

목차

- [생물 정보학 프로그래밍](#)
 - [1일 차: 표준 입출력과 조건문 사용에 대한 프로그래밍 연습 I](#)
 - [2일 차: 표준 입출력과 조건문 사용에 대한 프로그래밍 연습 II](#)
 - [3일 차: 파일 핸들 및 Fasta 포맷 파일 파싱](#)
 - [4일차: GenBank 포맷 파싱, 문자열 포매팅, 서열 핸들\(상보적서열, 역서열, 등 생성\), 폴드 핸들](#)
 - [5일차: 6-frame translation](#)

생물 정보학 프로그래밍

1일 차: 표준 입출력과 조건문 사용에 대한 프로그래밍 연습 I

1. "This is sequence file"을 다음과 같이 출력한다.

```
1 | $ python bioinformatics_1_1.py
2 | This is sequence file
3 | $
```

2. 문자열 입력받고 받은 문자열을 다음과 같이 출력하는 .

```
1 | $ python bioinformatics_1_2.py
2 | Enter a string: Good day, everyone!
3 | Good day, everyone!
4 | $
```

3. 두 정수를 입력받아 합을 구하여 다음과 같이 출력한다.

```
1 | $ python bioinformatics_1_3.py
2 | Enter a integer: 10
3 | Enter another: 9
4 | Sum: 19
5 | $
```

4. 두 정수를 입력받아 큰 수를 출력하는 프로그램을 작성하시오.

```
1 | $ python bioinformatics_1_3.py
2 | Enter a integer: 10
3 | Enter another: 9
4 | Sum: 19
5 | $
```

2일 차: 표준 입출력과 조건문 사용에 대한 프로그래밍 연습 II

1. 구구단의 단 수를 입력받아 출력하는 프로그램을 작성하시오. 단, 단 수는 1과 9사이의 수만을 입력받아야 한다.

```
1 | $ python bioinformatics_2_1.py
2 | Which times table: 5
3 | 5 * 1 = 5
4 | 5 * 2 = 10
5 | 5 * 3 = 15
6 | 5 * 4 = 20
7 | 5 * 5 = 25
8 | 5 * 6 = 30
9 | 5 * 7 = 35
10 | 5 * 8 = 40
11 | 5 * 9 = 45
12 | $
13 | $ python bioinformatics_2_1.py
14 | Which times table: 12
15 | What?
16 | $
```

2. 문자열을 입력받아 포함된 문자의 갯수를 출력하는 프로그램을 완성하시오.

```
1 | $ python bioinformatics_2_2.py
2 | Enter a string: Good day, everyone!
3 | The string length is 19.
4 | $
```

3. 문자열과 숫자를 입력받아 입력받은 문자열을 입력한 수만큼 반복하여 출력하는 프로그램을 완성하시오.

```
1 | $ python bioinformatics_2_3.py
2 | Enter a string: I am a biologist.
3 | How many times to repeat: 3
4 | I am a biologist.I am a biologist.I am a biologist.
5 | $
```

4. 두 개의 문자열 입력받아 서로 다르다면 두 문자열을 연결하여 출력하고, 같으면 "Two strings are identical."를 출력

하는 프로그램을 완성하시오.

```
1 | $ python bioinformatics_2_4.py
2 | Enter a string: I am a computer scientist.
3 | Enter another: I am a biologist.
4 | I am a computer scientist.I am a biologist.
5 | $
6 | $ python bioinformatics_2_4.py
7 | Enter a string: I am a bioinformatician.
8 | Enter another: I am a bioinformatician.
9 | Two strings are identical.
10 | $
```

