

Untitled

LuchaoQi

February 28, 2019

Q1

1a

Compute the edit distance:

Alpha: EALERMFLSFPTTKTYFPHFDLSHGSAQVK

Beta: EALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVK

```
str1 = 'EALERMFLSFPTTKTYFPHFDLSHGSAQVK'
str2 = 'EALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVK'
m = len(str1)
n = len(str2)
dp = [[0 for x in range(n+1)] for x in range(m+1)]
for i in range(m+1):
    for j in range(n+1):
        if i == 0:
            dp[i][j] = j
        elif j == 0:
            dp[i][j] = i
        elif str1[i-1] == str2[j-1]:
            dp[i][j] = dp[i-1][j-1]
        else:
            dp[i][j] = 1 + min(dp[i][j-1], dp[i-1][j], dp[i-1][j-1])
for i in range(len(dp)):
    print(dp[i])
```

Results:

So the edit distance is 22.

Sequences:

EALERMFLSFPTTKTYFPHF - DLS - - - - HGSAQVK
EALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKVK

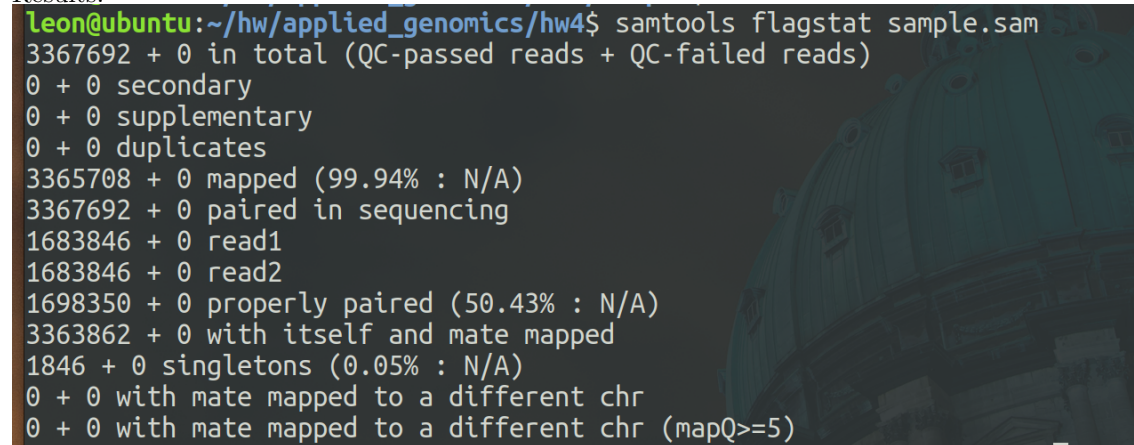
Q2

2a

Commands:

```
bowtie2-build chr22.fa chr22
bowtie2 -x chr22 -1 sample/pair.1.fq -2 sample/pair.2.fq > sample.sam
samtools flagstat sample.sam
```

Results:



```
leon@ubuntu:~/hw/applied_genomics/hw4$ samtools flagstat sample.sam
3367692 + 0 in total (QC-passed reads + QC-failed reads)
0 + 0 secondary
0 + 0 supplementary
0 + 0 duplicates
3365708 + 0 mapped (99.94% : N/A)
3367692 + 0 paired in sequencing
1683846 + 0 read1
1683846 + 0 read2
1698350 + 0 properly paired (50.43% : N/A)
3363862 + 0 with itself and mate mapped
1846 + 0 singletons (0.05% : N/A)
0 + 0 with mate mapped to a different chr
0 + 0 with mate mapped to a different chr (mapQ>=5)
```

So 3365708 reads align to the chr22 reference. 1984 reads did not align. 1846 reads had AKA singletons.

2b

Commands:

```
samtools view -c -f 16 sample.sam
```

We can get 1678219 reads mapped to the reverse strand.

