

Sistemas de Gestión de Datos y de la Información
Máster en Ingeniería Informática, 2020-21
Práctica Recuperación Información

Fecha de entrega: domingo 13 de diciembre de 2020, 23:55h

Entrega de la práctica

La práctica se entregará en un único fichero **GrupoXX.zip** mediante el Campus Virtual de la asignatura. El fichero contendrá una carpeta por cada uno de los apartados.

Lenguaje de programación

Python 3.6.

Ficheros

Colección de mensajes 20news-18828.tar.gz de un grupo de noticias, obtenido de <https://archive.ics.uci.edu/ml/datasets/Twenty+Newsgroups>. Listado charset.txt con la codificación de cada fichero de la colección.

Calificación

Además del correcto funcionamiento del código se valorará también su claridad, su documentación y su organización en funciones auxiliares reutilizables.

Declaración de autoría e integridad

Todos los ficheros entregados contendrán una cabecera en la que se indique la asignatura, la práctica, el grupo y los autores. Esta cabecera también contendrá la siguiente declaración de integridad:

Yo, *Nombre Alumno*, declaro que esta solución es fruto exclusivamente de mi trabajo personal. No he sido ayudado por ninguna otra persona ni he obtenido la solución de fuentes externas, y tampoco he compartido mi solución con nadie. Declaro además que no he realizado de manera deshonesto ninguna otra actividad que pueda mejorar mis resultados ni perjudicar los resultados de los demás.

No se corregirá ningún fichero que no venga acompañado de dicha cabecera.

1. Recuperación vectorial [5pt]

Completar el esqueleto `VectorialIndex.py` para implementar la clase `VectorialIndex` que procesa una colección de documentos, crea un índice invertido con pesos TF-IDF y lo usa para acelerar consultas vectoriales y booleanas (únicamente conjunción) sobre la colección. La clase debe contener al menos los siguientes métodos:

- `__init__(self, directorio)`
Recorre **recursivamente todos los ficheros** que hay en `directorio` y crea un índice invertido para relacionar las palabras con una lista de parejas (documento, peso). Los pesos se deben calcular usando la técnica TF-IDF, y los documentos dentro del índice deben representarse como números enteros en lugar de rutas completas para minimizar el espacio que ocupan en memoria. Para que todos consideremos la misma *definición* de palabra, se debe usar la función `extrae_palabras` que aparece en el esqueleto para extraer las palabras de una línea de texto. Es importante que el índice invertido almacene la **mínima información necesaria** para optimizar su tamaño.
- `consulta_vectorial(self, consulta, n=3)`.
Dada una consulta representada como una cadena de palabras no repetidas separadas por espacios, devuelve una lista de parejas (`ruta.fichero`, `relevancia`) con los `n` resultados más relevantes usando el modelo de recuperación vectorial ordenadas de mayor a menor relevancia.
- `consulta_conjuncion(self, consulta)`
Dada una consulta representada como una cadena de palabras no repetidas separadas por espacios que se entiende como una conjunción, devuelve una lista de nombres de fichero con todos los resultados en los que aparecen **todas las palabras** de la consulta.

En los métodos `consulta_vectorial()` y `consulta_conjuncion()` el índice invertido se debe utilizar de la manera más eficiente posible para resolver la consulta.

La clase `VectorialIndex` se utilizaría de la siguiente manera:

```
>>> i = VectorialIndex('20news-18828')
>>> i.consulta_vectorial('DES Diffie-Hellman', n=5)
[('20news-18828/sci.crypt/15991', 0.5096209025568682),
 ('20news-18828/sci.crypt/15826', 0.40780449792065276),
 ('20news-18828/sci.crypt/16141', 0.3327934131638604),
 ('20news-18828/sci.crypt/15894', 0.31276008256792764),
 ('20news-18828/sci.crypt/15592', 0.29572512521719024)]
>>> i.consulta_conjuncion('DES Diffie-Hellman')
['20news-18828/sci.crypt/15670',
 '20news-18828/sci.crypt/14831',
 '20news-18828/sci.crypt/15746',
 '20news-18828/sci.crypt/15493',
 '20news-18828/sci.crypt/15991',
 '20news-18828/sci.crypt/15644',
 '20news-18828/sci.crypt/15301',
 '20news-18828/sci.crypt/15241']
```

2. Índice invertido completo [5pt]

Completar el esqueleto `CompleteIndex.py` para implementar la clase `CompleteIndex` que procesa una colección de documentos, crea un índice invertido completo y lo usa para acelerar consultas de frase sobre la colección. La clase tiene un interfaz similar al índice vectorial y debe contener al menos los siguientes métodos:

- `__init__(self, directorio)`
Recorre **recursivamente todos los ficheros** que hay en `directorio` y crea un índice invertido completo para relacionar cada palabra con una lista de tuplas (`numero_doc`, `l_pos`), donde `numero_doc` es el número de documento y `l_pos` contiene una lista de posiciones. Como granularidad usaremos el número de palabra dentro del documento empezando en 1. Al igual que en el apartado anterior, se debe usar la función `extrae_palabras` que aparece en el esqueleto para extraer las palabras de una línea de texto.
- `consulta_frase(self, frase)`.
Dada una consulta representada como una cadena de palabras separadas por espacios, devuelve una lista con los nombres de los ficheros en los que aparece **exactamente** esa frase.

La clase `CompleteIndex` se utilizaría de la siguiente manera:

```
>>> i = CompleteIndex('20news-18828')
>>> i.consulta_frase('either terrestrial or alien')
['20news-18828/alt.atheism/53209',
 '20news-18828/alt.atheism/53222',
 '20news-18828/alt.atheism/53218']
>>> i.consulta_frase('is more complicated')
['20news-18828/comp.os.ms-windows.misc/9966',
 '20news-18828/comp.os.ms-windows.misc/10010',
 '20news-18828/alt.atheism/53564',
 '20news-18828/soc.religion.christian/20725']
```

Importante

- Es importante no fallar al abrir los ficheros debido a su codificación. Revisad la documentación de `open` en <https://docs.python.org/3/library/functions.html#open> para descubrir qué parámetros os pueden ayudar. También tenéis un listado `charset.txt` con la codificación detectada por el programa `file` en cada fichero de la colección.
- No tratéis de procesar la colección completa `20news-18828` desde el primer momento, ya que puede tardar mucho. Cread una colección mínima para hacer pruebas y probad la colección completa únicamente cuando estéis seguros de vuestro código.