

Refining core string edits and alignment

مقدمه

در بیوانفورماتیک و یا علم زیست شناسی هم تراز کردن توالی به روشهای مرتب کردن توالی های DNA و RNA و پروتئین گفته می شود به طوریکه مکانهای مشابه بین توالی ها را مشخص نماید. این مکانهای مشابه میتواند نشانگر ارتباط عملکردی، ساختاری و یا تکاملی بین توالی ها باشد. توالی های هم تراز شده را معمولا به شکل سطرهایی زیر هم درون یک ماتریس نشان می دهیم. در صورت لزوم در برخی مکانهای توالی فاصله ایجاد اضافه می کنیم تا در چند ستون پی در پی کاراکترهای یکسان زیر هم قرار گیرند.

رویکردهای محاسباتی برای هم تراز کردن توالی ها: 1- تراز بندی سراسری (Global Alignment)

2- تراز بندی محلی (Local Alignment)

تراز بندی سراسری

محاسبه ی هم تراز ی سراسری، شکلی از بهینه سازی سراسری است که به هم تراز ی فشار می آورد تا در کل طول توالی های مورد جستجو گسترده شود. در واقع تلاش می کند هر پس ماندی را در هر توالی هم تراز کند که این برای وقتی بهترین کاربرد را دارد که اندازه ی دورشته (توالی) یکسان باشد.

تراز بندی محلی

هم تراز کردن محلی، نواحی مشابه درون توالی های بلند را که معمولا در طول توالی بسیار متفاوتند، تشخیص می دهد. در واقع هم تراز کردن محلی که نشان دهنده ی وضعیت فاصله هاست به نحوی که توالی ها به اندازه ی کافی مشابه باشند.

مثال

Global FTFTALILLAVAV

F--TAL-LLA-AV

Local FTFTALILL-AVAV

--FTAL-LLAAV--

حال به مساله خود می پردازیم.

Needleman_wunsch الگوریتم

مرحله 1- تعریف مساله

الگوریتم Needleman_wunsch یک تراز بندی سراسری روی دو رشته s1 و s2 می باشد که

این الگوریتم در تراز بندی DNA ها و پروتئین ها در علم زیست شناسی رایج است.

این الگوریتم اولین الگوریتم کاربردی در برنامه ریزی پویا برای مقایسه ی توالی هاست. (در علم زیست شناسی)

ماتریس مشابهت (Similarity matrix)

زمانی که می خواهیم میزان مشابهت دو کاراکتر را بررسی کنیم به ماتریس S مراجعه می کنیم که این میزان مشابهت در این ماتریس مقدار دهی شده است.

مثلاً $s(a,b)$ مقدار مشابهت کاراکترهای a,b است.

	A	G	C	T	-
A	10	-1	-3	-4	-5
G		7	-5	-3	-5
C			9	0	-5
T				8	-5
-					-5

:

مثلاً برای رشته های تراز بندی شده ی روبرو:

AGACTAGTTAC

CGA - - -GACGT

$$S(A,C)+S(G,G)+S(A,A)+S(C,-)+S(T,-)+S(A,-) \\ +S(G,G)+S(T,A)+S(T,C)+S(A,G)+S(C,T)=$$

$$-3+7+10-(3*5)+7-4+0-1+0=1$$

برای یافتن ماکزیمم میزان مشابهت در تراز بندی ما نیاز به یک ماتریس V داریم که هر درایه در سطر i ام و ستون j ام این ماتریس با $V(i,j)$ نشان داده می شود و از طرفی سطرهای ماتریس بیانگر رشته ی s_1 و ستونهای آن بیانگر کاراکترهای رشته ی s_2 میباشد.

که در این صورت می دانیم با توجه به اینکه رشته ی s_1 دارای n کاراکتر و رشته ی s_2 دارای m کاراکتر است میزان فضای حافظه $O(mn)$ است.

مرحله 2- راه حل بازگشتی

همانطور که قبلا در رایه ی edit distance ذکر شده

$$V(0,j) = \sum_{1 \leq k \leq j} s(-, s_2(k))$$

$$V(i,0) = \sum_{1 \leq k \leq i} s(s_1(k), -)$$

و رابطه بازگشتی

$$V(i,j) = \max[v(i-1,j-1) + s(s_1(i), s_2(j)),$$

$$v(i-1,j) + s(s_1(i), -),$$

$$v(i,j-1) + s(-, s_2(j))]$$

مرحله 3- محاسبه بیشترین مقدار مشابهت ها $[V_{mn}]$ تراز بندی بهینه

شبه کد برای محاسبه ی ماتریس v :

1. For $i=0$ to $\text{length}(s_1)$
2. $V(i,0) \leftarrow v(i-1,0) + s(-, s_1(i))$
3. For $j=0$ to $\text{length}(s_2)$
4. $V(0,j) \leftarrow v(0,j-1) + s(s_2(j), -)$
5. For $i=1$ to $\text{length}(s_1)$
6. For $j=1$ to $\text{length}(s_2)$
7. (
8. Match $\leftarrow v(i-1,j-1) + s(s_1(i), s_2(j))$
9. Delete $\leftarrow v(i-1,j) + s(s_1(i), -)$

10. $\text{Insert} \leftarrow v(i,j-1)+s(-,s2(j))$
11. $V(i,j) \leftarrow \max(\text{Match}, \text{Delete}, \text{Insert})$
12.)

