

1. Diseñar un pequeño programa que lea el fichero de datos y calcule la valoración media, el número de votos medios considerando la colección All-Capítulos) así como el número de personajes, en media, que tiene los datos en la colección Guiones. Se recomienda utilizar openCSV (o cualquier otra librería para el acceso correcto a los datos). Además los datos deben pasar por un proceso de limpieza y mecanismo para tratar las excepciones.

Definimos dos rutas de directorio carpetaCU y carpetaC que representan las carpetas que contienen los archivos CSV de dos colecciones de datos: "CapitulosUnidos" y "Capitulos."

Inicializamos las variables valoraciones, votos, personajesTotal, totalCU, y totalC para mantener un seguimiento de las valoraciones, votos, número total de personajes y el número total de archivos procesados.

Creamos objetos File para representar las carpetas y luego obtiene la lista de archivos CSV en cada carpeta que cumplan con la extensión ".csv."

Recorremos la lista de archivos CSV en la carpeta "CapitulosUnidos" (representada por csvUnidos), y para cada archivo, realizamos lo siguiente:

- Abre el archivo CSV y lee sus contenidos.
- Extrae la valoración (valor numérico) y el número de votos de cada registro en el archivo CSV.
- Acumula las valoraciones y votos en las variables valoraciones y votos.
- Incrementa el contador totalCU para realizar un seguimiento del número de archivos procesados.

Recorremos la lista de archivos CSV en la carpeta "Capitulos" (representada por csvCaps), y para cada archivo, realizamos lo siguiente:

- Abre el archivo CSV y lee sus contenidos.
- Extrae el número de personajes de cada registro en el archivo CSV.
- Acumula el número de personajes en la variable personajesTotal.
- Incrementa el contador totalC para realizar un seguimiento del número de archivos procesados.

Calculamos las medias de valoración y votos si se procesaron archivos en la carpeta "CapitulosUnidos." Luego, muestra estos valores en la consola.

Calculamos la media del número de personajes por capítulo si se procesaron archivos en la carpeta "Capitulos." Luego, muestra este valor en la consola.

El programa utiliza la librería openCSV para leer y analizar los archivos CSV y realiza un manejo básico de excepciones para manejar errores en la lectura de archivos

```

9  import java.util.Iterator;
10
11
12  public class ej1{
13
14      Run | Debug
15      public static void main(String[] args){
16          String carpetaCU = "/home/anne/Escritorio/RI/RI/P3/CapitulosUnidos";
17          String carpetaC = "/home/anne/Escritorio/RI/RI/P3/Capitulos";
18
19          double valoraciones=0;
20          double votos=0;
21          int personajesTotal=0, totalCU=0, totalC=0;
22
23          File carpetaCapUnid = new File(carpetaCU);
24          File carpetaCaps = new File(carpetaC);
25
26          File[] csvUnidos= carpetaCapUnid.listFiles((dir,nombre)->nombre.endsWith(".csv"));
27          File[] csvCaps= carpetaCaps.listFiles((dir,nombre)->nombre.endsWith(".csv"));
28
29          if(csvUnidos!=null){
30              for(int i=0; i< csvUnidos.length;i++){
31                  File archivo= csvUnidos[i];
32                  Reader reader = new FileReader(archivo);
33                  if(!reader.ready()) {
34                      System.out.println("Error en el fichero de datos");
35                  }
36
37                  CSVReader csvReader= new CSVReader(reader);
38                  String[] nextRecord;
39                  String[] firstLine;
40
41                  firstLine=csvReader.readNext();
42
43                  while ((nextRecord = csvReader.readNext()) != null){
44                      if(nextRecord.length>=4){
45                          double valoracion = Double.parseDouble(nextRecord[3]);
46                          int votacion= Integer.parseInt(nextRecord[4]);
47
48                          valoraciones+=valoracion;
49                          votos+=votacion;
50                      }
51                  }
52                  totalCU++;
53              }
54          }
55          if(csvCaps!=null){

