

Structures en arbre pour le traitement de textes

Tries

Trie: du mot anglais "retrieval"

Tries 1

Preprocessing Strings

- Un trie est une structure de données pour représenter un ensemble de chaînes de caractères (par exemple, les mots dans un texte)
- Utile pour la recherche rapide
- Aussi pour compresser un texte.

Tries 2

Standard Trie (1)

- Le trie de base est un ensemble de mots dans un arbre tel que:
 - Chaque noeud sauf la racine est un caractère
 - Les enfants sont en ordre alphabétique
 - Les différents chemins de la racine aux feuilles mènent à un mot de S
- Exemple:
S = { bear, bell, bid, bull, buy, sell, stock, stop }

Tries 3

Trie de base (2)

- Un trie utilise $O(n)$ espace et requiert une complexité de $O(dm)$ pour l'insertion, la recherche et le retrait.
 - n nombre de mots dans S
 - m longueur d'un mot
 - d taille de l' alphabet

Tries 4

Trie de base (3)

Pourquoi $O(dm)$ pour la recherche?
 par exemple "ZZZZZZZ"
 A chaque niveau, il faut parcourir tout l'alphabet, donc $7 \cdot 26$ comparaisons!

Tries 5

Recherche de mots avec un Trie

- Chaque feuille enregistre les positions d'occurrences de chaque mot d'un texte

s	e	e	a	b	e	a	r	?	s	e	l	l	s	t	o	c	k	!					
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
s	e	e	a	b	u	l	l	?	b	u	y	s	t	o	c	k	!						
24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	
b	i	d	s	t	o	c	k	!	b	i	d	s	t	o	c	k	!						
47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68		
h	e	a	r	t	h	e	b	e	l	l	?	s	t	o	p	!							
69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88				

Tries 6

Trie compressé

- Dans un trie compressé, chaque noeud doit avoir un degré d'au moins 2
- On élimine les branches redondantes

Tries 7

Représentation compacte

- Chaque noeud stocke le triplet (i, l, k) avec i le mot, l et k sont les indices, e.g. $(6, 1, 2)$ correspond à $S[6]$ de 1 à 2 ou "id".

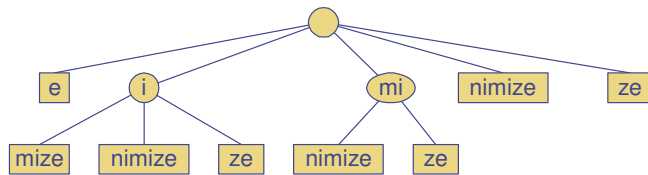
$S[0] =$	0 1 2 3 4	$S[4] =$	0 1 2 3	$S[7] =$	0 1 2 3
	s e e		b u l l l		h e a r
$S[1] =$	b e a r	$S[5] =$	b u y	$S[8] =$	b e l l l
$S[2] =$	s e l l l	$S[6] =$	b i d	$S[9] =$	s t o p
$S[3] =$	s t o c k				

Tries 8

Trie Suffixe (1)

- Un trie suffixe d'un mot X est le trie compressé de tout ses suffixes

m i n i m i z e
0 1 2 3 4 5 6 7



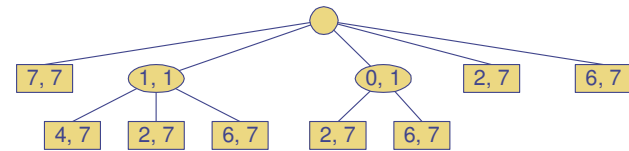
Tries

9

Trie Suffixe (2)

- La représentation compacte du trie suffixe d'un mot X de longueur n à partir d'un alphabet de taille d
 - Utilise un espace $O(n)$
 - Supporte une recherche de motifs avec une complexité $O(dm)$ ou e m est la longueur du motif

m i n i m i z e
0 1 2 3 4 5 6 7



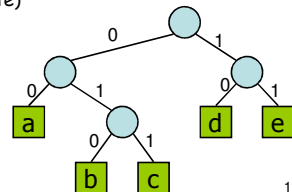
Tries

10

Codage avec Trie (1)

