

الیمنت Sequence Alignment

فصل سی و یکم از سری کتب الکترونیکی رایگان
پرتال بیوانفورماتیک ایرانیان www.ibp.ir

درج کننده مطلب: بابک باباعباسی

منتج: کارگاه الکترونیکی بیوانفورماتیک، انجمن بیوتکنولوژی به مدیریت دکتر ملابویی

مقایسه دو توالی

در دهه ۸۰، یک محقق هیچ برنامه رایانه‌ای برای این که بتواند بین تعدادی توالی مشابه توالی را به توالی خود پیدا کند نداشت. بنابراین زبده ترین دانشمندان آن روزگار نیز مجبور بودند این کار را به صورت دستی انجام دهند. به طور مثال، اگر قرار بود از بین چهار توالی زیر مشابه ترین توالی را به توالی خود انتخاب می کرد. باید تک تک توالی ها را با توالی الگو مقایسه می کرد و میزان شباهت‌ها را در هر مقایسه به دست می آورد. امروزه به این هم‌ردیفی دوگانه (Pairwise Alignment) می‌گویند.

AATTGGCCTTGC

AATGCGCCGTGC

AAGCCGAAGTGA

TTTTGCCCAAGG

ساده ترین راه برای مقایسه کردن دو توالی این است که هر بار دو توالی را در زیر هم قرار دهیم و یک به یک بازها را با هم مقایسه کنیم تا شبیه ترین توالی را پیدا کنیم (شکل‌های A تا D زیر). ولی سوالی که پیش می‌آید این است که چگونه و با چه معیاری دو توالی مشابه تر را انتخاب می‌کنیم؟ در این جاست که بحث امتیاز دهی (Scoring) مطرح می‌شود. به عنوان مثال، ساده ترین نوع امتیاز دهی این گونه می‌تواند باشد که اگر دو بازی که زیر هم قرار می‌گیرند از یکسان باشند، امتیاز ۱ و اگر همسان نباشند، امتیاز صفر داده شود. با این روش می‌توان مشابه ترین توالی و درجه شباهت سایر توالی‌ها را با توالی الگو به دست آورد.

حال توالی‌های بالا را با توالی الگوی داده شده با همین روش هم‌ردیفی دو گانه می‌کنیم. به نظر می‌رسد مورد C که امتیاز ۹ گرفته است، همسانی بیشتری با توالی ما دارد.

A

```

AATTGGCCTTGC
TTTTGCCCAAGG
001110110010 =6
    
```

B

```

AATTGGCCTTGC
AAGCCGAAGTGA
110001000100 =4
    
```

C

```

AATTGGCCTTGC
AATGGCCCGTGC
111010110111 =9
    
```

D

```

AATTGGCCTTGC
AATTGGGCTTGC
111111011100 =7
    
```

در هم‌ردیفی دوگانه مسأله این است که دو توالی چه قدر به هم شبیه هستند. زمانی که ما برای تعیین میزان همسانی از امتیاز دهی و عدد استفاده می‌کنیم، در واقع از روش‌های ریاضی برای حل مسأله‌ی زیستی استفاده می‌کنیم. از آن جا که در دنیای زیست‌شناسی پارامترهای دخیل بسیار زیاد و در بسیاری از موارد ناشناخته هستند، بنابراین برای حل مسأله‌های زیستی با استفاده از الگوریتم‌های ریاضی و رایانه‌ای، همواره با مشکل عدم تطبیق کامل مدل ریاضی با واقعیت زیستی روبرو هستیم. تفاوت راه‌حل‌ها و الگوریتم‌ها با هم در این است که جواب کدام یک به واقعیت زیستی که مشاهده می‌شود، نزدیک‌تر است و آن را بهتر توجیه می‌کند.

