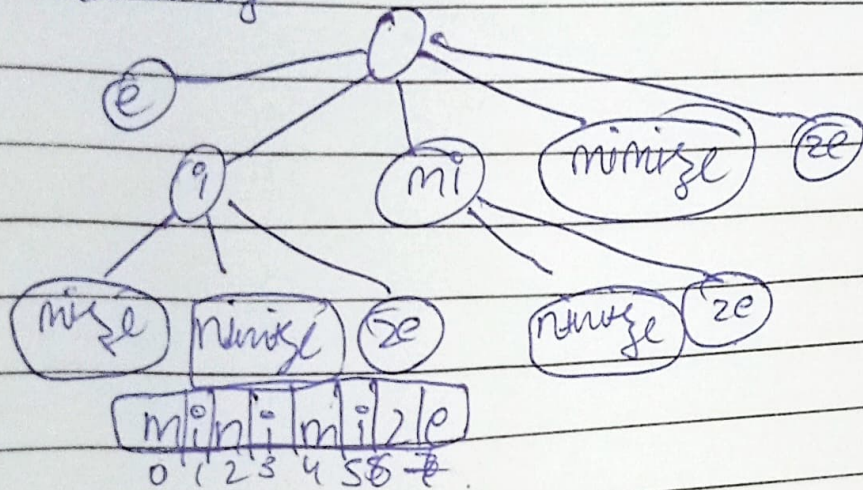


compressed tree is similar to standard tree but it enforces that each internal node in the tree has at least two children. It enforces this rule by compressing chains of single child nodes into individual edges. Let T be a standard tree. We say that an internal node v of T is redundant if it has one child and is not the root.



$S = \{ \text{bear, bell, bid, bull, buy, sell, stock, stop} \}$

Suffix trees : one of the primary applications for this is for the case when the strings in the collection S are all the suffixes of a string x , such a tree is called Suffix tree of string x .



Using a suffix tree allows us to save space over a standard tree by using several space compression techniques.

The suffix tree T for a string x can be used to efficiently perform pattern matching queries on text x .