- Un code fait correspondre chaque caractère d'un alphabet avec un mot binaire
- Un code préfixe est un code binaire tel que aucun mot de code est le préfixe d'un autre (par exemple le code Morse)
- Un codage avec Trie produit un code préfixe
 - Chaque feuille est un caractère (un symbole)
 - On doit donc produire un arbre binaire ayant n feuilles
 - Le mot de code associé à un caractère est obtenu en suivant le chemin de la racine jusqu'à ce caractère (0 pour l'enfant de gauche et 1 pour l'enfant de droite)

00	010	011	10	11
a	b	c	d	e

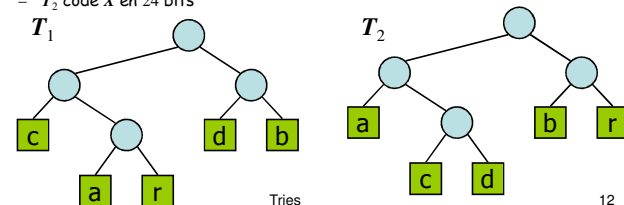


Tries

11

Codage avec Trie (2)

- Soit une chaîne X , nous voulons trouver un code préfixe pour les caractère de X donnant lieu à un codage compact
 - Les caractères les plus fréquents devraient avoir des mots de codes très courts
 - Les caractères rares peuvent avoir des codes plus long
- Exemple
 - $X = \text{abracadabra}$
 - T_1 code X en 29 bits
 - T_2 code X en 24 bits



Tries

12

Algorithme de Huffman

- Soit une chaîne de caractères (ou de symboles) X , l'algorithme de Huffman construit un code préfixe qui **minimise la taille du code pour X**
- Il s'exécute en temps $O(n + d \log d)$, avec n la longueur de X et d le nombre de caractères distinct de X
- Une file à priorité est utilisée comme structure auxiliaire

Tries

13

Algorithme de Huffman

- Insérer chacun des caractères distincts de la chaîne X dans une file à priorité en utilisant la fréquence d'apparition de chaque caractère comme clé; le calcul des fréquences est $O(n)$.
- Retirer 2 noeuds de la file à priorité; $O(\log d)$.
- Combiner ces 2 noeuds et créer un arbre binaire dont la racine est la somme des deux clés; $O(1)$.
- Insérer ce nouvel arbre dans la file à priorité (insérer un arbre correspond à insérer sa racine); $O(\log d)$.
- Répéter jusqu'à ce qu'il y ait un seul arbre dans la file; $O(d)$.

Tries

14

Huffman's Algorithm

Algorithm *HuffmanEncoding*(X)

Input string X of size n

Output optimal encoding trie for X

$C \leftarrow \text{distinctCharacters}(X)$

$\text{computeFrequencies}(C, X)$

$Q \leftarrow$ new empty heap

for all $c \in C$

$T \leftarrow$ new single-node tree storing c

$Q.\text{insert}(\text{getFrequency}(c), T)$

while $Q.\text{size}() > 1$

$f_1 \leftarrow Q.\text{minKey}()$

$T_1 \leftarrow Q.\text{removeMin}()$

$f_2 \leftarrow Q.\text{minKey}()$

$T_2 \leftarrow Q.\text{removeMin}()$

$T \leftarrow \text{join}(T_1, T_2)$

$Q.\text{insert}(f_1 + f_2, T)$

return $Q.\text{removeMin}()$

Tries

15

Exemple

$X = \text{abracadabra}$

Frequencies

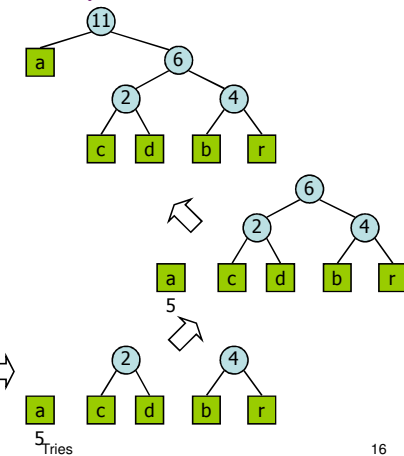
a	b	c	d	r
5	2	1	1	2

a	b	c	d	r
5	2	1	1	2

a	b	c	d	r
5	2	1	1	2

a	b	c	d	r
5	2	1	1	2

a	b	c	d	r
5	2	1	1	2



16