بامحاسبه ی ماتریس V_{mn} می توان ماکزیم مقدار راز ترازبندی های ممکن بدست آورد.
 برای محاسبه ی تراز بندی از آخرین خانه سمت راست V_{mn} شروع به حساب کردن می کنیم
 ومقدارها را باسه مقدار ممکن که ازسه عمل (Match, Delete, Insert) بدست می
 آید،مقایسه می کنیم.

مرحله 4- ساختن ترازبندی بهینه

شبه کد برای محاسبه ی $s1, s2$ ترازبندی شده با استفاده از ماتریس
 V (محاسبه مسیر بهینه)

1. Alignment $s1 \leftarrow \text{“ “}$
2. Alignment $s2 \leftarrow \text{“ “}$
3. $i \leftarrow \text{length}(s1)$
4. $j \leftarrow \text{length}(s2)$
5. while($i>0$ and $j>0$)
6. {
7. Score $\leftarrow v(i,j)$
8. Score diag $\leftarrow v(i-1,j-1)$
9. Score up $\leftarrow v(i,j-1)$
10. Score left $\leftarrow v(i-1,j)$
- 11.If(score == score diag + $s(s1(i),s2(j))$)
- 12.{
- 13.Alignment $s1 \leftarrow s1(i)+\text{alignment } s1$
- 14.Alignment $s2 \leftarrow s2(j)+\text{alignment } s2$

```

15.  i ← i-1
16.  j ← j-1
17.}
18.Else if (score == score left+ s(s1(i),_))
19.{
20. Alignment s1 ← s1(i) + alignment s1
21. Alignment s2 ← “_” + alignment s2
22.  i ← i-1
23.}
24.Other wise (score == score up+ s(_,s2(j)))
25.{
26.Alignment s1 ← “_” + alignment s1
27.Alignment s2 ← s2(j)+alignment s2
28.  j ← j-1
29.  }
30.  }
31.While (i>0)
32.{
33. Alignment s1 ← s1(i)+alignment s1
34. Alignment s2 ← “_” + alignment s2
35.  i ← i-1
36.}
37. While (j>0)
38.{
39.Alignment s1 ← “_” + alignment s1
40.Alignment s2 ← s2(j)+alignment s2
41.j ← j-1
42.}

```

مرحله 5- بهبود برنامه

همانطور که مشاهده شد الگوریتم بالا به $O(mn)$ فضا و همچنین به $O(mn)$ زمان نیاز دارد، برای بهبود برنامه می خواهیم فضای $O(mn)$ را به $O(\min(m,n))$ کاهش دهیم. یعنی فضا را به یک فضای خطی کاهش دهیم.

الگوریتم هیشبرگ فضا را به $O(\min(m,n))$ کاهش میدهد ولی زمان را تقریباً به دو برابر افزایش می دهد.

حال مساله ی جدیدی به وجود آمده است که با استفاده از برنامه ریزی پویا به صورت زیر ارائه می شود.

مرحله 1- تعریف مساله

الگوریتم هیشبرگ می تواند به عنوان یک نسخه ی تقسیم و غلبه

(divide and conquer) از الگوریتم needleman_wunsch معرفی شود.

الگوریتم هیشبرگ برای محاسبات زیستی رایج است تا بهترین تراز بندی سراسری بین توالی های DNA و پروتئین ها را محاسبه کند.

دورشته ی زیر موجود می باشد:

S1 _____ $n - i$ n

S2 _____ $m - j$ m

تعریف: برای هر رشته s ، وجود دارد s^r که به معکوس رشته اشاره می کند.

مثال:

$$S=abc \Rightarrow s^r=cba$$

در نتیجه داریم:

$$s1^r \dots \dots \dots i \text{-----} n$$

$$s2^r \dots \dots \dots j \text{-----} m$$

مرحله 2- راه حل بازگشتی

حال برای حل این مساله نیاز به شکستن مساله به زیر مساله هاست.

برای این کار ابتدا مساله را به دوزیر مساله تبدیل می کنیم، بدین صورت که یکی از رشته ها مثلا رشته s_1 اول را نصف میکنیم، پس در ادامه s_2 را به دنبال یک مقدار بهینه برای شکستن رشته s_1 دوم می باشیم (K)

$$S1 \text{-----} n/2 \text{-----} n$$

$$S2 \text{-----} k \text{-----} m$$

بنابراین k بهینه مقدار V بهینه یعنی v ماکزیمم را به ما می دهد.