5. 두 개의 문자열 `s1` 과 `s2` 를 입력받아 `s1` 의 길이가 홀수이고 `s2` 보다 짧으면 `s1` , `s2` 의 순서로 출력하고, 그렇지 않으면 반대 순서로 출력하는 프로그램을 완성하시오.

```
1 | $ python bioinformatics_2_5.py
2 | Enter s1: 12345
3 | Enter s2: ABCDEF
4 | 12345ABCDEF
5 | $
6 | $ python bioinformatics_2_5.py
7 | Enter s1: 12345
8 | Enter s2: ABCDE
9 | ABCDE12345
10 | $
11 | $ python bioinformatics_2_5.py
12 | Enter s1: 123456
13 | Enter s2: ABCDE
14 | ABCDE123456
15 | $
```

6. 문자열을 입력받아 역순으로 출력하는 프로그램을 작성하시오.

```
1 | $ python bioinformatics_2_6.py
2 | Enter a string: I am a biologist.
3 | Reversed string: .tsigoloib a ma I
4 | $
```

7. 다음과 같이 3-base codon과 아미노산을 입력받아 테이블을 만든다. 단, codon에 `xxx` 를 입력하면 테이블 만들기를 종료한다. (입력에 필요한 codon과 이에 해당되는 아미노산은 다음의 테이블을 참조한다.) 이제 codon을 입력받으면 테이블에서 입력한 codon에 해당되는 아미노산을 출력한다. 단, 3-base codon에 `xxx` 를 입력하면 종료한다.

```

1 | $ python bioinformatics_2_7.py
2 | Enter three-base codon to build: GCT
3 | Enter amino acid: A
4 | Enter three-base codon to build: GCC
5 | Enter amino acid: A
6 | ...
7 | Enter three-base codon to build: TAG
8 | Enter amino acid: *
9 | Enter three-base codon to build: XXX
10 | Okay, I will switch.
11 | Enter three-base codon to search: CCG
12 | Amino acid for CCG: P
13 | Enter three-base codon to search: GAC
14 | Amino acid for GAC: D
15 | Enter three-base codon to search: XXX
16 | Okay, I will stop.
17 | $

```

Amino acid	codons	Amino acid	codons
Ala/A	GCT , GCC , GCA , GCG	Leu/L	TTA , TTG , CTT , CTC , CTA , CTG
Arg/R	CGT , CGC , CGA , CGG , AGA , AGG	Lys/K	AAA , AAG
Asn/N	AAT , AAC	Met/M	ATG
Asp/D	GAT , GAC	Phe/F	TTT , TTC
Cys/C	TGT , TGC	Pro/P	CCT , CCC , CCA , CCG
Gln/Q	CAA , CAG	Ser/S	TCT , TCC , TCA , TCG , AGT , AGC
Glu/E	GAA , GAG	Thr/T	ACT , ACC , ACA , ACG
Gly/G	GGT , GGC , GGA , GGG	Trp/W	TGG
His/H	CAT , CAC	Tyr/Y	TAT , TAC
Ile/I	ATT , ATC , ATA	Val/V	GTT , GTC , GTA , GTG
Met/M	ATG	Stop/*	TAA , TGA , TAG

3일 차: 파일 핸들 및 Fasta 포맷 파일 파싱

1. "This is a sequence file."이라는 내용을 가진 `title.txt` 파일을 만든다.

```
1 | $ python bioinformatics_3_1.py
2 | $ cat title.txt
3 | This is a sequence file.
4 | $
```

2. 위에서 만든 `title.txt` 파일의 첫째줄을 읽어 다음과 같이 출력한다.

```
1 | $ python bioinformatics_3_2.py
2 | The first line is: This is a sequence file.
3 | $
```

3. <https://www.ncbi.nlm.nih.gov/protein/AAB07223.1?report=fasta> 에서 protein FASTA 파일을 받아 `sequence.protein.fasta` 란 이름으로 저장한 후 그 내용을 읽어 다음과 같이 출력한다.

```
1 | $ python bioinformatics_3_3.py
2 | >AAB07223.1 BRCA2 [Homo sapiens]
3 | MPIGSKERPTFFEIFKTRCNKADLGPISLNWFEELSSEAPPYNSEPAEESSEHKNNNYEPNLFKTPQRKPS
4 | ...
5 | QFISVSESTRTAPTSSSEDYLRLLKRRCTTSLIKEQESSQASTEECEKNKQDTITTKKYI
6 | $
```