اما همانطور که در شکل زیر دیده می‌شود، این دو توالی را به طریق دیگری هم می‌توان هم‌ردیف کرد. لذا سوالاتی هنوز باقی است مانند آن که آیا هم‌ردیفی دیگری ممکن است؟ کدام هم‌ردیفی بهتر است؟ هم‌ردیفی بهتر یعنی چه؟ کدام هم‌ردیفی گویای اتفاقات زیستی است؟ آیا در هم‌ردیفی‌ها روندهای تکاملی قابل ردیابی است؟ تا چه حد؟ و چگونه می‌توان هم‌ردیفی‌ها در این راستا به کار گرفت؟

اینها سوالات عمیقی است که پایه‌های اساسی داده‌پردازی زیستی را تشکیل می‌دهند. لکن در این دوره سعی بر آن است که اصول همردیفی تا اندازه‌ای آموزش داده شود تا به توان از آن برای جستجوی توالی‌های مشابه و قضاوت در مورد میزان مشابهت و درک مفهوم خانواده‌های ژنی و پروتئینی استفاده نمود .

AATTGGCCTTGC
AATGGCCCGTGC
111010110111 =9

AATTGGCCT-TGC
AATGGCCCGTGC
101111100111 =9

روش Dot plot

به عنوان راهی برای شناسایی تمامی همردیفی‌های ممکن محققین، روشی گرافیکی به نام دات پلات (dot plot) به کار بردند. در این روش دو توالی به صورت عمود بر هم روی محور X ها و Y ها در یک صفحه قرار داده می شوند و در هر نقطه ای که شبیه هم باشند عدد یک قرار داده می شود. اگر دو توالی کاملاً شبیه باشند

	A	A	T	T	G	G	C	C	T	T	G	C
A	1	1	0	0	0	0	0	0	0	0	0	0
A	1	1	0	0	0	0	0	0	0	0	0	0
T	0	0	1	1	0	0	0	0	1	1	0	0
G	0	0	0	0	1	1	0	0	0	0	1	0
G	0	0	0	0	1	1	0	0	0	0	1	0
C	0	0	0	0	0	0	1	1	0	0	0	1
C	0	0	0	0	0	0	1	1	0	0	0	1
C	0	0	0	0	0	0	1	1	0	0	0	1
G	0	0	0	0	1	1	0	0	0	0	1	0
T	0	0	1	1	0	0	0	0	1	1	0	0
G	0	0	0	0	1	1	0	0	0	0	1	0
C	0	0	0	0	0	0	1	1	0	0	0	1

در نهایت، از رسم نقاط یک خط اوریج بدون شکستگی را می توان از انتهای سمت راست صفحه به انتهای سمت برآست بالای صفحه رسم کرد. همردیفی در واقع مشخص کردن رابطه ی بین نوکلئوتید های یک توالی با توالی دیگر است. اگر دو توالی در مثال فوق را به صورت دات پلات در آوریم جدول زیر به دست خواهد آمد. اگر دور بیش از دو عدد ۱ به دنبال هم خط بکشیم، منظره بالا ظاهر خواهد شد. حال بر اساس این خطوط اریب می‌توان کلیه همردیفی‌های دو گانه را استخراج کرده و به صورت خطی نوشت. در این جدول، دو رابطه خطی C و E در شکل قبل با خطوط آبی و قرمز نشان داده شده‌اند .

مقایسه همردیفی‌ها

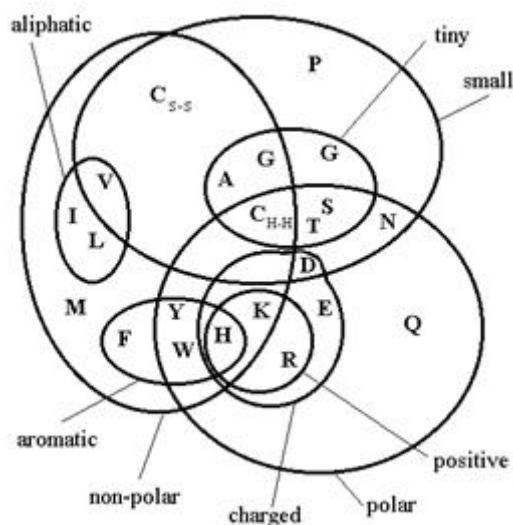
در مثال‌های فوق روش امتیاز دهی صفر ویک را می توان یک نوع الگوریتم به حساب آورد که بر اساس آن همردیفی با بالاترین امتیاز را به عنوان بهترین همردیفی انتخاب کردیم. ولی در عمل می‌توان همردیفی‌هایی را مثال زد که با وجود امتیاز مساوی یا حتی بالاتر صحیح نبوده و با دانسته‌های قبلی تطبیق نمی‌کنند. بنابراین تلاش زیادی برای طراحی و بکارگیری الگوریتم‌هایی دقیق‌تر صورت می‌گیرد که تا هر چه بیشتر دربرگیرنده واقعیات زیستی و اصول حاکم بر حیات باشد .