برای بدست آوردن مقادیر v ، قسمت اول مساله را می توان از جدول v محاسبه کرد ولی در مورد تکه s_2 دوم چون حدپایین آنرا نداریم نمی توان محاسبه نمود. به همین دلیل از معکوس آن استفاده میکنیم. یعنی V^r

بنابراین از فرمول بازگشتی زیر استفاده می کنیم:

$$V(m,n)=\max [v(n/2,k) + v^r(n/2,m-k)]$$

$$0 \leq k \leq m$$

پس مساله به دو زیرساختار بهینه v, v^r تبدیل شد که با بهترین مقدار k که همان k^* است بدست می آید.

اثبات بهینه بودن زیرمساله ها:

اگر k^* در تکه ی اول مساله یعنی قسمت v باشد $\{s1 \dots n/2\}$

$$s2 \dots k^*$$

آنگاه زیر ساختار اول که بهینه است در مورد زیر ساختار دوم هم اگر فرض کنیم که یک k بهتری در آن موجود می باشد که مقدار v را ماکزیم می کند، در نتیجه نهایتاً یک $v(m,n)$ بدست می آید که مقدار بیشتری از $v \max$ دارد که با حکم ما در تناقض استدر نتیجه فرض خلف باطل است.

در مورد قسمت دوم اثبات هم دقیقاً به همین ترتیب، با این تفاوت که k بهینه را در تکه ی دوم مساله فرض می کنیم و...

سپس مساله را به همین ترتیب میشکنیم، یعنی در هر بار رشته ی اول را به دو قسمت تقسیم کرده ($\dots \rightarrow n/8 \rightarrow n/4 \rightarrow n/2$) تا به ابتدای رشته برسیم، سپس از نقطه ی $(0,0)$ تا (n,m) در هر مرحله k بهینه را بدست می آوریم و از این طریق مسیر بهینه ی ما ساخته می شود.

مشخص است که با این روش دیگر نیازی به نگه داشتن ماتریس v به طور کامل نیست و ماتنها به خود درایه وسط قبل از آن احتیاج داریم وبقیه ی سطرهای قبل از آن دورریز می شود (یعنی به جز خود درایه ی مورد نظر به درایه ی قبلی آن در ستون قبلی ودرایه ی قطری قبلی ودرایه ی بالای آن در سطر قبلی نیاز است، بنابراین علاوه بر خود سطر تنها به سطر قبلی نیاز است) در نتیجه فضا به فضای خطی کاهش می یابد.

.....

.....
مرحله ی 3 و 4- محاسبه و ساختن مسیر بهینه

$V(s1,s2)$

$n=length(s1)$, $m=length(s2)$

If $n \leq 1$ or $m \leq 1$

Output alignment $(s1,s2)$ using Needleman_wunsch

Else

$Mid=floor(n/2)$

$S1_left = s1(1,mid)$

$S1_right = s1(n,n-mid)$

$Max=-\infty$

For $k=1$ to $k=m$ {

$S2_left=s2(1,k)$

$S2_right = s2(m, m-k)$

$F = v(s1_left, s2_left) + v(s1_right, s2_right)$

If $F > \max$

$\text{Max} = F$

$k^* = \text{index max} \{ v(s1_left, s2_left) + v(s1_right, s2_right) \}$

Return max

تحلیل زمان:

در تحلیل فضایی گفته شد چون بعضی از اطلاعات در هر مرحله دورریز می شود فضای ما به فضای خطی کاهش می یابد، ولی باید توجه داشت چون همه ی محاسبات باید انجام شود و این بار دو ماتریس V, V^T باید محاسبه شود پس زمان افزایش می یابد:

تعداد زیر مساله ها :

$n/2 \quad 2m$

$n/2^2 \quad 2m$

.

$$\rightarrow 2m \log_2^n$$

•
•

$$n/2^i \quad 2m$$

بنابراین:

$$n/2 \times 2m + n/2^2 \times 2m + \dots + n/2^i \times 2m$$

$$o(2m \log_2^n)$$

$$o(m)$$

$$\Rightarrow \Rightarrow o(m^2 \log_2^n)$$

تحلیل زمانی که در کتاب ذکر شده:

$$cnm \rightarrow cnm/2 \rightarrow cnm/4 \rightarrow \dots$$

در مرحله ی i ام از بازگشت $cnm/2^{i-1} \Rightarrow$

از طرفی تعداد فراخوانی ها $\log_2 n =$

$$\sum_{i=1}^{\log_2 n} cnm/2^{i-1} = cnm \sum_{i=1}^{\log_2 n} 1/2^{i-1} \quad (*)$$

تصاعد هندسی:

$$a(1-q^n)/1-q = (1-1/2^{\log_2 n})/1/2 = 1-2^{\log_2 n-1}/1/2 = 1-2^{\log_2 n-1} \cdot 2 = 1-2^{\log_2 n} = 1-n$$

$$1/n/1/2 = 2(n-1)/n$$

$\Rightarrow \Rightarrow$

$$(*) = cnm \cdot 2(n-1)/n = 2cnm - 2m \leq 2cnm \Rightarrow o(nm)$$