4. 위의 `sequence.protein.fasta` 파일의 내용을 `sequence.protein.2.fasta` 파일로 출력한다.

```
1 | $ python bioinformatics_3_4.py
2 | $ cat sequence.protein.2.fasta
3 | >AAB07223.1 BRCA2 [Homo sapiens]
4 | MPIGSKERPTFFEIFKTRCNKADLGPISLNWFEELSSEAPPYNSEPAEESSEHKNNNYEPNLFKTPQRKPS
5 | ...
6 | QFISVSESTRTAPTSSSEDYLRLLKRRCTTSLIKEQESSQASTEECEKNKQDTITTKKYI
7 | $
```

5. 위에서 만든 `sequence.protein.2.fasta` 의 둘째줄을 다음과 같이 출력한다.

```
1 | $ python bioinformatics_3_5.py
2 | The second line is: MPIGSKERPTFFEIFKTRCNKADLGPISLNWFEELSSEAPPYNSEPAEESSEHKNN
3 | $
```

6. `sequence.protein.2.fasta` 의 첫째줄은 `title` 이란 변수에, 둘째줄부터 끝까지는 개행문자를 없앤

후 `seq` 이란 변수에 저장하고 다음과 같이 출력한다.

```
1 $ python bioinformatics_3_6.py
2 title: >AAB07223.1 BRCA2 [Homo sapiens]
3 seq: MPIGSKERPT...QDTITTKKYI
4 $
```

7. 위에서 만든 `seq` 문자열의 위치를 물어보고, 입력된 위치에서 시작하는 3개의 아미노산을 다음과 같이 출력한다.
`xxx` 를 입력할 때 까지 이 작업을 반복한다.

```
1 $ python bioinformatics_3_7.py
2 Position: 15
3 Three amino acids: FKT
4 Position: 20
5 Three amino acids: NKA
6 Position: XXX
7 Okay, I will stop.
8 $
```

8. 찾고자 하는 아미노산을 물어보고, 입력된 아미노산을 가진 `seq` 문자열의 모든 위치를 다음과 같이 출력한다.
`xxx` 를 입력할 때 까지 이 작업을 반복한다.

```
1 $ python bioinformatics_3_8.py
2 Searching for: S
3 Found at: 5, 28, 36, 37, 44, ..., 3376, 3389, 3396, 3397, 3400
4 Searching for: Q
5 Found at: 66, 73, 84, 92, 126, ..., 3350, 3361, 3394, 3398, 3409
6 Searching for: XXX
7 Okay, I will stop.
8 $
```

4일차: GenBank 포맷 파싱, 문자열 포매팅, 서열 핸들(상보적서열, 역서열, 등 생성), 폴드 핸들

1. <https://www.ncbi.nlm.nih.gov/protein/AAB07223.1?report=genpept> 에서 protein GenBank 파일을 받아 `sequence.protein.gb` 란 이름으로 저장한 후 첫째줄은 `title` 이란 변수에, 'ORIGIN'이라는 단어로 시작되는 행을 찾아, 그 다음 행부터 끝까지 읽어 `seq` 이라는 변수에 저장하고 다음과 같이 출력한다.

```

1 | $ python bioinformatics_4_1.py
2 | title: LOCUS          AAB07223          3418 aa          linear  PRI
3 | seq:          1 mpigskerpt ffeifktrcn kadlgpisln wfeelsseap pynsepaees ehknn
4 | ...
5 |          3361 qfisvsestr taptssedyl rlkrrcttsl ikeqessqas teeceknkqd tittkkyi
6 | //
7 | $

