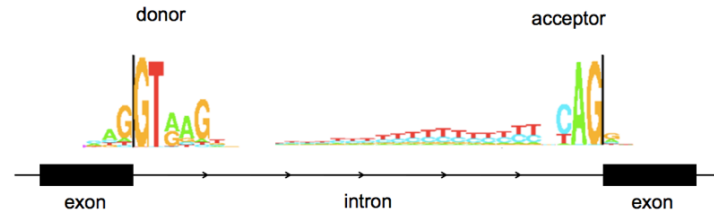
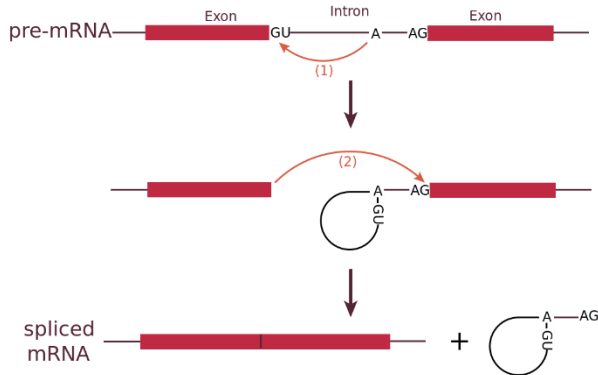


BIOINFORMATICS SESSION 13. PRACTICE

2023-11-27

Finding genes: In the world of snurps

RNA splicing and its sequence features



Identification of splice sites with a PSSM

	1	2	3	4	5	6	7	8	9
S1	C	A	G	G	T	A	G	G	G
S2	C	A	G	G	T	T	A	C	A
S3	A	A	G	G	T	A	T	G	T
S4	G	A	G	G	T	G	A	G	C
S5	G	A	G	G	T	A	A	A	C
S6	A	G	A	G	T	A	A	G	G
S7	C	G	G	G	T	G	G	G	T
S8	G	T	G	G	T	G	A	T	T
S9	A	C	A	G	T	A	A	C	T
S10	C	T	T	G	T	A	A	G	T

Frequency

	1	2	3	4	5	6	7	8	9
A	3	5	2	0	0	6	7	1	1
T	0	2	1	0	10	1	1	1	5
C	4	1	0	0	0	0	0	2	2
G	3	2	7	10	0	3	2	6	2

Pseudocount +1

	1	2	3	4	5	6	7	8	9
A	4	6	3	1	1	7	8	2	2
T	1	3	2	1	11	2	2	2	6
C	5	2	1	1	1	1	1	3	3
G	4	3	8	11	1	4	3	7	3

Identification of splice sites with a PSSM

	1	2	3	4	5	6	7	8	9
A	4	6	3	1	1	7	8	2	2
T	1	3	2	1	11	2	2	2	6
C	5	2	1	1	1	1	1	3	3
G	4	3	8	11	1	4	3	7	3

$$\frac{\text{Number of As at position 1}}{\text{Total number of sequences at position 1}} = \frac{4}{14} = 0.29$$



Probability

	1	2	3	4	5	6	7	8	9
A	0.29	0.43	0.21	0.07	0.07	0.50	0.57	0.14	0.14
T	0.07	0.21	0.14	0.07	0.79	0.14	0.14	0.14	0.43
C	0.36	0.14	0.07	0.07	0.07	0.07	0.07	0.21	0.21
G	0.29	0.21	0.57	0.79	0.07	0.29	0.21	0.50	0.21



Log observed/expected ratio

$$\log\left(\frac{\text{Actual probability of A at position 1}}{\text{Expected probability of A at position 1}}\right) = \log\left(\frac{0.29}{0.25}\right) = 0.13$$

	1	2	3	4	5	6	7	8	9
A	0.13	0.54	-0.15	-1.25	-1.25	0.69	0.83	-0.56	-0.56
T	-1.25	-0.15	-0.56	-1.25	1.15	-0.56	-0.56	-0.56	0.54
C	0.36	-0.56	-1.25	-1.25	-1.25	-1.25	-1.25	-0.15	-0.15
G	0.13	-0.15	0.83	1.15	-1.25	0.13	-0.15	0.69	-0.15

TCCTGTCC**CAGGTAGGG**AA



score = 5.11

Basic Shell Commands

```
$ cd [User_Folder]  
$ mkdir session13  
$ cd session13
```

Constructing a PSSM

```
$cp /home/biguser/tutor/session13/splice5.txt .  
$less splice5.txt
```

```
CAGGTAGGG  
CAGGTAACA  
AAGGTAAGT  
GAGGTGAGC  
GAGGTAAAC  
AAAGTAAGG
```

