

Übungen zum Propädeutikum Programmierung in der Bioinformatik

Blatt 8

Termin: Dienstag, 10. Dezember 2019

Achte bei den folgenden Aufgaben darauf, die Exceptions mit `throws` weiterzugeben und/oder mit einem `try-catch/try-with-resources` Konstrukt zu behandeln, da deine Programme sonst eventuell nicht kompilieren.

Übung 1 *Fasta einlesen mit `BufferedReader`*

Schreibe eine Klasse `Proteome` die mehrere Proteinsequenzen direkt aus einer Fasta-Datei ausliest und speichert.

Die Klasse soll folgende Elemente enthalten:

1. eine Variable `ArrayList<Protein> proteins`, in der alle Proteine aus dem Fasta abgelegt werden. Du kannst natürlich deine in den vorherigen Übungen geschriebene `Protein`-Klasse verwenden. Falls du diese nicht hast, schreibe eine kleine Klasse die nur eine Protein-ID und die Sequenz enthält.
2. einen Konstruktor `Proteome(String fastaFile)`. In diesem Konstruktor soll mit einem `BufferedReader` die Fasta-Datei eingelesen werden, und für jede Sequenz darin ein `Protein` zu `proteins` hinzugefügt werden. Dabei soll die ID des Proteins auf den Header gesetzt werden (*ohne* das `>`-Symbol am Zeilenbeginn). Zum Testen steht auf der Website `Vibrio_cholerae.fasta` zum Download bereit.
3. eine Methode `int getProteinCount()` (gibt die Anzahl der geladenen Proteine zurück)
4. eine Methode `Protein getProtein(int i)` (gibt das *i*-te Protein-Objekt zurück)

Übung 2 *Statistiken schreiben mit `BufferedWriter`*

Erweitere `Proteome` jetzt um eine Methode `void writeStatistics(String outputFile)`. Diese soll für das gesamte Proteom folgende Statistiken in eine Datei `outputFile` schreiben:

1. minimale Sequenzlänge
2. maximale Sequenzlänge
3. Mittelwert aller Sequenzlängen

Achte darauf, dass du beim Schreiben der Datei keine anderen Dateien überschreibst.

Übung 3 *Einlesen einer TSV-Datei*

Erstelle nun in `Proteome` eine Variable `HashMap<Character,Double> massMap`. Diese Variable soll von einer Methode `void readAminoAcidMassFile(String file)` befüllt werden, welche eine TSV-Datei als Input erhält die die Massen der einzelnen Aminosäuren enthält. Eine TSV-Datei ist letztendlich nur eine Tabelle, in der die Spalten durch einen Tab getrennt sind (TSV = "Tab-separated-values").

`mass.tsv` enthält in der ersten Spalte die Buchstaben für jede Aminosäure, in der zweiten die zur jeweiligen Aminosäure gehörige Masse. Folgende Informationen könnten hilfreich für das **Parse**n dieser Datei sein:

- Ein Tab wird in Java (und vielen anderen Kontexten) durch die Zeichenfolge `\t` beschrieben.
- Mit der Methode `String[] split(String pattern)` kann man einen String an dem angegebenen `pattern` in Einzelstrings aufspalten.
- Die Methode `String trim()` returned eine Kopie des Strings mit allem eventuell vorhandenen **Whitespace** um den String herum entfernt.

Nun könntest du mit Hilfe von `massMap` bequem eine Methode schreiben die die Masse einer Proteinsequenz aufsummiert.

Lösung 3

```
1 import java.io.*;
2 import java.util.ArrayList;
3 import java.util.HashMap;
4
5 public class Proteome {
6     // Aufgabe 1
7     private ArrayList<Protein> proteins;
8     private HashMap<Character,Double> massMap; // Aufgabe 3
9
10    public Proteome(String fastaFile) {
11        this.proteins = new ArrayList<>();
12        this.massMap = new HashMap<>(); // Aufgabe 3
13        try (BufferedReader br = new BufferedReader(new FileReader(fastaFile))) {
14            String line = "";
15            String id = "";
16            StringBuilder seq = new StringBuilder();
17            while ((line = br.readLine()) != null) {
18                if (line.startsWith(">")) {
19                    if (!id.isEmpty()) {
20                        proteins.add(new Protein(id, seq.toString()));
21                        seq = new StringBuilder();
22                    }
23                    id = line.substring(1);
24                } else {
25                    seq.append(line.strip());
26                }
27            }
28            proteins.add(new Protein(id, seq.toString()));
29        } catch (FileNotFoundException e) {
30            e.printStackTrace();
31        } catch (IOException e) {
32            e.printStackTrace();
33        }
34    }
35
36    public int getProteinCount() {
37        return proteins.size();
38    }
39
40    public Protein getProtein(int i) {
41        return proteins.get(i);
42    }
43
44    // Aufgabe 2
45    public void writeStatistics() {
46        int min = Integer.MAX_VALUE;
47        int max = Integer.MIN_VALUE;
48        int total = 0;
49        for (Protein p : proteins) {
50            int length = p.getSequence().length();
51            total += length;
52            if (length > max) max = length;
53            if (length < min) min = length;
54        }
55        try (BufferedWriter bw = new BufferedWriter(new FileWriter("outputFile"))) {
56            bw.write(min + " " + max + " " + (total*1.0)/this.getProteinCount());
57        } catch (IOException e) {
```

```

58         e.printStackTrace();
59     }
60 }
61
62 // Aufgabe 3
63 public void readAminoAcidMassFile(String file) {
64     try (BufferedReader br = new BufferedReader(new FileReader(file))) {
65         String line = "";
66         while ((line = br.readLine()) != null) {
67             String[] stringParts = line.trim().split("\t");
68             massMap.put(stringParts[0].charAt(0), Double.parseDouble(stringParts[1]));
69         }
70     } catch (FileNotFoundException e) {
71         e.printStackTrace();
72     } catch (IOException e) {
73         e.printStackTrace();
74     }
75 }
76
77 public static void main(String[] args) {
78     Proteome p = new Proteome("Vibrio_cholerae.fasta");
79     System.out.println(p.getProteinCount()); // 3504
80     p.writeStatistics(); // 26 4558 325.80422374429224
81     p.readAminoAcidMassFile("mass.tsv");
82 }
83 }
84 o

```