```

```

56              for(int i=0; i< csvCaps.length;i++){
57                  File archivo= csvCaps[i];
58                  Reader reader = new FileReader(archivo);
59                  if(!reader.ready()) {
60                      System.out.println("Error en el fichero de datos");
61                  }
62
63                  CSVReader csvReader= new CSVReader(reader);
64                  String[] nextRecord;
65                  String[] firstLine;
66
67                  firstLine=csvReader.readNext();
68
69                  while ((nextRecord = csvReader.readNext()) != null){
70                      if(nextRecord.length>=4){
71                          int personajes=Integer.parseInt(nextRecord[3]);
72                          personajesTotal+=personajes;
73                      }
74                  }
75                  totalC++;
76              }
77          }
78      }
79
80      if(totalCU>0){
81          double valoracionMedia= valoraciones/totalCU;
82          double mediaVotos= votos/totalCU;
83
84          System.out.println("Valoración media de todos los capitulos --> "+ valoracionMedia);
85          System.out.println("Media del número de votaciones --> "+ mediaVotos);
86      }
87
88      if(totalC>0){
89          double mediaPersonajes = personajesTotal/totalC;
90          System.out.println("Media del número de personajes por capítulo --> "+ mediaPersonajes);
91      }
92  }
93
94
95
96
97
98
99
100 }

```

2. Identificar para cada campo de los distintos ficheros qué analizador utilizarías, y en caso de necesitar uno específico implementarlo. Recordad que el tipo de analizador dependerá de cómo vayáis a realizar las búsquedas por lo que es necesario plantearos antes esta cuestión.

En la carpeta capitulos podemos diferenciar los siguientes campos en los CSV

- episode_id : keyWordAnalyzer
- number : keyWordAnalyzer
- timestamp_in_ms : keyWordAnalyzer
- raw_character_text : CustomAnalyzer
- raw_location_text : simpleAnalyzer
- spoken_words : standarAnalyzer

En la carpeta capitulos unidos diferenciamos los siguientes campos

- episode_id : keyWordAnalyzer
- spoken_words : stopAnalyzer
- raw_character_text : CustomAnalyzer
- imdb_rating : keyWordAnalyzer
- imdb_votes : keyWordAnalyzer
- numer_in_season : keyWordAnalyzer
- original_air_date : keyWordAnalyzer
- season : keyWordAnalyzer
- title : stopAnalyzer
- us_viewers_in_millions : keyWordAnalyzer
- views : keyWordAnalyzer

CustomAnalyzer

El analizador creado es un analizador que pone las palabras a minúscula y posteriormente reemplaza los espacios en blanco por _ para facilitar el posterior procesamiento de los datos

```
Analyzer analyzer = CustomAnalyzer.builder()
    .withTokenizer("standard")
    .addTokenFilter("lowercase")
    .addTokenFilter("patternreplace",
        "pattern" , "\\s+", // Eliminar espacios vacíos
        "replacement" , "_" // Reemplazar por _ para posterior procesamiento de datos
    )
    .build();
```

3. Diseñar una consulta, para ello debemos de proporcionar los términos de la misma (con cuatro o mas términos) así como qué es lo que debe tener un documento para ser considerado relevante. Para encontrar los documentos (diálogos/capítulos) relevantes en la colección se debe lanzar la consulta sobre dos índices (uno para cada colección).

La consulta que he seleccionado ha sido Homer Simpson Duff Beer

Los resultados han sido los siguientes

```
home > phuertass > Documentos > RI_PABLO > RI > P3 > ≡ queryC.qrels
1  Query: ASPH
2  Homer Simpson Duff Beer
3  Description:
4  Busca diálogos en los que Homer Simpson hable o esté relacionado con Duff Beer.
5  Index:
6  C
7  Relevant documents:
8  5
9  cap 75 l 1
10 cap 272 l 1
11 cap 263 l 1
12 cap 72 l 1
13 cap 179 l 1
14
```

```
home > phuertass > Documentos > RI_PABLO > RI > P3 > ≡ queryD.qrels
1  Query: ASPH
2  Homer Simpson Duff Beer
3  Description:
4  Busca diálogos en los que Homer Simpson hable o esté relacionado con Duff Beer.
5  Index:
6  D
7  Relevant documents:
8  4
9  cap 72 l 207
10 cap 409 l 169
11 cap 63 l 30
12 cap 75 l 47
```