به طور مثال، در همردیفی توالی‌های نوکلئوتیدی می‌توان بین امتیازات جایگزینی‌های از نوع جانشینی (Substitution) و انتقال (Transition) تفاوت قایل شد. زیرا با توجه به

ساختمان دو رشته‌ای DNA احتمال جایگزینی بازهای پورینی با هم و باز های پیریمیدینی با هم بیشتر است. در حالی که الگوریتم قبلی تفاوتی را بین این دو حالت قائل نبود. بنابراین جواب هایی هم که با الگوریتم قبلی به دست آمد کمتر به واقعیت نزدیک است .

این مشکل در توالی‌های پروتئینی به طور جدی‌تری مطرح است. در این توالی‌های همه جایگزینی‌ها اشکال ساختاری و عملکردی زیادی ایجاد نمی‌کنند. به عبارت دیگر، برخی اسید آمینه‌ها خواص فیزیکوشیمیایی مشابهی دارند (شکل زیر) و می‌توانند با حداقل تغییر خواص جایگزین یکدیگر شوند .

نکته ی دیگر این است که در هم‌ردیفی‌های بالا فرض برابر بودن طول توالی‌هاست. در حالی که در تکامل توالی‌ها هم پدیده ی اضافه شدن را داریم و هم پدیده حذف اتفاق می‌افتد . بنابراین این انتظار می‌رود در بسیاری از موارد دو توالی را با هم مقایسه کنیم که دارای طول یکسانی نباشند .



جمع‌بندی این مقدمات نشان می‌دهد می‌توان با کمی‌نمودن (امتیازدهی) نتایج هم‌ردیفی آنها را باهم مقایسه نمود. البته برای کمی نمودن هم‌ردیفی‌ها حداقل دو نوع امتیاز دهی را بایستی منظور کرد .

امتیاز دهی جایگزینی‌ها امتیاز دهی حذف و اضافه شدن توالی‌ها

به این ترتیب، امتیاز هر هم‌ردیفی جمع جبری کلیه امتیازات جایگزینی‌ها و حذف یا اضافه‌ها خواهد بود.

با درک این که روش صفر و یک کفایت نمی‌کند و جایگزینی نوکلئوتیدها یا اسید آمینه‌ها با یکدیگر امتیاز منفی یا مثبت مساوی ندارند، متخصصین امر در پی تهیه جداول امتیازدهی جایگزینی‌ها (Substitution Scoring Matrices) بوده‌اند. به طوری که تا حد امکان واقعیت‌های زیستی را منعکس نماید .

برای توالی‌های نوکلئوتیدی کار چندان دشوار نیست زیرا هر گونه جایگزینی منجر به جهش می‌شود که اثر آن در رمزدهی پروتئین‌ها ممکن است مشاهده شود. یعنی در این مولکول‌های بحث ساختار و عمل چندان مطرح نیست. البته با توجه به ساختمان دو رشته‌ای DNA متخصصین تکامل زیستی بین جانشینی نوکلئوتید پورین و پیریمیدین و انتقال از پورین به پیریمیدین یا بالعکس تفاوت قائلند. جدول زیر نمونه‌ای از جداول امتیازدهی برای هم‌ردیفی دو توالی نوکلئوتیدی را نشان می‌دهد. در این جدول به طور ساده‌ای کلیه همسانی‌ها امتیاز +5 و برای غیر جفت شدگی امتیاز -4 در نظر گرفته شده است. سایر حروف در صورت وجود انتخاب برای دو نوکلئوتید یا بیشتر در هر موقعیت از توالی کاربرد دارند .

