

Algorithm Engineering

Md.Atiquur Rahman
Rabin-Karp Algorithm

Rabin Karp Algorithm

The main idea behind Rabin-Karp is to convert strings into numeric hashes so we can compare numbers instead of characters. This makes the process much faster.

We first compute the hash of the pattern, then compare it with the hashes of all substrings of the text. If the hashes match, we do a quick character-by-character check to confirm (to avoid errors due to hash collisions and this step will incur high time complexity if there are enough collisions).

Hash function is simple (for now). Let's say I have string "abc" its hash value is,

$\text{Hash}(\text{"abc"}) = 31^{\{2\}} * 1 + 31^{\{1\}} * 2 + 3$, here 'a' has the value 1 and similarly 'z' has the value 26, 31 is the base and $31^{\{2\}} == 31 * 31$.

For the simplicity, we will ignore all the modular functions **[FOR NOW]**.

HOW THE ALGORITHM WORKS

Let's say I have text as "abcd" and pattern as "bc". Then, the hash values for the pattern is,
 $\text{Hash}(\text{"bc"}) = 31 \cdot 2 + 3$.

Now, should we calculate a single hash value for "abcd"? Obviously, not. We prepare hash values for every available prefix of the "abcd".

$$\text{Hash}[0] = \text{Hash}(\text{"a"}) = 1$$

We don't need to compute Hash of every prefix since we can use previous computed hash.

$$\text{Hash}[1] = \text{Hash}(\text{"ab"}) = \text{Hash}[0] \cdot 31 + 2 = (1) \cdot 31 + 2 = 31 \cdot 1 + 2$$

$$\text{Hash}[2] = \text{Hash}(\text{"abc"}) = \text{Hash}[1] \cdot 31 + 3 = (31 \cdot 1 + 2) \cdot 31 + 3 = 31^{\{2\}} \cdot 1 + 31 \cdot 2 + 3$$

$$\begin{aligned} \text{Hash}[3] &= \text{Hash}(\text{"abcd"}) = \text{Hash}[2] \cdot 31 + 4 = (31^{\{2\}} \cdot 1 + 31 \cdot 2 + 3) \cdot 31 + 4 \\ &= 31^{\{3\}} \cdot 1 + 31^{\{2\}} \cdot 2 + 31 \cdot 3 + 4 \end{aligned}$$

We now compare each substring of the text, equal in length to the pattern, with the pattern itself.

Continued

At first, we match `text[0..1]` with pattern since pattern is of size 2.

$$\text{Hash}[1] = \text{Hash}(\text{text}[0..1]) = \text{Hash}(\text{"ab"}) = 31*1 + 2 \neq \text{Hash}(\text{pattern}) = 31*2 + 3$$

So, no match. Now, we go for, `text[1..2]` and here we need to do some calculation,

$$\text{Hash}[2] = \text{Hash}(\text{text}[0..2]) = \text{Hash}(\text{"abc"}) = 31^{\{2\}}*1 + 31*2 + 3$$

$$\text{Hash}[1-1] = \text{Hash}[0] = \text{Hash}(\text{text}[0..0]) = \text{Hash}(\text{"a"}) = 1$$

$$\text{Hash}[1-1]*31^{\{2-1+1\}} = 31^{\{2\}}*1$$

$$\begin{aligned} \text{Hash}[1..2] &= \text{Hash}[2] - \text{Hash}[1-1]*31^{\{2-1+1\}} = (31^{\{2\}}*1 + 31*2 + 3) - (31^{\{2\}}*1) = 31*2 + 3 = \\ &\text{Hash}(\text{pattern}) \end{aligned}$$

It's a match. Since, `text[1..2] = "bc"` and pattern is "bc" so it should match.

ANOTHER EXAMPLE

My student id is “01234567” and the pattern is always the last two digits of the student id, in my case it is “67”. Now, ‘0’ has the value 1 and ‘9’ has the value 10. Since, I already know there’s a match, hence let’s directly calculate the hash values and try to see if it’s a match.

$\text{Hash}[L..R] = \text{Hash}[R] - \text{Hash}[L-1] * \text{base}^{\{R-L+1\}}$, why do we need this equation ? because, we need to remove the hash value for the $0..L-1$ prefix. But, if we just subtract the $\text{Hash}[L-1]$ it will not work since, that prefix is already multiplied by base for $R-L+1$ times to reach $\text{Hash}[R]$.

$$\begin{aligned}\text{Hash}[6..7] &= \text{Hash}[7] - \text{Hash}[5] * 31^{\{2\}} \\ &= (31^{\{7\}} * 1 + 31^{\{6\}} * 2 + 31^{\{5\}} * 3 + 31^{\{4\}} * 4 + 31^{\{3\}} * 5 + 31^{\{2\}} * 6 + 31 * 7 + 8) - \\ &\quad (31^{\{5\}} * 1 + 31^{\{4\}} * 2 + 31^{\{3\}} * 3 + 31^{\{2\}} * 4 + 31 * 5 + 6) * 31^{\{2\}}\end{aligned}$$

Since after 5 we have more 2 characters at position 6 and 7, all the values of the previous characters will be multiplied by 31 for 2 times.

$$\begin{aligned}&= (31^{\{7\}} * 1 + 31^{\{6\}} * 2 + 31^{\{5\}} * 3 + 31^{\{4\}} * 4 + 31^{\{3\}} * 5 + 31^{\{2\}} * 6 + 31 * 7 + 8) - \\ &\quad (31^{\{7\}} * 1 + 31^{\{6\}} * 2 + 31^{\{5\}} * 3 + 31^{\{4\}} * 4 + 31^{\{3\}} * 5 + 31^{\{2\}} * 6) = 31 * 7 + 8, \text{ it's a match.}\end{aligned}$$

Class Work with actual SID

My actual student id is “180041123”, and the pattern will be “23”. Now I need to calculate Hash[6..7] and Hash(“23”), to show both the values are same.

Write the same in your notebook by hand, scan the page(s), convert them to a PDF, and submit it in the classroom to receive today's attendance and classwork marks.