



دانشگاه صنعتی امیرکبیر
(پلیتکنیک تهران)
دانشکده مهندسی کامپیوتر

درس بیوانفورماتیک تمرین چهارم ۴

امیرمهدی زرین نژاد

۹۷۳۱۰۸۷

سوال ۱)

الف)

می‌توان از ماتریس **blosum** استفاده کرد. به این صورت که برای هر کاراکتر رشته داده‌شده، مقادیر ماتریس بلوسام مربوط به آن کاراکتر را قرار دهیم (ستون امتیازات مربوط به آن کاراکتر در ماتریس بلوسام). به این صورت ماتریس از امتیازات بلوسام کاراکترها پر می‌شود.

ب) profile drift:

یک مشکلی است که در **psi blast** می‌تواند رخ بدهد. و آن اینست که چون روند به صورت اتوماتیک انجام میشود و پروفایل‌ها به صورت خودکار ساخته می‌شوند، اگر در مرحله‌ای اشتباه صورت گیرد و توالی‌های نامناسب انتخاب شوند (تعدادی سیکوینس که هومولوگ نیستند اضافه شوند-سیکوئنس **false-positive** داشته باشیم) و یا پروفایلی غلط تولید شود، این مشکل بسط پیدا می‌کند و در پیمایش‌های بعدی نیز باقی می‌ماند و تجمیع و بیشتر می‌شود و هی سیکوئنس‌های مشابه آن سیکوئنس‌های غلط اضافه می‌شود. نهایتاً این اتفاق باعث تجمیع خطا و **low selectivity** می‌شود.

ج)

با کاهش آستانه در قسمت جست و جوی این الگوریتم، مشکل **profile drift** بیشتر می‌شود. زیرا اگر حد آستانه را کاهش دهیم، حساسیت برای میزان هومولوگ رشته‌های انتخاب شده کمتر می‌شود و امکان اضافه شدن توالی‌های **FP** (توالی غیر هومولوگ) و متفاوت‌تر بیشتر می‌شود. در نتیجه مشکل **profile drift** شدتش بیشتر می‌شود.

سوال ۲)

الف) ماتریس PSSM

تا ۳ رقم اعشار، $\text{Pseudocount} = 0$, $\text{background prob} = \text{randome chance} = 0.25$

ATGCCG
AAGATT
TACTCA
CTGAGG
CACCTG

	1	2	3	4	5	6	Overall freq
A	2/5	3/5	-	2/5	-	1/5	8/30
C	2/5	-	2/5	2/5	2/5	-	8/30
G	-	-	3/5	-	1/5	3/5	7/30
T	1/5	2/5	-	1/5	2/5	1/5	7/30

Normalizing: values/background prob(0.25)

	1	2	3	4	5	6	Overall freq
A	1.6	2.4	-	1.6	-	0.8	8/30
C	1.6	-	1.6	1.6	1.6	-	8/30
G	-	-	2.4	-	0.8	2.4	7/30
T	0.8	1.6	-	0.8	1.6	0.8	7/30

Log2

	1	2	3	4	5	6	Overall freq
A	0.678	1.263	-	0.678	-	-0.322	8/30
C	0.678	-	0.678	0.678	0.678	-	8/30
G	-	-	1.263	-	-0.322	1.263	7/30
T	-0.322	0.678	-	-0.322	0.678	-0.322	7/30

ب) logo

تنها برای ستون اول دنباله ها

$$H_u = - \sum_a p_{u,a} \cdot \log_2 p_{u,a}$$

$$I_u = \log_2 \frac{1}{p_u} - H_u$$

$$\Rightarrow H_u = (p_{u,A} \times \log_2 p_{u,A} + p_{u,C} \times \log_2 p_{u,C} + p_{u,T} \times \log_2 p_{u,T})$$

$$\Rightarrow H_u = \left(\frac{8}{30} \log_2 \frac{8}{30} + \frac{8}{30} \log_2 \frac{8}{30} + \frac{7}{30} \log_2 \frac{7}{30} \right) = 1.522$$

$$\Rightarrow I_u = \log_2 \frac{1}{p_u} - 1.522 = 0.1478$$

$$p_{u,a} \times I_u =$$

$$A: \frac{8}{30} \times 0.1478 = 0.1912$$

$$C: \frac{8}{30} \times 0.1478 = 0.1912$$

$$G: \frac{7}{30} \times 0.1478 = 0.1912$$

$$T: \frac{7}{30} \times 0.1478 = 0$$

سوال ۷ ب)
 سوال ۸
 سوال ۹

Profile (الف)

Pseudocount = 1, background = overall freq, تا 3 رقم اعشار

AT - G - CCG

AA - G - CTT

T - ACT - CA

CTGACGGA

	1	2	3	4	5	6	7	8	Overall freq
A	2/4	1/4	1/4	1/4	-	-	-	2/4	7/32
C	1/4	-	-	1/4	1/4	2/4	2/4	-	7/32
G	-	-	1/4	2/4	-	1/4	1/4	1/4	6/32
T	1/4	2/4	-	-	1/4	-	1/4	1/4	6/32
-	-	1/4	2/4	-	2/4	1/4	-	-	6/32

