

# High throughput

## The second homework

Name: Chengzhi Dong

Student ID: 19331027

### 1 Code

#### 1.1 Directory operations

1. Create a directory named learning\_python.

```
1 mkdir learning_python
```

2. Switch to learning\_python directory

```
1 cd learning_python
```

3. Put the chrM.fa file in this directory.

```
1 wget http://222.200.186.169/highput/data/chrM.fa
```

4. Open vim editor with *vim count.py*

```
1 vim count.py
```

#### 1.2 Python in biology

Edit your python code in vim editor, then save and exit. The program should include the following:

1. Calculate the percentage of each base of ATCG in the chrM chromosome.

```
1 #Question 1
2 def percent_cal_1(DNA_seq):
3     count_A = 0
4     count_T = 0
5     count_C = 0
6     count_G = 0
7     for i in DNA_seq:
8         if i == "A":
9             count_A += 1
10        elif i == "T":
11            count_T += 1
12        elif i == "C":
13            count_C += 1
14        elif i == "G":
15            count_G += 1
16    all_count = len(DNA_seq)
17    precent_A = count_A*100/all_count
18    precent_T = count_T*100/all_count
19    precent_C = count_C*100/all_count
```

```

20     precent_G = count_G*100/all_count
21
22     print("A | Frequency:{0} | Percentage:{1:.2f}%".format(count_A,precent_A))
23     print("T | Frequency:{0} | Percentage:{1:.2f}%".format(count_T,precent_T))
24     print("C | Frequency:{0} | Percentage:{1:.2f}%".format(count_C,precent_C))
25     print("G | Frequency:{0} | Percentage:{1:.2f}%".format(count_G,precent_G))
26
27 def percent_Cal_2(DNA_seq):
28     count_ls = set(DNA_seq)
29     all_count = len(DNA_seq)/100
30     for x in count_ls:
31         frequency = DNA_seq.count(x)
32         percent = frequency/all_count
33         print("{0} | Frequency:{1} | Percentage:{2:.2f}%".format(x,frequency,percent))
34
35 def percent_Cal_3(DNA_seq):
36     count_ls = dict()
37     all_count = len(DNA_seq)/100
38     for x in DNA_seq:
39         if x in count_ls:
40             count_ls[x] += 1
41         else:
42             count_ls[x] = 1
43     for key in count_ls:
44         count = count_ls[key]
45         percent = count/all_count
46         print("{0} | Frequency:{1} | Percentage:{2:.2f}%".format(key,count,percent))
47
48
49 def text_to_DNA(text_file_name):
50     text_file = open(text_file_name)
51     lines = text_file.readlines()
52     text_file.close()
53     line = [x.strip() for x in lines[1:]]
54     lines = ''.join(line)
55     DNA_seq = lines.upper()
56     return DNA_seq
57
58 text_file_name = "chrM.fa"
59 DNA_seq = text_to_DNA(text_file_name)
60 print("Question1: Calculate the percentage of each base of ATCG in the chrM chromosome by
        Funtion1.1:")
61 percent_Cal_1(DNA_seq)
62 print("Question1: Calculate the percentage of each base of ATCG in the chrM chromosome by by
        Funtion1.2:")
63 percent_Cal_2(DNA_seq)
64 print("Question1: Calculate the percentage of each base of ATCG in the chrM chromosome by by
        Funtion1.3:")
65 percent_Cal_3(DNA_seq)

```

2. Calculate insulin = "GIVEQCCTSICSLYQLENYCNFVNQHLCGSHLVEALYLVCGERGFFYTPKT"  
amino acid frequency in insulin sequence.

```

1 #Question 2
2 insulin = "GIVEQCCTSICSLYQLENYCNFVNQHLCGSHLVEALYLVCGERGFFYTPKT"
3 print("Question2: Calculate insulin = GIVEQCCTSICSLYQLENYCNFVNQHLCGSHLVEALYLVCGERGFFYTPKT amino
        acid frequency in insulin sequence by Funtion1.2:")
4 percent_Cal_2(insulin)
5 print("Question2: Calculate insulin = GIVEQCCTSICSLYQLENYCNFVNQHLCGSHLVEALYLVCGERGFFYTPKT amino
        acid frequency in insulin sequence by Funtion1.3:")
6 percent_Cal_3(insulin)