2c

Commands:

```
samtools sort sample.sam > sample.bam
samtools index sample.bam
freebayes -f chr22.fa sample.bam > sample.vcf
bcftools filter -e "QUAL<20" sample.vcf > filtered.vcf
bcftools stats filtered.vcf > stats.txt
vim stats.txt
```

Results:

```
leon@ubuntu: ~/hw/applied_genomics/hw4
File Edit View Search Terminal Help
15 # number of multiallelic sites .. number of rows with multiple alternate alleles
16 # number of multiallelic SNP sites .. number of rows with multiple alternate alleles, all SNPs
17 #
18 # Note that rows containing multiple types will be counted multiple times, in each
19 # counter. For example, a row with a SNP and an indel increments both the SNP and
20 # the indel counter.
21 #
22 # SN [2]id [3]key [4]value
23 SN 0 number of samples: 1
24 SN 0 number of records: 16536
25 SN 0 number of no-ALTs: 0
26 SN 0 number of SNPs: 14965
27 SN 0 number of MNPs: 152
28 SN 0 number of indels: 1403
29 SN 0 number of others: 18
30 SN 0 number of multiallelic sites: 6
31 SN 0 number of multiallelic SNP sites: 4
32 # TSTV, transitions/transversions:
33 # TSTV [2]id [3]ts [4]tv [5]ts/tv [6]ts (1st ALT) [7]tv (1st ALT) [8]ts/tv (1st ALT)
34 TSTV 0 5342 9612 0.56 5339 9610 0.56
35 # SiS, Singleton stats:
36 # SiS [2]id [3]allele count [4]number of SNPs [5]number of transitions [6]number of transversions [7]
number of indels [8]repeat-consistent [9]repeat-inconsistent [10]not applicable
37 SiS 0 1 12990 4472 8518 1231 0 0 1231 25,1 1%
```

So there are 14965 single nucleotides with QUAL>20 and 1403 indels.

```
810 # IDD [2]id [3]length (deletions negative) [4]count
811 IDD 0 -6 1
812 IDD 0 -5 4
813 IDD 0 -4 10
814 IDD 0 -3 60
815 IDD 0 -2 122
816 IDD 0 -1 510
817 IDD 0 1 483
818 IDD 0 2 144
819 IDD 0 3 49
820 IDD 0 4 12
821 IDD 0 5 6
822 IDD 0 9 2
```

Of the indels, there are 707 deletions, 696 insertions.

Q3

3a

Python:

```
import pysam
chr22 = pysam.FastaFile("chr22.fa")
seq = chr22.fetch("chr22",21000000,22000000)

sub = open("sub.fa","w")
sub.write(">chr22:21000000,22000000 \n")
sub.write(seq + '\n')
sub.close()

reads = open("reads.fa","w")
for i in range(len(seq)-35+1):
    number = ">" + str(i+1) + "#pos in chr22"
    subseq = seq[i:i+35]
    reads.write(number + "\n")
    reads.write(subseq + "\n")
reads.close()
```

Commands:

```
python3 reads.py #get the ref and pseudo reads
bowtie2-build sub.fa sub #build up the index
bowtie2 -x sub -f -U reads.fa -S sub.sam
```

Results:

```
(base) leon@ubuntu:~/hw/applied_genomics/hw4$ bowtie2 -x sub -f -U reads.fa -S sub.sam
999966 reads; of these:
  999966 (100.00%) were unpaired; of these:
    0 (0.00%) aligned 0 times
    595713 (59.57%) aligned exactly 1 time
    404253 (40.43%) aligned >1 times
100.00% overall alignment rate
```

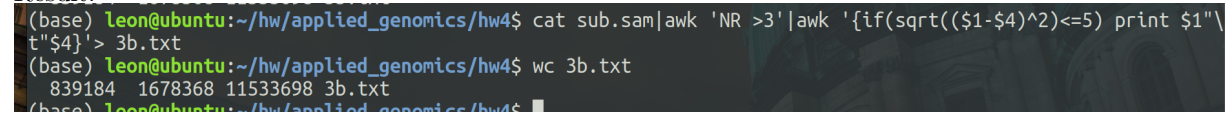
So 404253 aligned more than one times. The alignment rate is 100%. 0 read did not align.

