

## Árboles Generales 6.5 (96 pdf)

árboles - est. de datos

- conjunto de nodos y aristas que conectan pares de nodos. El Futuro se moldea

Características de un árbol con raíz:

- 1) la raíz es 1 nodo
  - 2) todos los nodos excepto la raíz tienen un padre
  - 3) camino único desde raíz hasta cada nodo.
- longitud de camino - cant. de aristas

hojas - nodos sin hijos

hermanos - nodos con el mismo padre

Ej de SO que los usan: DOS, VMS, Unix

la estructura de un fichero se asemeja a un árbol pero no es un árbol en sí

Operaciones: insertar, buscar y eliminar nodos

árboles de expresiones. las hojas se evalúan a sí mismas

# Data Structure and Algorithms in Java



## 13.1 Abundance of Digitalized Text

Al buscar un patrón de texto en un documento se le llama:  
"The pattern machine problem" método → "brute-force method"

↓  
Poco eficiente, fácil de aplicar

Tamaño muy grande de datos de texto → problema de compresión  
se resuelve con

↓  
"greedy method"  
efectivo para este tipo de problema

dynamic programming → técnica de programación usada para  
resolver problemas difíciles de forma rápida y eficiente en  
tiempo polinomial

### 13.3 Tries

- \* Se preprocesa el texto a lugar del patrón
- \* útil si el texto es fijo + muchas búsquedas diferentes



Trie: estructura de datos en forma de árbol → guarda cadenas eficientemente

↳ uso ppal → recuperación de información palabras con cierto prefijo

↳ 2 Operaciones → Búsqueda de palabras completas  
→ Búsqueda de cadenas que comienzan con cierto prefijo

#### 13.3.1 Árbol Estándar ( $\Sigma$ : alfabeto)

\* Conjunto  $S$  de  $s$  cadenas hechas con  $\Sigma$

\* Ninguna cadena debe ser prefijo de otra  $C_j = \{car, cerd, cat\}$   
no funcionaria

Características: → c/nodo excepto la raíz tiene una letra de  $\Sigma$   
→ los hermanos tienen  $\neq$  etiquetas ( $\neq$  letras) - ordenados  $\Sigma$   
→ El hoja representa una cadena de  $S$   
camino raíz → hoja (se obtiene una cadena original)  
→ máximo de nodos  $n+1$ : long. (no  $n$  fal)

\* Si hay cadenas prefijos de otras se debe agregar un carácter especial al final para diferenciarlas

\* Cada nodo puede tener entre 1 hijo y  $|\Sigma|$  hijos



### 13.1.1 Notaciones para cadenas de Strings

Ejemplos:

$s = \text{"CGAAGD"} \rightarrow \text{DNA}$

$t = \text{"http://www.colgla.com"} \rightarrow \text{URL}$  El Futuro se moldea

Un string es un substring de sí mismo.

algoritmos para procesar datos  $\rightarrow$  <sup>se usan</sup> cadenas de caracteres. Formados por símbolos del alfabeto ( $\Sigma$ ), el tamaño puede variar  $\rightarrow$  esto afecta el rendimiento de algoritmos que trabajen con texto.

String  $\rightarrow$  cadenas inmutables

StringBuilder  $\rightarrow$  cadenas mutables

$\text{char}[] \rightarrow S[i]$  es vez de  $S.\text{charAt}(i) \rightarrow$  + simple

Subcadena  $\rightarrow$  parte de cadena  $P[i \dots j]$

si  $i > j \rightarrow$  cadena vacía (longitud 0)

Subcadena  $\rightarrow$  Prefijo: índice empieza en 0  $P[0 \dots j]$   
 $\rightarrow$  Sufijo: termina en el último índice  $P[i \dots n-1]$

"La cadena completa es una subcadena de sí misma"  
\* una subcadena puede ser toda la cadena

\* una subcadena propia es  $\neq$  a la cadena completa

\* la cadena vacía es tanto un prefijo como un sufijo

\* 2 operaciones de los string: modificar cadena propia y devolver información, sin modificarlo.

\* Un trie almacena de forma eficiente los prefijos comunes entre cadenas, desde el inicio hasta donde sea \*

\* cuanto más profundo estamos en el árbol, las probabilidades de que más cadenas compartan el mismo prefijo disminuyen. El Futuro se moldea

\* Un nodo puede tener 0, 1 o más hijos

Altura: longitud de la cadena más larga de sus cadenas

Número de nodos: num máximo es  $n+1$ ,  $n$ : longitud total de todas las cadenas en  $S$

Hojas: tiene  $s$  hojas,  $s$ : número de cadenas del conjunto  $S$

\* Peor caso: ninguna cadena comparte prefijo  $\rightarrow$  Similar a lista enlazada  
c/nodo tiene 1 solo hijo

\* Para buscar seguimos los caracteres de  $S$ , si termina en una hoja  
 $X$  es una cadena de  $S$ , si termina en un nodo interno  $X$  no está en  $S$ .

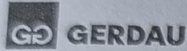
\* Orden de tiempo de ejecución  $O(m \cdot |S|)$   $\rightarrow$  Verifiquemos como máximo

$m+1$  nodos  
c/nodo podemos  $O(|S|)$

\* Si usamos una tabla secundaria de búsqueda o tabla hash  
en c/nodo el tiempo para encontrar a su hijo puede ser  $O(\log |S|)$   
o  $O(1)$  con una tabla tamaño  $|S|$  si  $|S|$  es pequeño como las cadenas ADN

\* longitud  $m$  es  $O(m)$

word matching - verifica si un patrón coincide con una palabra completa o un texto. No coincide con subcadenas



El futuro se moldea

construcción - insertar las cadenas una a una, al insertar una se sigue el camino, si no se encuentra entonces se agregan los nodos necesarios.

inserción:  $O(m \cdot |\Sigma|)$   $\left. \begin{array}{l} m: \text{longitud de la cadena} \\ |\Sigma|: \text{tamaño del alfabeto} \end{array} \right\} \text{con tables Hash} \rightarrow O(m)$

\* construcción completa:  $O(n)$   $n$ : longitud total de las cadenas

Problema de nodos solo c/1 hijo  $\rightarrow$  ineficiente  $\rightarrow$  solución: - Trie comprimido  
- Patricia Trie

2 partes de un word matching usando Trie:

- 1) Texto original  $\rightarrow$  keywords - las que se buscan en el trie
- 2) Trie estándar  $\rightarrow$  c/haya tiene una lista de posiciones donde la palabra comienza en el texto original

Se eliminan los stop words: a, an, the



### 13.3.2 Árboles Comprimidos

Es como un árbol estándar pero se asegura de que c/nodo interno tenga al menos 2 hijos. Se elimina el desperdicio de espacio.

\* cant nodos proporcional a cant de strings

\* Si un nodo (y raíz) tiene 1 solo es redundante

El Futuro se moldea

\* Si hay una cadena de redundantes se comprime en un solo borde

Ej:  $s \rightarrow t \rightarrow o \rightarrow c \rightarrow n$   
→ "Stock"

Este se etiqueta con todos los caracteres concatenados entre separados.

tree estándar → chodo 1 sola letra  
tree comprimido → chodo secuencia de letras

Tamaño de un tree comprimido: cant de cadenas

Propiedades → el nodo interno tiene de 2 a d hijos (d = tamaño del alfabeto)

→ S hijos, 1 x c/palabra

→ número de nodos proporcional a S,  $O(S)$

\* No siempre conviene usarlos → etiquetas + largos = + espacio

\* + útiles como estructuras auxiliares para búsquedas rápidas, no para almacenar todos los caracteres.