

מבנה נתונים ואלגוריתמים - הרצאה 18

1 בינואר 2012

טבלאות גיבוב

פונק' גיבוב $i \rightarrow f(x)$ כאשר $x \in V$, $V \in \{0, \dots, N\}$ מרחב רב מידי ו $\alpha < 0.5$. רוצים שההתפלגות של הפונק' תהיה אחידה. אז, בהינתן x בודקים את המיקום $f(x)$, ויש שתי אפשרויות:

1. כל מקום מכיל מצביע לרשימה ונינו להכניס מספר ערכים לאותו מקום בטבלה i . (טבלת גיבוב פתוחה, Open Addressing).

2. כל מקום מכיל את הנתונים עצמו, ואם המקום תפוס שעשים Rehash (טבלת גיבוב סגורה, Closed Addressing).

מגדירים α להיות היחס בין כמות הנתונים לבין גודל הטבלה. הזמן הטיפוסי לפעולה בטבלה פתוחה הוא $\alpha + 1$ ובטבלה סגורה $\frac{1}{1-\alpha}$.

עלויות אתחול

צורך לאתחל את המערך שלנו בהתחלתו. נניח שאתחלו אותו בגודל k , אלו רוצים שיתקיים $\alpha < 0.5$, אך כאשר $\alpha > 0.5$ נדרש להגדיל את המערך כדי להחזיר את α לערך קטן. נגיד שנכפיל את המערך לגודל $2k$, אך צריך להעביר את הערכים כי פונק' הגיבוב משתנה. דרך אחת לעשות זאת היא להחזיק את פונק', החדש והישנה, ואז שמחפשים ערך ממחפשים אותו בשני המיקומות, אחד של הפונק' החדש והאחד של הישנה. אפשרות שנייה, היא לעבור על כל המערך ולהעביר כל מקום למקום החדש. זה אمنם לוקח יותר זמן, אך זה בדר"כ מה שעושים כי לאחר נctrיך להחזיק הרבה פונק'.

בדר"כ בטבלאות גיבוב עושים גם את ההפק - מחלקים את המערך פי 2 אם α קינה מספיק. יכולם אם $\alpha > \alpha_0$ מגדילים את המערך פי 2 ואם $\alpha_0 < \alpha$ מקטינים אותו פי 2. נשים לב שמתקיים $\frac{\alpha_0}{2} \ll 1$, כדי ששינויים קטנים בכמות הנתונים לא יגרמו לנו להכפיל ולהלך הרבה.

Rehash - שיטה ל-Cuckoo Hashing

מציאותים 2 פונק' גיבוב, f_1 ו- f_2 . אם אחד מהם לא תפוס, מכניסים שם את x . אם שניהם קרושרים, מכניסים את x באחד מהם ואת מה שהוציאנו מנסים להכניס למקום אחר, לפי הפונק' השנייה, וכך הלאה.

השווות רצפים

נניח שיש לנו טקסט T באורך L ויש לנו מילה M באורך K . אנו רוצים לבדוק האם ואייפה נמצא המילה בטקסט. האלגוריתם הנאיבי זה לעבור אותן אות בטקסט ולבדוק האם היא מתאימה למילה, אך זו הולות היא $O(L \cdot K)$.

שיטת Rabin Karp

1. נגדיר פונק' גיבוב.

2. חשב $f(M)$.

3. עבור בכל מקום בטקסט ונחשב $z_i = f(T(i : i + k - 1))$.

.4. אם $i + +$ אז השווה את המילים אותן אותן. אחרת, $f(M) = z_i$.

שם הפונק' הוא לוקחת $O(L \cdot K)$ בניתוח.

Rolling Hash

פונק' גיבוב היא מתגלגת אם ניתן לחשב אותה באמצעות הפונק' במילה מסוימת אחת לאחר. למשל, $f(x) = \sum \text{Ascii}(x_j)$ היא פונק' גיבוב מתגלגת (אך היא לא טובה כי היא לא מתפלגת טוב). פונק' אחרת, יותר טובה, היא $f(x) = \sum \text{Ascii}(x_j) \cdot A^j$ כאשר A מס' ראשוני גדול. במקרה, כדי "לגלגל", צריך להוריד את האות הראשונה של המילה הקודמת, לחלק הכל ב- A ולהוסיף את האות الأخيرة כפול A^{k-1} .

$$f(x_{i+1}) = \frac{f(x_i) - \text{Ascii}(0)}{A} + \text{Ascii}(i+k) \cdot A^{k-1}$$

עם פונק' גיבוב כזו, אלגוריתם Rabin-Karp פועל בעלות של $O(L)$.

אלגוריתם Knuth-Morris-Pratt

אלגוריתם 1 אלגוריתם Knuth-Morris-Pratt

```
j=0
i = 0
while (i < L-k)
    while j<k && i<L-k:
        if T(i + j)=M(j):
            j++
            if (j==k):
                output("victory!")
            endif
        else
            i += j-C(j)
            j = max(C(j), 0)
        endif
    endwhile
    j=C(k)
    i=i+k-C(k)
endwhile
```

כאשר נגידר את $C(j)$ כאורך הסיפה הגדולה ביותר שהיא גם רישא של המילה $M(0 : j - 1)$ נגידר $C(0) = -1$ בצורה מלאכותית כדי שיסתדר עם האלגוריתם.

$C(1) = 0$ אם $M(j - 1) = M(C(j - 1) + 1)$ או $C(j) = C(j - 1) + 1$. אחרת, $C(j) = 0$ זה לא בדיק נכון, יש דוגמאות שזה לא נכון, נראה בשיעור הבא.

Algorithm Boyer-Moore

זה אלגוריתם שעבוד ברוב המקרים ב- $O(\frac{L}{k})$. כתוב אותו מסודר בשיעור הבא.