3b

Commands:

```
samtools sub.sam -o sub.sam  
cat sub.sam|awk 'NR >3'|awk '{if(sqrt(($1-$4)^2)<=5) print $1"\t"$4}'> 3b.txt  
wc 3b.txt
```

Result:

A terminal window screenshot showing the execution of the commands. The first command is 'cat sub.sam|awk 'NR >3'|awk '{if(sqrt((\$1-\$4)^2)<=5) print \$1"\t"\$4}'> 3b.txt'. The second command is 'wc 3b.txt', which outputs '839184 1678368 11533698 3b.txt'.

```
(base) leon@ubuntu:~/hw/applied_genomics/hw4$ cat sub.sam|awk 'NR >3'|awk '{if(sqrt(($1-$4)^2)<=5) print $1"\t"$4}'> 3b.txt  
(base) leon@ubuntu:~/hw/applied_genomics/hw4$ wc 3b.txt  
839184 1678368 11533698 3b.txt  
(base) leon@ubuntu:~/hw/applied_genomics/hw4$
```

So 839184 reads aligned correctly(within 5bp)

3c

Commands:

```
cat sub.sam|awk 'NR>3 {print $1"\t"$5}' >3c.txt
```

Python:

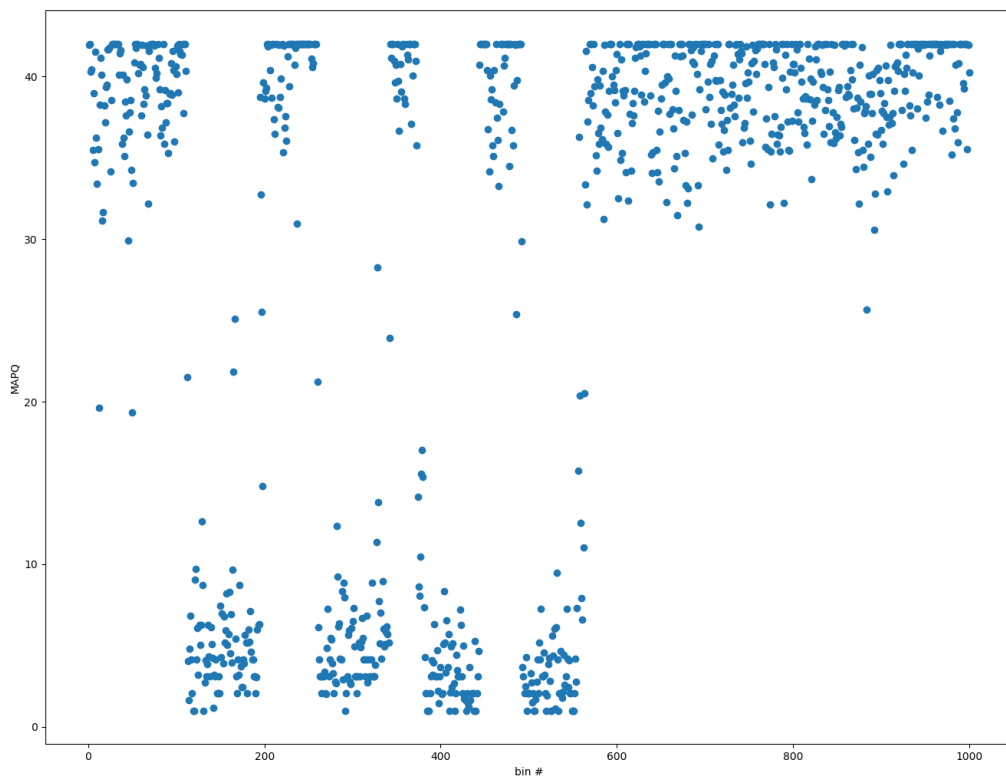
```
import matplotlib.pyplot as plt

f = open('3c.txt','r')
dic = {}
for lines in f:
    line = lines.strip().rstrip('\n').split('\t')
    a = int(line[0])//1000 + 1
    if a in dic:
        dic[a] = dic[a] + int(line[1])
    else:
        dic[a] = int(line[1])
for i in dic:
    dic[i] = dic[i]/1000

f.close()

plt.scatter(dic.keys(),dic.values())
plt.xlabel('bin #')
plt.ylabel('MAPQ')
plt.show()
```