Constructing a PSSM

```
$ vi make_matrix5.py
```

```
1 import sys, math
2
3 splice5 = sys.argv[1] #splice5.txt
4 number_of_sequences = 0
5
6 for line in open(splice5):
7     line = line.rstrip()
8     if number_of_sequences == 0:
9         msa_matrix = [[]]
10    if number_of_sequences > 0:
11        msa_matrix.append([])
12    for j in range(0,9):
13        msa_matrix[number_of_sequences].append(line[j])
14    number_of_sequences += 1
15
16 print(msa_matrix)
```

Constructing a PSSM

```
$ python make_matrix5.py splice5.txt
```

```
[biguser@R440 session13]$ python make_matrix5.py splice5.txt  
[['C', 'A', 'G', 'G', 'T', 'A', 'G', 'G', 'G'], ['C', 'A', 'G', 'G', 'T', 'A', 'A', 'C', 'A'], ['A', 'A', 'G', 'G', 'T', 'A',  
'A', 'G', 'T'], ['G', 'A', 'G', 'G', 'T', 'G', 'A', 'G', 'C'], ['G', 'A', 'G', 'G', 'T', 'A', 'A', 'A', 'C'], ['A', 'A', 'A',  
'G', 'T', 'A', 'A', 'G', 'G']]
```



1	2	3	4	5	6	7	8	9
C	A	G	G	T	A	G	G	G
C	A	G	G	T	A	A	C	A
A	A	G	G	T	A	A	G	T
G	A	G	G	T	G	A	G	C
G	A	G	G	T	A	A	A	C
A	A	A	G	T	A	A	G	G

Constructing a PSSM

```
1 bases = ['A', 'T', 'C', 'G']
2 pssm = [[]]
3
4 for i in range(0, 4):
5     if i > 0:
6         pssm.append([])
7     for j in range(0, 9):
8         pssm[i].append(1.0) #pseudocount
9         for k in range(0, number_of_sequences):
10            if msa_matrix[k][j] == bases[i]:
11                pssm[i][j] += 1
12 print(pssm)
```

Constructing a PSSM

```
$ python make_matrix5.py splice5.txt
```

```
[biguser@R440 session13]$ python make_matrix5.py splice5.txt  
[[3.0, 7.0, 2.0, 1.0, 1.0, 6.0, 6.0, 2.0, 2.0], [1.0, 1.0, 1.0, 1.0, 7.0, 1.0, 1.0, 1.0, 2.0], [3.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 2.0, 3.0], [3.0, 1.0, 6.0, 7.0, 1.0, 2.0, 2.0, 5.0, 3.0]]
```



	1	2	3	4	5	6	7	8	9
A	3	7	2	1	1	6	6	2	2
T	1	1	1	1	7	1	1	1	2
C	3	1	1	1	1	1	1	2	3
G	3	1	6	7	1	2	2	5	3

Constructing a PSSM

```
1 for i in range(0, 4):
2     for j in range(0, 9):
3         pssm[i][j] = math.log(pssm[i][j] / (number_of_sequences + 4)*4 ) / math.log(2)
4         print(round(pssm[i][j], 2), end = ' ')
5     print('')
```

$$\log_2\left(\frac{\text{Actual probability of nucleotide at position } J}{\text{Expected probability of nucleotide at position } J}\right) =$$
$$\log_2\left(\frac{\text{Actual probability}}{0.25}\right) = \log_2(\text{Actual probability} * 4)$$

Syntax

```
math.log(x, base)
```

Parameter Values

Parameter	Description
<i>x</i>	Required. Specifies the value to calculate the logarithm for. If the value is 0 or a negative number, it returns a ValueError. If the value is not a number, it returns a TypeError
<i>base</i>	Optional. The logarithmic base to use. Default is 'e'

$$\frac{\log B}{\log A} = \log A B$$

Constructing a PSSM

```
$ python make_matrix5.py splice5.txt
```

```
[biguser@R440 session13]$ python make_matrix5.py splice5.txt  
0.26 1.49 -0.32 -1.32 -1.32 1.26 1.26 -0.32 -0.32  
-1.32 -1.32 -1.32 -1.32 1.49 -1.32 -1.32 -1.32 -0.32  
0.26 -1.32 -1.32 -1.32 -1.32 -1.32 -1.32 -0.32 0.26  
0.26 -1.32 1.26 1.49 -1.32 -0.32 -0.32 1.0 0.26
```