Matrix Structure: Nucleotides

	A	T	G	C	S	M	R	Y	K	H	B	V	D	N
A	5	-4	-4	-4	-4	1	1	-4	-4	1	-4	-1	-1	-2
T	-4	5	-4	-4	-4	1	-4	1	1	-4	-1	-4	-1	-2
G	-4	-4	5	-4	1	-4	1	-4	1	-4	-1	-1	-4	-2
C	-4	-4	-4	5	1	-4	-4	1	-4	1	-1	-1	-4	-2
S	-4	1	1	1	5	-4	-4	-4	-4	-4	1	1	-4	-2
M	1	1	-4	-4	-4	5	-4	-4	-4	-4	1	-4	-1	-2
R	1	-4	1	-4	-4	-4	5	-4	-4	-4	1	-4	-1	-2
Y	-4	1	-4	1	-4	-4	-4	5	-4	-4	1	-4	-1	-2
K	-4	1	1	-4	-4	-4	-4	-4	5	-4	1	-4	-1	-2
H	1	-4	-4	1	-4	-4	-4	-4	-4	5	1	-4	-1	-2
B	-4	-1	-1	-1	-1	-4	-4	-4	-4	-4	5	1	-4	-2
V	-1	-4	-1	-1	-1	-4	-4	-4	-4	-4	-4	5	1	-2
D	-1	-1	-4	-1	-1	-4	-1	-1	-1	-1	-4	-4	5	-1
N	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	5

- Simple match/mismatch scoring scheme:

Match	+ 5
Mismatch	- 4
- Assumes each nucleotide occurs 25% of the time

در تهیه جداول امتیازات جایگزینی توالی‌های پروتئینی خواص فیزیکوشیمیایی اسیدهای آمینه و تاثیر جایگزینی آنها در ساختار و عمل پروتئین‌ها مطرح است. در گذشته، پژوهشگران به گروه‌بندی اسیدهای آمینه بر اساس خواص آنها (شکل صفحه قبل) مراجعه کرده و میزان مشابهت را به صورت توصیفی (و نه عددی) بیان می‌کردند. در دو دهه اخیر روش‌های تهیه جداول امتیازدهی جایگزینی مبتنی بر داده‌های موجود در طبیعت بوده است. با فرض بر این که اگر دو اسید آمینه دارای خواص فیزیکوشیمیایی مشابهی هستند بایستی در طول تکامل جایگزینی آنها تحمل شده باشد، پژوهشگران نسبت به جمع‌آوری توالی‌ها، هم‌ردفی آنها با هم و محاسبه فراوانی جایگزینی‌ها در بین پروتئین‌های هم‌خانواده اقدام نمودند .

در اولین تلاش، با هم‌ردیفی ۱۵۷۲ توالی پروتئینی در ۷۱ درخت از ۳۴ خانواده پروتئینی گروه‌بندی شدند. سپس فراوانی جایگزینی یک اسید آمینه با اسید آمینه دیگر در فرمول زیر بکار گرفته شد :

$$PAM_n = \frac{\log \text{Probability of one substitution}}{\log \text{Probability of occurring by chance}} * 100$$

در این فرمول یک واحد (Point Accepted Mutation) PAM معادل تغییر ۱ در یک توالی صدتایی از اسید آمینه‌هاست. داده‌های حاصل در جدول PAM ثبت می‌شود. از آنجا که در طول تکامل ممکن است اسید آمینه در یک موقعیت چندین بار جایگزین شود، جدول حاصل را چندین بار در خود ضرب می‌کنند. به طور مثال، برای تهیه جدول PAM₁₀₀ را ۱۰۰ بار در خودش ضرب می‌کنند (جدول a در صفحه بعد).