Results:



3d

Commands:

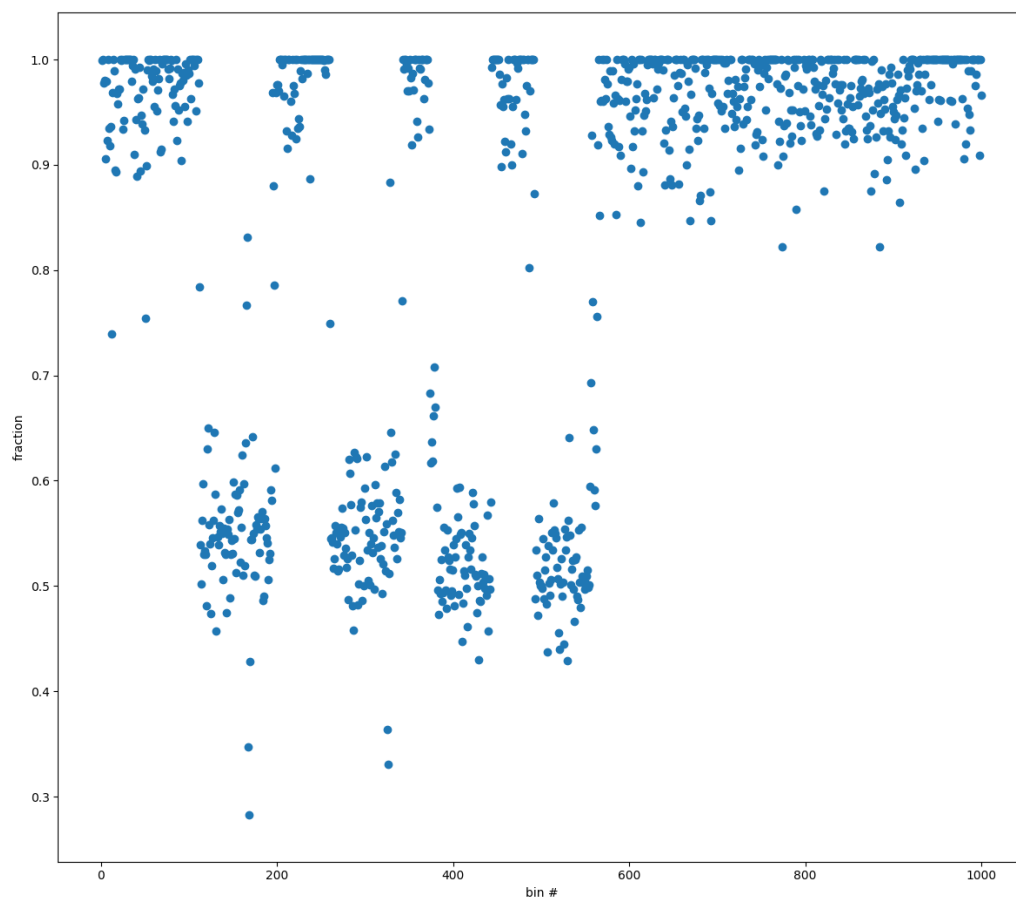
```
cat sub.sam|awk 'NR >3'|awk '{print $1"\t"$4}'> 3d.txt
```

Python:

```
import matplotlib.pyplot as plt
f = open('3d.txt','r')
dic = {}
for lines in f:
    line = lines.strip().rstrip('\n').split('\t')
    a = int(line[0])//1000 + 1
    if a in dic:
        dic[a] = dic[a] + (line[0] == line[1])
    else:
        dic[a] = (line[0]==line[1])
for i in dic:
    dic[i] = dic[i]/1000

f.close()
plt.scatter(dic.keys(),dic.values())
plt.xlabel('bin #')
plt.ylabel('fraction')
plt.show()
```