	1	2	3	4	5	6	7	8	9
A	0.26	1.49	-0.32	-1.32	-1.32	1.26	1.26	-0.32	-0.32
T	-1.32	-1.32	-1.32	-1.32	1.49	-1.32	-1.32	-1.32	-0.32
C	0.26	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-0.32	0.26
G	0.26	-1.32	1.26	1.49	-1.32	-0.32	-0.32	1.0	0.26

```
$ python make_matrix5.py splice5.txt > matrix5.txt
```

Scoring with a PSSM

```
$ cp /home/biguser/tutor/Week13/session13/amyloid.fa .  
$ less amyloid.fa
```

```
>amyloid
```

```
agccatcacttgtctctaataaataactcccattgattttccagctcagggctcaccact  
ccttacgtaagcgcaggaggagactggaaaatcactcacatattattgggtgctcttcct  
ccccatcctcacccaaggtgcatataaacctgaataacctgaagtctaaggcatgaa  
tatcagacgctagggggacagccactgtgttgtctgctaccctcatcctggtcactgctt  
ctgctataacagccctaggccaggaatatgaacaagccgctgctttggatctctgtcctc  
accagcctcctggaagcctttgctcacacaggtgaaggaggtgaaggaatgggtcaagaatc  
ataaagtgagaaaataggttgaagctgagatatctttccctgcatttatactgaaggtc  
attatctttctttctttatcccgcagacctcagtggggaaggtgtttgtatttcctagaga
```

Scoring with a PSSM

```
$ vi score5.py
```

```
1 # score5.py
2
3 import sys, re
4
5 matrix5 = sys.argv[1] #matrix5.txt
6 amyloid = sys.argv[2] #amyloid.fa
7
8 i = 0
9
10 for line in open(matrix5):
11     line = line.rstrip()
12     if i == 0:
13         pssm = [[]]
14     if i > 0:
15         pssm.append([])
16     col = line.split()
17     for j in range(0, 9):
18         pssm[i].append(float(col[j]))
19     i += 1
20
21 seq = ''
22 for line in open(amyloid):
23     if not re.search('>', line):
24         line = line.rstrip()
25         seq += line
26
27 print('pos\tscore') # print header
28
29 seq = seq.upper() # covert the sequence to upper case letters
30 bases = ['A', 'T', 'C', 'G']
31
32 for k in range(0, len(seq) - 8):
33     test = seq[k:k+9]
34     score = 0
35     for j in range(0, 9):
36         base = test[j]
37         for b in range(0, 4):
38             if bases[b] == base:
39                 score += pssm[b][j]
40     score = 2 ** score # convert the log2 to real values, ** : exponential operator
41     pos = k + 3 # print the position next to the exon-intron junction
42
43     print(pos, '\t', score)
```

Scoring with a PSSM

```
$ python score5.py matrix5.txt amyloid.fa
```

```
[biguser@R440 session13]$ python score5.py matrix5.txt amyloid.fa
pos      score
3        0.004742948767168147
4        0.13304627280666997
5        0.018971795068672588
6        0.016630784100833736
7        0.007921558435859595
8        0.022250784306204228
9        0.3120826372254029
10       0.002371474383584075
11       0.1560413186127015
12       0.0031728609232665435
```

```
$ python score5.py matrix5.txt amyloid.fa > score5.txt
```