```

2. 위에서 만든 `seq` 문자열에서 알파벳만 저장하여 다음과 같이 출력한다.

```

1 | $ python bioinformatics_4_2.py
2 | title: LOCUS          AAB07223          3418 aa          linear  PRI
3 | seq: mpigskerptffeifktrcnkadlgpisln...ikeqessqasteceknkqdtittkkyi
4 | $

```

3. 위에서 만든 `seq` 문자열을 한 줄에 70개의 문자를 포맷팅하여 다음과 같이 출력한다.

```

1 | $ python bioinformatics_4_3.py
2 | mpigskerptffeifktrcnkadlgpislnwfeelsseappynsepaeesehknnnyepnlfktpqrkps
3 | ...
4 | qfisvsestrtaptssedylrlkrrcttslikeqessqasteceknkqdtittkkyi
5 | $

```

4. https://www.ncbi.nlm.nih.gov/nuccore/NM_001276699.1?report=fasta 에서 nucleotide FASTA 파일을 받아 `sequence.nucleotide.fasta` 란 이름으로 저장하고, https://www.ncbi.nlm.nih.gov/nuccore/NM_001276699.1?report=genbank 에서 nucleotide GenBank 파일을 받아 `sequence.nucleotide.gb` 란 이름으로 저장한 후 입력파일명, 출력파일명, 옵션을 입력받아 아래와 같이 동작하는 프로그램을 작성하시오. (option-3의 경우 FASTA format의 header line은 GeneBank format의 첫 번째 줄로 한다.)

```

1 $ python bioinformatics_4_4.py
2 Enter input file: sequence.nucleotide.fasta
3 Enter output file: reverse.sequence.nucleotide.fasta
4 Option-1) Read a FASTA format DNA sequence file and make a reverse sequence
5 Option-2) Read a FASTA format DNA sequence file and make a reverse compleme
6 Option-3) Convert GenBank format file to FASTA format file.
7 Select the option: 1
8 $ cat reverse.sequence.nucleotide.fasta
9 >NM_001276699.1 Homo sapiens tumor protein p53 (TP53), transcript variant 7
10 GTGGGGAGTCTGTGTGTCCACCGTCGTTTCAAATAACATTTTATTCTCTAGCTATATTTTACCCTATA
11 ...
12 TCGGAGGTAGAGGACCGGAGT
13 $
14 $ python bioinformatics_4_4.py
15 Enter input file: sequence.nucleotide.fasta
16 Enter output file: reverse.complement.sequence.nucleotide.fasta
17 Option-1) Read a FASTA format DNA sequence file and make a reverse sequence
18 Option-2) Read a FASTA format DNA sequence file and make a reverse compleme
19 Option-3) Convert GenBank format file to FASTA format file.
20 Select the option: 2
21 $ cat reverse.sequence.nucleotide.fasta
22 >NM_001276699.1 Homo sapiens tumor protein p53 (TP53), transcript variant 7
23 CACCCCTCAGACACAGGTGGCAGCAAAGTTTTATTGTAAAATAAGAGATCGATATAAAAATGGGATAT
24 ...
25 AGCCTCCATCTCCTGGCCTCA
26 $
27 $ python bioinformatics_4_4.py
28 Enter input file: sequence.nucleotide.gb
29 Enter output file: new.sequence.nucleotide.fasta
30 Option-1) Read a FASTA format DNA sequence file and make a reverse sequence
31 Option-2) Read a FASTA format DNA sequence file and make a reverse compleme
32 Option-3) Convert GenBank format file to FASTA format file.
33 Select the option: 3
34 $ cat new.sequence.nucleotide.fasta
35 >LOCUS          NM_001276699          2331 bp    mRNA    linear    PRI 28-DEC
36 tgaggccaggagatggaggctgcagtgagctgtgatcacaccactgtgctccagcctgagtgacagagca
37 ...
38 cacctgtgtgtctgaggggtg
39 $

