



جستجو در پایگاه داده

مسعود صدیقین

۱۲ بهمن ۱۳۹۳

چکیده

حجم داده‌ها در پایگاه‌های داده‌های زیستی به سرعت در حال افزایش است. همچنین تعداد پرسش‌ها نیز با سرعت زیادی رو به افزایش می‌باشد. برای مثال حدود ۴۰۰۰۰ پرسش در یک روز انجام می‌شود. بنابراین نیاز به روش‌های کارایی است تا بتوان در پایگاه‌های داده جستجو انجام داد. در این جلسه به بررسی این روش‌ها می‌پردازیم.

۱ تعریف مساله

مساله‌ی جستجو در پایگاه داده به طور کلی به این شکل می‌باشد:

تعریف ۱ جستجو در پایگاه داده: یک پایگاه داده D از رشته‌های ژنوم و یا پروتئینی در اختیار داریم. با داشتن یک پرسش مانند Q می‌خواهیم رشته‌ای مانند S را درون D بیابیم که به Q نزدیک‌تر است.

می‌توان برای نزدیک بودن، دو معنی در نظر گرفت:

• End-space free

S

Q

• Local alignment

S

Q

همچنین برای مقایسه‌ی دو الگوریتم جستجو می‌توانیم معیارهای زیر را در نظر بگیریم:

- Sensitivity: به این معنی که الگوریتم می‌تواند جواب درست را پیدا کند.
 - Specificity: به این معنی که الگوریتم کمتر دچار خطا شده و جواب نادرست را نشان می‌دهد.
- در ابتدای کار تنها به مساله تطابق محلی^۱ خواهیم پرداخت. روش‌های موجود برای این کار به شرح زیر است:

۱. جستجویی تمام و کمال: که می‌توان برای جستجو در هر رشته از الگوریتم Smith-Waterman استفاده نمود.

۲. روش‌های مکاشفه‌ای^۲: روش‌های مکاشفه‌ای متعددی برای این کار وجود دارد که از جمله‌ی آن می‌توان به روش‌های *FASTA*, *BLAST/BLAT*, *PatternHunter* اشاره نمود.

۳. روش‌های فیلتر کردن و بهبود دادن: مانند روش‌های *LSH*, *QUASAR* اشاره کرد.

۴. روش *BWT – SW*

از میان روش‌های فوق، الگوریتم *Smith – Waterman* حساس‌ترین الگوریتم بوده و البته نیاز به محاسبات $O(|Q||D|)$ دارد که با توجه به حجم پایگاه داده و تعداد پرسش‌ها، عملی نمی‌باشد. بنابراین به سراغ روش‌های مکاشفه‌ای می‌رویم.

۲ الگوریتم‌های *FASTA* و *BLAST*

۱.۲ تاریخچه

در جدول زیر، تاریخچه مختصری از این دو الگوریتم ارائه شده است:

Year	Algorithm	Author	-
1985	FASTP	Limpman and Pearson	Global gapped alignment
1988	FASTA	Pearson and Limpman	Local gapped alignment
1990	BLAST1	Altschul et al	Local ungapped alignment
1996	WU-BLAST2	Gish	Local gapped alignment
1997	NCBI-BLAST2	Altschul et al	Local gapped alignment

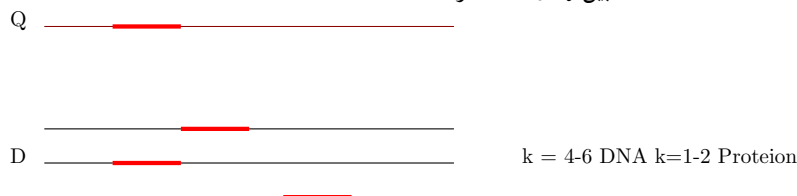
۲.۲ روش *FASTA*

این روش بر پایه آن بنا شده است که اگر بتوان دو رشته را با یکدیگر هم‌ردیف نمود و هم‌ردیفی دو رشته امتیاز بالایی دارد، آنگاه زیررشته‌های مشترکی بایستی بین آن‌ها موجود باشد. این الگوریتم دارای ۵ مرحله می‌باشد:

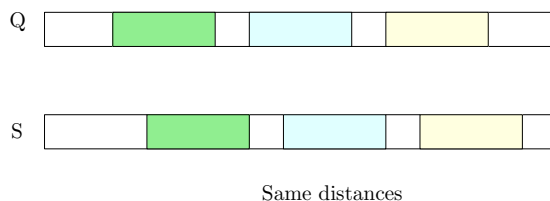
^۱ Local Alignment

^۲ Heuristic

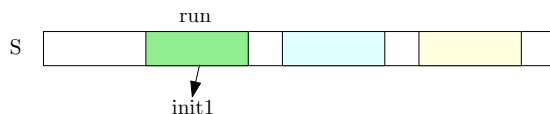
۱. در مرحله‌ی اول، ابتدا *Hot - spot* ها را پیدا می‌کنیم. یک *Hot - spot* در حقیقت یک *k - mer* است که بین *Q* و *D* مشترک است.