Update with pseudocount:

(Values+pseudocount)/(N + B*pseudocount), N=4, B=5=4+1(gap)

	1	2	3	4	5	6	7	8	Overall freq
A	2+1/4+ 5	1+1/4+ 5	1+1/4+ 5	1+1/4+ 5	0+1/4+ 5	0+1/4+ 5	0+1/4+ 5	2+1/4+ 5	7+8/32+5* 40 = 15/72
C	1+1/4+ 5	0+1/4+ 5	0+1/4+ 5	1+1/4+ 5	1+1/4+ 5	2+1/4+ 5	2+1/4+ 5	0+1/4+ 5	7+8/32+5* 8 = 15/72
G	0+1/4+ 5	0+1/4+ 5	1+1/4+ 5	2+1/4+ 5	0+1/4+ 5	1+1/4+ 5	1+1/4+ 5	1+1/4+ 5	6+8/32+5* 8 = 14/72
T	1+1/4+ 5	2+1/4+ 5	0+1/4+ 5	0+1/4+ 5	1+1/4+ 5	0+1/4+ 5	1+1/4+ 5	1+1/4+ 5	6+8/32+5* 8 = 14/72
-	0+1/4+ 5	1+1/4+ 5	2+1/4+ 5	0+1/4+ 5	2+1/4+ 5	1+1/4+ 5	0+1/4+ 5	0+1/4+ 5	6+8/32+5* 8 = 14/72

=

	1	2	3	4	5	6	7	8	Overall freq
A	0.333	0.222	0.222	0.222	0.111	0.111	0.111	0.333	0.208
C	0.222	0.111	0.111	0.222	0.222	0.333	0.333	0.111	0.208
G	0.111	0.111	0.222	0.333	0.111	0.222	0.222	0.222	0.194
T	0.222	0.333	0.111	0.111	0.222	0.111	0.222	0.222	0.194
-	0.111	0.222	0.333	0.111	0.333	0.222	0.111	0.111	0.194

Normalize with overall freq:

	1	2	3	4	5	6	7	8	Overall freq
A	1.601	0.534	0.534	0.534	1.067	1.067	1.067	1.601	0.208
C	0.534	1.601	1.601	1.067	1.067	0.534	0.534	1.067	0.208
G	1.144	1.144	1.144	0.572	1.716	1.144	0.572	0.572	0.194
T	1.144	1.144	0.572	1.144	0.572	0.572	1.716	1.144	0.194
-	0.572	0.572	1.144	1.716	0.572	1.716	1.144	0.572	0.194

Log2

	1	2	3	4	5	6	7	8	Overall freq
A	0.679	-0.905	-0.905	-0.905	0.094	0.094	0.094	0.679	
C	-0.905	0.679	0.679	0.094	0.094	-0.905	-0.905	0.094	
G	0.194	0.194	0.194	-0.805	0.779	0.194	-0.805	-0.805	
T	0.194	0.194	-0.805	0.194	-0.805	-0.805	0.779	0.194	
-	-0.805	-0.805	0.194	0.779	-0.805	0.779	0.194	-0.805	

ب) احتمال دنباله AA-CTCTG

$$0.679+0.094+0.779+0.094+0.194+0.679+0.194+0.194 = 2.907$$

یعنی احتمال این دنباله، $2^{2.907} = 7.5$ تقریباً ۷ و نیم برابر رندم چنس است.

ج) محتمل ترین دنباله

در هر ستون (جایگاه در دنباله) کاراکتری انتخاب می شود که بیش ترین مقدار امتیاز را دارد.

AT-G-CCA

	A	T	C	G (الف)
S_1	$1/4 \times 1/2 = 0.125$	$1/4 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$	$1/4 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$	$1/4 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$
S_2	$1/2 \times 0.1 = 0.05$	$1/2 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$	$1/2 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$	$1/2 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$
S_3	$0.1 \times 0.125 = 0.0125$	$0.1 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$	$0.1 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$	$0.1 \times (0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125 + 0.125 \times 0.125) = 0.125$

احتمال سبب ATCG = $0.00125 + 0.00125 + 0.00125 + 0.00125 = 0.005$

	C	G	A	T (ب)
S_1	$1/4 \times 0.125 = 0.03125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$
S_2	$1/2 \times 0.1 = 0.05$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$
S_3	$0.1 \times 0.125 = 0.0125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$	$\max(0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125, 0.125 \times 0.125) = 0.125$

حاصل کردن شماره صحت: $S_1 \rightarrow S_2 \rightarrow S_3 \rightarrow S_4$

ماکسیم مقدار در ستون آخر مربوط به S_3 است (0.0012) که اگر مسیرش به عقب را دنبال کنیم محتمل ترین دنباله حالت را می یابیم. (مسیرها توسط فلش مشخص شده اند همچنین استیت های مربوط به مقادیر ماکسیم در جلوی شان آمده است)