

A opção mais simples para construção do índice é a seguinte:

d_1 = “a casa é verde”

d_2 = “a casa não é vermelha”

Chave	Valor
a	[< $d_1,1d_2,1$
casa	[< $d_1,1d_2,1$
é	[< $d_1,1d_2,1$
verde	[< $d_1,1$
não	[< $d_2,1$
vermelha	[< $d_2,1$

Neste exemplo, para os documentos d_1 e d_2 , o índice foi armazenado em um dicionário onde a chave é o termo e o valor é uma lista ocorrência representado pelas tuplas <id_doc, freq> onde id_doc é o id do documento e freq é a frequência do termo no documento.

Outra opção é a seguinte, para o mesmo exemplo de documentos:

Chave	Valor
a	<1,0,2>
casa	<2,2,2>
é	<3,4,2>
verde	<4,6,1>
não	<5,7,1>
vermelha	<6,8,1>

<1,d ₁ ,1>	<1,d ₂ ,1>	<2,d ₁ ,1>	<2,d ₂ ,1>	<3,d ₁ ,1>	<3,d ₂ ,1>	<4,d ₁ ,1>	<5,d ₂ ,1>	<6,d ₂ ,1>
0	1	2	3	4	5	6	7	8

Desta vez, há um dicionário onde a chave é o termo e o valor é uma tripla <termo_id,pos_inicial,num_oc>. Onde termo_id é o id do termo, pos_inicial e num_oc são valores que possibilitam saber as posições do vetor de ocorrências que correspondem à este termo: pos_inicial é a posição inicial e num_oc é o número de ocorrências deste termo (ou seja, o número de documentos que o termo ocorre). A lista de ocorrências é representada por <termo_id, doc_id, freq>. Por exemplo, para o termo “a” o valor da tripla é <1,0,2>, isso significa que, no vetor de ocorrências, existe duas ocorrências do termo “a” sendo que a primeira ocorrência está na posição 0 e a segunda na posição 1 indicando que o termo “a” ocorre uma vez no documento d_1 e uma vez no

Caso prefira, é possível armazenar as ocorrências em arquivo. Dessa forma, o dicionário terá uma tupla $\langle \text{termo_id}, \text{offset}, \text{num_oc} \rangle$ onde offset é o deslocamento no arquivo que é necessário para buscar a primeira ocorrência da palavra. Lembrar que, como a ordenação é por termo, o armazenamento do vetor de ocorrências em memória secundária implica na necessidade de ordenação externa (o que dificulta um pouco mais o trabalho):

<u>Chave</u>	Valor
a	$\langle 1, 0, 2 \rangle$
<u>casa</u>	$\langle 2, 6, 2 \rangle$
<u>é</u>	$\langle 3, 12, 2 \rangle$
<u>verde</u>	$\langle 4, 18, 1 \rangle$
<u>não</u>	$\langle 5, 21, 1 \rangle$
<u>vermelha</u>	$\langle 6, 24, 1 \rangle$

offset	<u>Arquivo</u>
0	$\langle 1, d_1, 1 \rangle$
3	$\langle 1, d_2, 1 \rangle$
6	$\langle 2, d_1, 1 \rangle$
9	$\langle 2, d_2, 1 \rangle$
12	$\langle 3, d_1, 1 \rangle$
15	$\langle 3, d_2, 1 \rangle$
18	$\langle 4, d_1, 1 \rangle$
21	$\langle 5, d_2, 1 \rangle$
24	$\langle 6, d_2, 1 \rangle$