```

5. 위의 `sequence.nucleotide.gb` 파일의 내용 중 'Title'에 해당하는 문자열을 읽어 들인 후, 읽은 순서대로 각 내용당 한 줄로 만들어 다음과 같이 출력한다.


```

1 | $ python bioinformatics_4_5.py
2 |     TITLE      G-quadruplex structures in TP53 intron 3: role in alternative s
3 |                 Delta160p53 is a novel N-terminal p53 isoform encoded by Delta1
4 |                 ...
5 |                 Isolation and characterization of a human p53 cDNA clone: expre
6 | $

```

5일차: 6-frame translation

1. 위의 `sequence.nucleotide.fasta` 파일의 `sequence`를 한 줄에 3-base codon 단위로 다음과 같이 출력한다. (끝에 codon 이 만들어지지 않는 부분은 출력하지 않는다.)

```

1 | $ python bioinformatics_5_1.py
2 | TGA
3 | GGC
4 | CAG
5 | ...
6 | TGA
7 | GGG
8 | GTG
9 | $

```

2. 위의 `sequence.nucleotide.fasta` 파일의 `sequence`를 한 줄에 60 bases 단위로 출력하고 각 줄당 해당하는 아미노산으로 translation 하여 다음과 같이 출력한다.

```

1 | $ python bioinformatics_5_2.py
2 | TGAGGCCAGGAGATGGAGGCTGCAGTGAGCTGTGATCACACCACTGTGCTCCAGCCTGAG
3 | *  G  Q  E  M  E  A  A  V  S  C  D  H  T  T  V  L  Q  P  E
4 | TGACAGAGCAAGACCCTATCTCAAAAAAAAAAAAAAAAAAGAAAAGCTCCTGAGGTGTAG
5 | *  Q  S  K  T  L  S  Q  K  K  K  K  K  R  K  A  P  E  V  *
6 | ...
7 | ATCTCTTATTTTACAATAAACTTTGCTGCCACCTGTGTGTCTGAGGGGTG
8 | I  S  Y  F  T  I  K  L  C  C  H  L  C  V  *  G  V
9 | $

```

3. 위의 `sequence.nucleotide.fasta` 파일의 `sequence`를 이용하여 translation을 6개의 reading frame으로 다음과 같이 출력한다. (출력순서 : forward 1, forward 2, forward 3, sequence, reverse complement sequence, reverse 1, reverse 2, reverse 3)

```

1 | $ python bioinformatics_5_3.py
2 | * G Q E M E A A V S ...L C C H L C V * G V
3 | E A R R W R L Q * A ...F A A T C V S E G
4 | R P G D G G C S E L... L L P P V C L R G
5 | TGAGGCCAGGAGATGGAGGCTGCAGTGAGC...CTTTGCTGCCACCTGTGTGTCTGAGGGGTG
6 | ACTCCGGTCCTCTACCTCCGACGTCACG...GAAACGACGGTGGACACACAGACTCCCCAC
7 | S A L L H L S C H A... K A A V Q T D S P H
8 | P W s I s A A T L ... S Q Q W R H T Q P T
9 | L G P S P P Q L S ...V K S G G T H R L P
10 | $

```

4. 위의 `sequence.nucleotide.gb` 파일을 읽어 CDS 영역의 정보를 찾아 `sequence`에서 해당 서열을 translation 한 결과를 다음과 같이 출력한다.

```

1 | $ python bioinformatics_5_4.py
2 | MAIYKQSQHMTVEVVRCPHHERCSDSDGLAPPQHILIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDC
3 | TTIHYNMCMNSSCMGGMNRRPILTIITLEDSSGNLLGRNSFEVRVCACPRDRRTEENLRKKGEPHHEL
4 | PPGSTKRALPNNTSSSPQPKKKPLDGEYFTLQMLDLRWCYFLINSS*
5 | $

```