```

3. Run 5bp windows in the sequence "PRQTEINSEQWENCE" and count the number of occurrences of each window in the sequence.

```

1 #Question 3
2 def percent_Cal_3_1(DNA_seq):
3     count_ls = dict()
4     for i in range(len(DNA_seq)-4):
5         seq_short = DNA_seq[i:i+5]
6         if seq_short in count_ls:
7             count_ls[seq_short] += 1
8         else:
9             count_ls[seq_short] = 1
10    for key in count_ls:
11        count = count_ls[key]
12        print("{0} | Count:{1}".format(key, count))
13
14 print("Question3: Run 5bp windows in the sequence PRQTEINSEQWENCE and count the number of
15     occurrences of each window in the sequence by Funtion3.1:")
16 sequence = "PRQTEINSEQWENCE"
17 percent_Cal_3_1(sequence)

```

4. Calculate GC percentage in chrM.fa

```

1 # Question4
2 print("Question4: Calculate GC percentage in chrM.fa:")
3 count_GC = (DNA_seq.count("G")+DNA_seq.count("C"))/len(DNA_seq)
4 print("GC percentage in chrM: {0:.2f}".format(count_GC))

```

5. Use *python3 count.py* to run your program from the command line

```

1 python3 count.py

```

## 2 Result

1. Question1

```

Question1: Calculate the percentage of each base of ATCG in the chrM chromosome by Funtion1.1:
A | Frequency:5113 | Percentage:30.86%
T | Frequency:4086 | Percentage:24.66%
C | Frequency:5192 | Percentage:31.33%
G | Frequency:2180 | Percentage:13.16%
Question1: Calculate the percentage of each base of ATCG in the chrM chromosome by by
Funtion1.2:
A | Frequency:5113 | Percentage:30.86%
T | Frequency:4086 | Percentage:24.66%
G | Frequency:2180 | Percentage:13.16%
C | Frequency:5192 | Percentage:31.33%
Question1: Calculate the percentage of each base of ATCG in the chrM chromosome by by
Funtion1.3:
G | Frequency:2180 | Percentage:13.16%
A | Frequency:5113 | Percentage:30.86%
T | Frequency:4086 | Percentage:24.66%
C | Frequency:5192 | Percentage:31.33%

```

图 1: The percentage of each base of ATCG in the chrM chromosome.

2. Question2

```

Question2: Calculate insulin = GIVEQCCTSIQSLYQLENYCNFVNQHLCGSHLVEALYLVCGERGFFYTPKT amino acid
frequency in insulin sequence by Funtion1.3:
G | Frequency:4 | Percentage:7.84%
I | Frequency:2 | Percentage:3.92%
V | Frequency:4 | Percentage:7.84%
E | Frequency:4 | Percentage:7.84%
Q | Frequency:3 | Percentage:5.88%
C | Frequency:6 | Percentage:11.76%
T | Frequency:3 | Percentage:5.88%
S | Frequency:3 | Percentage:5.88%
L | Frequency:6 | Percentage:11.76%
Y | Frequency:4 | Percentage:7.84%
N | Frequency:3 | Percentage:5.88%
F | Frequency:3 | Percentage:5.88%
H | Frequency:2 | Percentage:3.92%
A | Frequency:1 | Percentage:1.96%
R | Frequency:1 | Percentage:1.96%
P | Frequency:1 | Percentage:1.96%
K | Frequency:1 | Percentage:1.96%

```

图 2: Insulin amino acid frequency in insulin sequence.

### 3. Question3

```

Question3: Calculate insulin = GIVEQCCTSIQSLYQLENYCNFVNQHLCGSHLVEALYLVCGERGFFYTPKT amino acid
frequency in insulin sequence by Funtion3.1:
PRQTE | Count:1
RQTEI | Count:1
QTEIN | Count:1
TEINS | Count:1
EINSE | Count:1
INSEQ | Count:1
NSEQW | Count:1
SEQWE | Count:1
EQWEN | Count:1
QWENC | Count:1
WENCE | Count:1

```

图 3: The number of occurrences of each window in the sequence.

### 4. Question4

```

Question4: Calculate GC percentage in chrM.fa:
GC percentage in chrM: 0.44

```

图 4: GC percentage in chrM.fa.