۲. در این مرحله بهترین *run* را پیدا می‌کنیم. به ازای هر رشته در *D*، ۱۰ تا انتخاب می‌گردد. برای امتیاز هی، به *Hot - spot* ها امتیاز مثبت داده می‌شود و به فاصله میان *Hot - spot* ها امتیاز منفی که مقدار آن با طول این فاصله، افزایش می‌یابد.



۳. مقدار امتیاز واقعی که بر حسب ماتریس امتیازدهی می‌باشد را برای مکان‌های به دست آمده محاسبه می‌کنیم. بهترین آن‌ها امتیاز *init1* را دارد. این امتیازدهی تنها بر پایه ماتریس امتیازدهی بوده و بنابراین *gap* را در نظر نمی‌گیرد. (فاصله یکی می‌باشند)



۴. برای هر رشته در *D*، آن محل‌هایی که دارای امتیازی پایین‌تر از یک مقدار مخصوص می‌باشند را بیرون می‌ریزیم. برای بقیه آن‌ها سعی می‌کنیم که با اجازه دادن *indel*، آن‌ها را هم‌ردیف کنیم. بهترین امتیاز در بین ترکیب‌های مختلف از طریق برنامه نویسی پویا قابل محاسبه است.

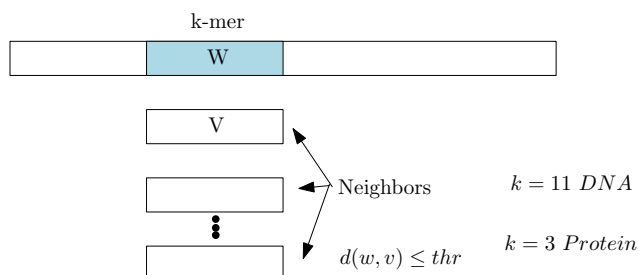
۵. رشته‌هایی که *initn* آن‌ها کمتر از یک مقدار مشخصی هست را بیرون می‌ریزیم و برای سایرین، الگوریتم *Smith - Waterman banded* را اجرا می‌کنیم و امتیاز بهینه را به دست می‌آوریم. رشته‌های باقیمانده را بر حسب *Opt* ردیف می‌نماییم.

به طور کلی این روش از روش *Smith - Waterman* سریعتر است. اما دارای حساسیت کمتری می‌باشد.

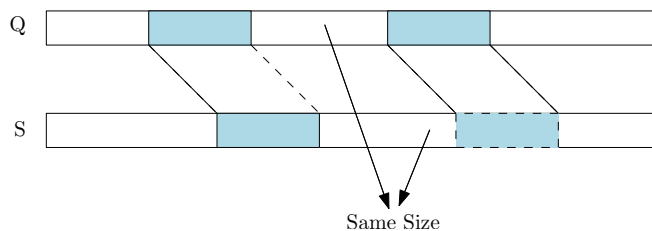
۳.۲ روش *BLAST*

این روش مخفف Basic Local Alignment Search Tool است. این روش نسبت به روش *FASTA* دارای سرعت بیشتری است. اما میزان حساسیت آن کمتر است. همانطور که در بخش تاریخچه به آن اشاره شد، الگوریتم *BLAST1* در سال ۱۹۹۰ ارائه شد که در آن شباهت‌های محلی بدون گپ محاسبه می‌شود. همچنین در الگوریتم *BLAST2* که در سال ۹۶ و ۹۷ ارائه شد، گپ نیز در نظر گرفته می‌شود. این الگوریتم نیز شامل ۴ مرحله می‌باشد:

۱. پردازش پرسش: برای هر مکان از پرسش، تمام k -mer هایی که دارای امتیازی بالاتر از یک مقدار آستانه می باشند را می یابیم. به این ها همسایه های رشته می گوئیم:



۲. پایگاه داده را اسکن کرده و Hit هایی که با رشته های همسایه به دست می آیند را پیدا می کنیم. هر Hit شامل دقیقاً یک مکان در رشته پرسش و یک مکان در پایگاه داده می باشد.
۳. برای آنکه یک Hit را بسط دهیم، بایستی یک Hit دیگر در فاصله ی کمتر از A از آن موجود باشد، که در آن $A = 40$ است. این Hit بایستی $Diagonal$ باشد.