Visualization with R

```
$ cp /home/biguser/tutor/session13/score3.txt .
$ cp /home/biguser/tutor/session13/amyloid.r .
$ vi amyloid.r
```

```
# plot the results of splice site prediction

# define some colours
rgb <- c("#009E73", "#D55E00", "#0072B2")

# make two graphs on top of each other
par(mfrow = c(2, 1))

# First consider the 5' splice site prediction and read the
# output from
# the Perl code
data <- read.table("score5.txt", sep = "\t", header = TRUE)

# for the plot, we need to know about the sequence length
seqlen <- max(data$pos)

# for the plot we need to know about the maximum score
max_score <- max(data$score)

# make a plot for the 5' splice site data
plot(0, type = "n", lwd = 2, xlim = c(0, seqlen),
     ylim = c(0, max_score * 1.1), main = "Splice site scoring",
     xlab = "Position", ylab = "Score")

# print a legend
legend(seqlen * 0.7, max_score, "5prime", col = rgb[2],
       lwd = 1)

# plot the splice site scores
for (i in (1:seqlen)) {
  lines(c(data$pos[i], data$pos[i]), c(0, data$score[i]), col = rgb[2],
    lw = 2)
}
```

```
# plot the location of exons (we do not know these from the
# prediction)
lines(c(268, 331), c(max_score/2, max_score/2), col = rgb[1],
      lw = 4)
lines(c(447, 1054), c(max_score/2, max_score/2), col = rgb[1],
      lw = 4)

# Now consider the 3' splice site and read the output from
# the Perl code
data <- read.table("score3.txt", sep = "\t", header = TRUE)
max_score <- max(data$score)

plot(0, type = "n", lwd = 2, xlim = c(0, seqlen),
     ylim = c(0, max_score * 1.1), xlab = "Position", ylab = "Score")

legend(seqlen * 0.7, max_score, "3prime", col = rgb[3],
       lwd = 1)

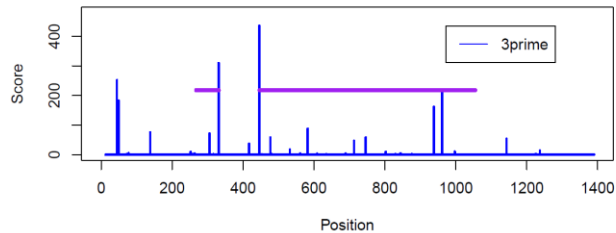
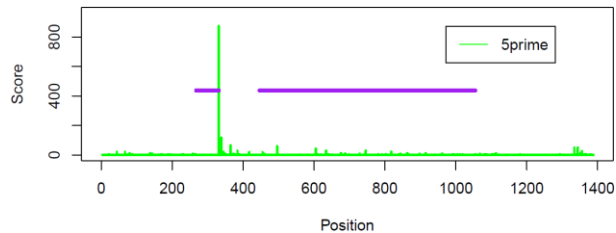
for (i in (1:seqlen)) {
  lines(c(data$pos[i], data$pos[i]), c(0, data$score[i]), col = rgb[3],
    lw = 2)
}

# plot the location of exons as in the previous graph
lines(c(268, 331), c(max_score/2, max_score/2), col = rgb[1],
      lw = 4)
lines(c(447, 1054), c(max_score/2, max_score/2), col = rgb[1],
      lw = 4)
```


Visualization with R

```
$ Rscript amyloid.r  
Open "Rplots.pdf"
```

Splice site scoring



Exercise

- When you construct a PSSM, you divided observed frequency by expected frequency ($M_{ij} = \log(F_{ij}/F_{exp})$, $F_{exp} = 0.25$). In real, however, those four bases (A, C, G, T) are not evenly distributed across a genome, which means that expected frequencies for each base are not equally 0.25 (1/4). It would be much more precise to use 'observed frequency across the matrix' ($M_{ij} = \log(F_{ij}/F_{exp})$, $F_{exp} = P_i/\text{Total } P$, $i=(A,C,G,T)$). Correct a Python script 'make_matrix5.py' to calculate M_{ij} which is divided by 'observed frequency across matrix'.

	1	2	3	4	5	6	7	8	9
A	0.29	0.43	0.21	0.07	0.07	0.50	0.57	0.14	0.14
T	0.07	0.21	0.14	0.07	0.79	0.14	0.14	0.14	0.43
C	0.36	0.14	0.07	0.07	0.07	0.07	0.07	0.21	0.21
G	0.29	0.21	0.57	0.79	0.07	0.29	0.21	0.50	0.21



	1	2	3	4	5	6	7	8	9
A	0.06	0.46	-0.23	-1.33	-1.33	0.62	0.75	-0.64	-0.64
T	-1.20	-0.11	-0.51	-1.20	1.19	-0.51	-0.51	-0.51	0.59
C	0.92	0.00	-0.69	-0.69	-0.69	-0.69	-0.69	0.41	0.41
G	-0.20	-0.49	0.49	0.81	-1.59	-0.20	-0.49	0.36	-0.49

Log-odd: $\log(F_{ij}/F_{exp})$

$F_{exp} = P_i/\text{Total } P$,
 $i=(A,C,G,T)$

Exercise

```
1 # exercise
2
3 for i in range(0, 4):
4     baseFreq = sum(pssm[i]) / (sum(pssm[0]) + sum(pssm[1]) + sum(pssm[2]) + sum(pssm[3]))
5     for j in range(0, 9):
6         pssm[i][j] = math.log(pssm[i][j] / (number_of_sequences + 4) / baseFreq) / math.log(2)
7         print(round(pssm[i][j], 2), end = ' ')
8     print('')
```