Results:

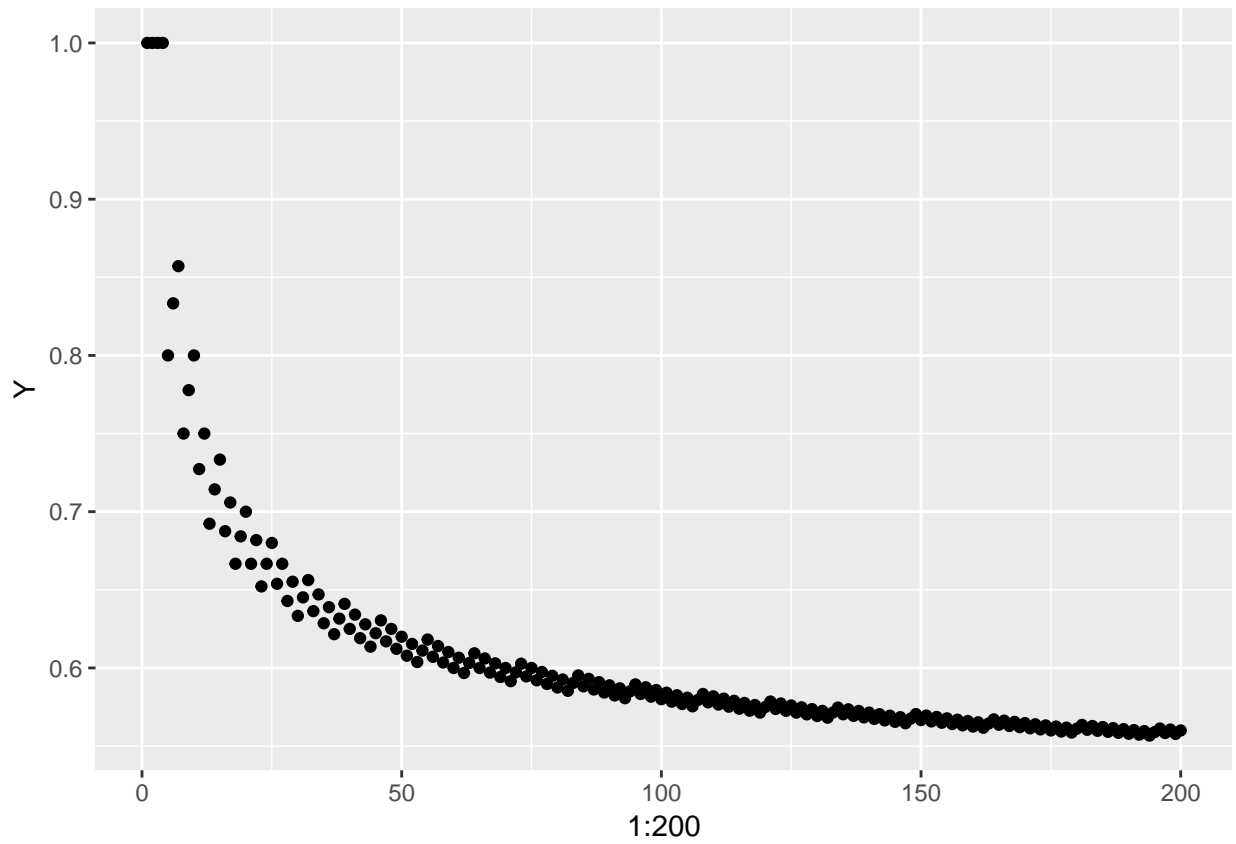


Looks similar to the plot in 3c.

Q4

4a

```
X = 1:200
Y = vector()
for (size in X)
{
  Y[size] = qbinom(0.95,size,0.5)/size
}
library(ggplot2)
ggplot(data = as.data.frame(Y),aes(x=1:200,y=Y))+geom_point()
```



4b

The plot seems to approach 0.5 because the larger the size is, according to the binomial distribution ($p=0.5$), the closer 95% quantile is to the 50% of the size.