۴. این Hit ها را انقدر گسترش بدون گپ می دهیم تا امتیاز نواحی اضافه شده از یک مقدار خاص مانند X کمتر گردد. از بین آن هایی که گسترش داده شده اند، آن هایی را انتخاب می کنیم که دارای مجموع امتیاز بالاتر یا مساوی S داشته باشند. به این نواحی HSP ^۳ می گوئیم.
۵. برای HSP های به دست آمده، گسترش با گپ ^۴ را انجام می دهیم. البته به شرطی که از یک میزان آستانه بیشتر باشند.

در حقیقت این روش، یک روش بهبود یافته از الگوریتم $Smith - Waterman$ است که در آن از وسط یک Hit شروع کرده و به دو طرف با استفاده از برنامه نویسی پویا حرکت می کنیم. هر زمان که میزان از یک امتیاز خاص کمتر گردید، متوقف می شویم.

۴.۲ NCBI - BLAST

$BLAST$ برنامه ای است که به طور متداول برای جستجو در بانک های اطلاعاتی مورد استفاده قرار می گیرد. که هدف آن بدینگونه است که رشته های مشابه یک پرسش را بتواند بیابد. برنامه $BLAST+$ در $NCBI$ شامل برنامه های زیر می باشد:

^۳ High Scoring Segment
^۴ gapped extension

- BLASTp
- BLASTn
- TBLASTn
- BLASTx
- TBLASTx

همچنین ابزارهایی مانند *makeblastdb* و *blastdbcmd* وجود دارند. می توان با اجرای *help* – متوجه عملکرد این برنامه ها شد. مثال:

`%makeblastdb -help`

۱.۴.۲ *makeblastdb*

این برنامه رشته ها را تبدیل به فرمت شناخته شده برای *BLAST* می کند. به عنوان مثال:

`%makeblastdb -in swissprot -dbtype prot -parse -seqids`

تعریف هر کدام از موارد بالا به این صورت است:

- *dbtype prot* –: تعیین می کند که رشته های پروتئینی را در بانک اطلاعاتی را می خواهیم *index* کنیم. اگر بخواهیم نوکلئوتید ها را *index* کنیم از *nucl* استفاده می کنیم.
 - *Swissprot* – *in* : بخش *in* – تعیین می نماید که پایگاه داده با *FASTA* فرمت در فایل *Swissprot* به عنوان ورودی می باشد.
 - *parse – seqids* –: خط تعریف در فرمت *FASTA* را *Parse* می نماید. اگر این *option* را در اختیار نماییم، آنگاه توسط دستور *blastdbcmd* می توانیم رشته ها را خارج نماییم.
- پس استفاده از دستور *makeblastdb* یک سری فایل با فرمت های *..pni*, *.phr*, *.psi*, *.pmd*. تولید می کند.

۲.۴.۲ *blastdbcmd*

این برنامه برای خارج نمودن رشته ها از پایگاه داده به کار می آید. پارامترهای مهم آن به شرح زیر است:

- *db*: نام پایگاه داده
- *dbtype*: نوع داده که می تواند پروتئین (*prot*) یا نوکلئوتید (*nucl*) باشد.
- *entry*: نام رشته های مورد نظر که با حرف ، از هم جدا می شود.
- *range*
- *strand*: plus/minus
- *out*: نام فایل خروجی

به عنوان مثال می توانید دستور زیر را در نظر بگیرید:

`blastdbcmd -entry ABL1-HUMAN -db swissprot -range 1-20`

۵.۲ *blastp, blastn, etc*

پارامترهای مهم برای برنامه‌های *blastp, blastn, tblastn, tblastx* عبارتند از:

- *outfmt*: فرمت خروجی
- *query*: فایل ورودی
- *db*: اسم پایگاه داده
- *out*: فایل خروجی
- *word – size*: عدد صحیح بزرگتر یا مساوی دو برای الگوریتم *word finder*
- *num – descriptions*
- *num – alignments*
- *seq*

به عنوان نمونه می‌توانید دستور زیر را در نظر بگیرید:

```
blastp -query bcrabl-human.fa -db swissprot -out bcrable.blastp
```

۶.۲ *PHYLIP*

یک پکیج می‌باشد که در آن بیش از ۳۰ برنامه برای ساختن درخت تکامل وجود دارد. در این پکیج فرمت ورودی به شکل *PHYLIP* می‌باشد که برنامه *Clustalx* آن را تولید می‌نماید.