بعدها جداول دیگری تدوین شدند که از یک نوع اصول پیروی می‌کردند. با این تفاوت که فراوانی جایگزینی‌ها تنها در مناطق حفاظت شده (Conserved Blocks) برای ساختن جدول وارد محاسبه می‌شدند. در آن هنگام، ۲۰۰۰ بلوک از ۵۰۰ خانواده پروتئینی در نظر گرفته شد. به طور مشابهی، فرمول زیر بکار رفت :

$$BLOSUM_{\%} = \frac{\log \text{Probability of substitution in block}}{\log \text{Probability of occurring by chance}}$$

این جداول را BLOSUM نامیدن که از اصطلاح Block Substitution Matrix برگرفته شده است. شماره جدول به نوع بلوک مورد استفاده برای محاسبه فراوانی و احتمال وقوع جایگزینی بستگی دارد. مثلاً BLOSUM62 یعنی این جدول بر مبنای فراوانی جایگزینی‌ها در بلوک‌های حاوی توالی‌های با همسانی ۶۲ درصد یا بیشتر تشکیل شده است (شکل زیر).

جند اصطلاح:

همسانی با (Identity) وقتی بکار می‌رود که منظور همردیفی یک نوکلئوتید یا اسیدآمینه با همانند آن منظور باشد.

مشابهت با (Similarity) وقتی بکار می‌رود که منظور همردیفی یک نوکلئوتید یا اسیدآمینه با معادل آن منظور باشد (شکل دو صفحه قبل و جداول صفحه بعد).

همولوژی (Homology) وقتی بکار می‌رود که منظور همردیفی یک نوکلئوتید یا اسیدآمینه با معادل آن منظور باشد و رابطه نیاکاتی بین توالی‌ها در نظر باشد.

(a)

G	A	V	L	I	P	S	T	D	E	N	Q	K	R	H	F	Y	M	C	B	Z	X	*
G	5																					G
A	1	2																				A
V	-1	0	4																			V
L	-4	-2	2	6																		L
I	-3	-1	4	2	5																	I
P	0	1	-1	-3	-2	6																P
S	1	1	-1	-3	-1	1	2															S
T	0	1	0	-2	0	0	1	3														T
D	1	0	-2	-4	-2	-1	0	0	4													D
E	0	0	-2	-3	-2	-1	0	0	3	4												E
N	0	0	-2	-3	-2	0	1	0	2	1	2											N
Q	-1	0	-2	-2	-2	0	-1	-1	2	2	1	4										Q
K	-2	-1	-2	-3	-2	-1	0	0	0	0	1	1	5									K
R	-3	-2	-2	-3	-2	0	0	-1	-1	-1	0	1	3	6								R
H	-2	-1	-2	-2	-2	0	-1	-1	1	1	2	3	0	2	6							H
F	-5	-3	-1	2	1	-5	-3	-3	-6	-5	-3	-5	-5	-4	-2	9						F
Y	-5	-3	-2	-1	-1	-5	-3	-3	-4	-4	-2	-4	-4	-4	0	7	10					Y
W	-7	-6	-6	-2	-5	-6	-2	-5	-7	-7	-4	-5	-3	2	-3	0	0	17				W
M	-3	-1	2	4	2	-2	-2	-1	-3	-2	-2	-1	0	0	-2	0	-2	-4	6			M
C	-3	-2	-2	-6	-2	-3	0	-2	-5	-5	-4	-5	-5	-4	-3	-4	0	-8	-5	12		C
B	0	0	-2	-3	-2	-1	0	0	3	3	2	1	1	-1	1	-4	-3	-5	-2	-4	3	B
Z	0	0	-2	-3	-2	0	0	-1	3	3	1	3	0	0	2	-5	-4	-6	-2	-5	2	Z
X	-1	0	-1	-1	-1	0	0	-1	-1	0	-1	-1	-1	-1	-2	-2	-4	-1	-3	-1	-1	X
*	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	-8	1*
G	A	V	L	I	P	S	T	D	E	N	Q	K	R	H	F	Y	M	C	B	Z	X	*

(b)

G	A	V	L	I	P	S	T	D	E	N	Q	K	R	H	F	Y	M	C	B	Z	X	*
G	6																					G
A	0	4																				A
V	-3	0	4																			V
L	-4	-1	1	4																		L
I	-4	-1	3	2	4																	I
P	-2	-1	-2	-3	-3	7																P
S	0	1	-2	-2	-2	-1	4															S
T	-2	0	0	-1	-1	-1	1	5														T
D	-1	-2	-3	-4	-3	-1	0	-1	6													D
E	-2	-1	-2	-3	-3	-1	0	-1	2	5												E
N	0	-2	-3	-3	-3	-2	1	0	1	0	6											N
Q	-2	-1	-2	-2	-3	-1	0	-1	0	2	0	5										Q
K	-2	-1	-2	-2	-3	-1	0	-1	-1	1	0	1	5									K
R	-2	-1	-3	-2	-3	-2	-1	-1	-2	0	0	1	2	5								R
H	-2	-2	-3	-3	-3	-2	-1	-2	-1	0	1	0	-1	0	8							H
F	-3	-2	-1	0	0	-4	-2	-2	-3	-3	-3	-3	-3	-3	-1	6						F
Y	-3	-2	-1	-1	-1	-3	-2	-2	-3	-2	-2	-1	-2	-2	2	3	7					Y
W	-2	-3	-3	-2	-3	-4	-3	-2	-4	-3	-4	-2	-3	-3	-2	1	2	11				W
M	-3	-1	1	2	1	-2	-1	-1	-3	-2	-2	0	-1	-1	-2	0	-1	-1	5			M
C	-3	0	-1	-1	-3	-1	-1	-3	-4	-3	-3	-3	-3	-3	-2	-2	-1	9				C
B	-1	-2	-3	-4	-3	-2	0	-1	4	1	3	0	0	-1	0	-3	-3	-4	-3	-3	4	B
Z	-2	-1	-2	-3	-3	-1	0	-1	1	4	0	3	1	0	0	-3	-2	-3	-1	-3	1	Z
X	-1	0	-1	-1	-1	-2	0	0	-1	-1	-1	-1	-1	-1	-1	-1	-2	-1	-2	-1	-1	X
*	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	-4	1*
G	A	V	L	I	P	S	T	D	E	N	Q	K	R	H	F	Y	M	C	B	Z	X	*

منابع زیادی در تهیه این کتاب الکترونیک استفاده شده است که در بخش معرفی کتب پرتال بیوانفورماتیک ایرانیان میتوانید مطالعه کنید www.ibp.ir

منبع اصلی این مطلب کارگاه بیوانفورماتیک به مدیریت دکتر ملبوبی میباشد

دوست ارجمند این مطلب آموزشی به منظور آشنایی هر چه بیشتر ایرانیان با علم بیوانفورماتیک تهیه شده است خواهشمند است از هر گونه سو استفاده از این مطالب پرهیزید زیرا که این مطالب تحت حمایت قانون کپی رایت میباشند

این مطالب به این دلیل به صورت کتابهای الکترونیکی تهیه شده اند تا

هرکس، در هرجا، در هر زمان و به صورت رایگان بتواند به این مطالب دست رسی داشته باشد پس اگر چنانچه تمایل به استفاده از این مطالب در سایت یا وبلاگ خود را دارید این کار با ذکر منبع بلامانع میباشد.

هیچ شخص حقیقی یا حقوقی حق ندارد تا از مطالب کتب الکترونیکی این پرتال به منظور چاپ و نشر کتاب استفاده کند چاپ و نشر کتاب کاری مقدس میباشد ولی فروش علم جزء ناشایسته ترین کارهاست که متأسفانه برخی با استفاده از نا آشنایی مردم با منبع سرشار اینترنت دست به چاپ کتب میزنند و آن را به قیمت بسیار بالایی به فروش میرسانند حرف من این است اگر عاشق علم هستید دانش خود را به صورت رایگان در اختیار همگان قرار دهید و بهترین راه آن توسط اینترنت میباشد.

به امید روزی که هیچ کتابی چاپ نشود و تمام مطالب به صورت رایگان در اینترنت قابل دست رسی باشند تا حد اقل در بخش استفاده از منابع علمی نامی از غنی و فقیر برده نشود .

پرتال بیوانفورماتیک ایرانیان www.ibp.ir
بابک باباعباسی